

Metodología Bayesiana para la Toma de Decisiones

**Aplicaciones en el Control Ambiental
de Actividades Portuarias**

**José-Miguel Bernardo
Universitat de València, España**

Informe Metodológico

Noviembre 2004

Documento producido para el Ente Público *Puertos del Estado* en virtud del convenio de colaboración suscrito en Junio de 2004 entre *Puertos del Estado* y la *Universitat de València*.

Valencia, Noviembre de 2004.

Índice

1. Introducción	
1.1. La toma de decisiones en el proyecto HADA.....	5
1.2. Contenido de este documento	7
2. Fundamentos Metodológicos	
2.1. Teoría de la decisión	9
2.2. Discrepancia e información	17
2.3. Métodos estadísticos bayesianos	23
2.4. Bibliografía comentada	33
3. Algunos Problemas Medioambientales Portuarios	
3.1. Red de control ambiental en La Coruña	39
3.2. Descarga de soja en Barcelona	42
3.3. Posible impacto ambiental en Huelva.....	44
3.4. Descarga de graneles sólidos en Santander	46
3.5. Descarga de carbón en Tarragona	49
Apéndice	
Bayesian statistics.....	52

Capítulo 1.

Introducción

1.1. LA TOMA DE DECISIONES EN EL PROYECTO HADA

La notable importancia que los ciudadanos actualmente conceden a los problemas medioambientales aparece reflejada en las directivas de calidad medioambiental, crecientemente exigentes, que son emitidas por la Unión Europea y transcritas por las Administraciones de los Estados miembros a sus respectivas legislaciones. Por otra parte, más allá de las exigencias legales, los actores económicos son ya conscientes de que una parte importante de su imagen depende del interés que demuestren en minimizar los posibles impactos negativos de sus actividades sobre el medio ambiente.

En este contexto, el Ente Público *Puertos del Estado* tomó la iniciativa de potenciar, en el contexto de los proyectos europeos LIFE-Medio Ambiente, el diseño de una *Herramienta Automática de Diagnóstico Ambiental* (HADA) que deberá permitir una gestión automatizada, en tiempo real, de los problemas medioambientales que puedan derivarse de las operaciones que se desarrollan en los puertos marítimos españoles.

El ámbito de intervención especificado en el proyecto HADA es muy amplio, incluyendo, entre otros aspectos, el desarrollo y ordenación del uso del territorio, la gestión del agua, el control de los impactos de las actividades económicas y la gestión de residuos, todo ello coordinado en el marco una política de producto integrado.

La gestión medioambiental de un sistema tan complejo como el de un puerto marítimo exige una consideración integrada de los problemas enfrentados, de sus probables interrelaciones, de los factores de los que dependen los resultados de sus

posibles soluciones, de la información disponible y de la valoración política de sus posibles consecuencias. Esto es posible en el marco de la *teoría de la decisión*, una construcción normativa que prescribe la *única* forma de tomar decisiones en ambiente de incertidumbre compatible con un comportamiento racional.

Naturalmente, las decisiones deben ser tomadas haciendo uso de toda la información relevante que resulte accesible, de forma que resulta crucial disponer de un mecanismo que permita la actualización inmediata de la información disponible en función de los datos que se van obteniendo. Los *métodos estadísticos bayesianos* constituyen la *única* forma compatible con los principios de comportamiento racional de incorporar información experimental adicional a la información inicialmente disponible.

La teoría bayesiana de la decisión proporciona una herramienta global, una visión de conjunto, con la que abordar *cualquier* problema de elección en ambiente de incertidumbre. Los métodos estadísticos bayesianos proporcionan la forma más general y más potente de inferencia estadística; permite utilizar la información inicial proporcionada por expertos—si se dispone de ella—y utiliza la teoría de la información para determinar distribuciones “noinformativas” que producen resultados *objetivos*—resultados que sólo dependen del modelo probabilístico asumido y de los datos experimentales obtenidos. Lamentablemente, sin embargo, tanto la teoría de la decisión como la metodología bayesiana de inferencia estadística son todavía relativamente poco conocidas fuera del ámbito académico.

En este contexto, *Puertos del Estado* y la *Universitat de València* firmaron en Junio de 2004 un convenio de colaboración para:

- (i) elaborar un informe metodológico divulgativo sobre el sistema bayesiano de toma de decisiones, con especial atención a los problemas de decisión que se plantean en el contexto del control ambiental de actividades portuarias, y
- (ii) diseñar para el Puerto de Tarragona—que aceptó actuar como puerto piloto—una estrategia de control ambiental, basada en la teoría de la decisión y en la metodología estadística bayesiana, que permita una gestión automatizada del control ambiental de sus actividades.

Este documento contiene el informe metodológico a que se refiere el primero de estos apartados. Proporciona una introducción elemental a la teoría de la decisión y a los métodos bayesianos de inferencia estadística, y describe la forma en la que la teoría de la decisión y la metodología estadística bayesiana pueden ser utilizadas para resolver algunos de los problemas de decisión en ambiente de incertidumbre a los que se enfrentan las Autoridades Portuarias.

1.2. CONTENIDO DE ESTE DOCUMENTO

Este documento contiene tres capítulos, de los que el primero es esta introducción, seguidos de un largo apéndice.

El capítulo segundo es una introducción elemental a los fundamentos metodológicos de la teoría matemática de la decisión y de la inferencia estadística bayesiana. En particular, en la Sección 2.1 se especifica la estructura formal de los problemas de decisión en ambiente de incertidumbre, y se describen brevemente los resultados más importantes de la teoría de la decisión. En la Sección 2.2 se introducen los conceptos básicos de divergencia intrínseca, de asociación intrínseca y de cantidad de información que puede esperarse de los resultados experimentales. En la Sección 2.3 se analizan las características básicas de los métodos estadísticos bayesianos. Finalmente, en la Sección 2.4 se ofrece una bibliografía comentada para quienes deseen profundizar en los fundamentos o en las aplicaciones de los métodos descritos.

En el capítulo tercero se describen, a título ilustrativo, algunos de los problemas de decisión en el contexto de gestión medioambiental planteados en los puertos integrados en el proyecto HADA poniendo de manifiesto la forma en la que la teoría de la decisión y la metodología estadística bayesiana podrían ser utilizadas para hacerles frente, y especificando el tipo de información que debería ser recogida y estudiada para su análisis racional. Se ha prestado especial atención a los problemas asociados al riesgo de contaminación atmosférica que puede derivarse de la manipulación de graneles sólidos pulverulentos. En particular, se consideran problemas de decisión planteados por la descarga de haba de soja en el puerto de Barcelona, y por el almacenamiento de carbón al aire libre en los puertos de Santander y de Tarragona; se analiza el problema de determinar la posible localización de estaciones de medida de partículas en el puerto de La Coruña, y se estudia el posible impacto ambiental del puerto de Huelva sobre el paraje protegido de las marismas del Odiel.

La memoria concluye con un largo apéndice. Se trata de una versión, abreviada pero actualizada, del artículo *Bayesian Statistics*, preparado por el autor de este informe para el volumen de Estadística de *ELOSS*, una enciclopedia científica que la UNESCO publicó en 2003. Este apéndice, redactado en inglés, constituye una extensión del material descrito en el Capítulo 2 que permite—a los lectores que así lo requieran—una introducción más detallada a los métodos estadísticos bayesianos desde la perspectiva unificadora de la teoría de la decisión.

Capítulo 2.

Fundamentos Metodológicos

2.1. TEORÍA DE LA DECISIÓN

Existe un problema de decisión cuando debe elegirse entre dos o más formas de actuar. Aunque sociólogos, historiadores y políticos han escrito a menudo sobre la forma en que determinadas decisiones se toman o han sido tomadas, se ha escrito relativamente poco sobre la forma en que *deberían* tomarse. La teoría de la decisión propone un determinado método de tomar decisiones y demuestra además que es el único método de decisión compatible con unos pocos principios básicos sobre la *elección coherente* entre opciones alternativas en ambiente de incertidumbre.

Conjunto de alternativas. Lo primero que hay que hacer ante un problema de decisión es considerar *todas* las formas de actuación posibles. No es necesario distinguir entre una decisión y la acción a que da lugar. En efecto, si la acción no llega a realizarse es porque algo lo ha impedido, dando lugar con ello a un nuevo problema de decisión. Generalmente, no resulta adecuado considerar únicamente una decisión y su negación como segunda decisión, formulando el problema con sólo dos alternativas. No es correcto, por ejemplo, plantearse si realizar o no la descarga de un buque granelero en determinadas condiciones climáticas. En efecto, además de detener la descarga del buque existen otras alternativas (utilizar instalaciones específicas, modificar el método de descarga, instalar pantallas protectoras,...) que deben ser consideradas.

Así pues, el primer paso para resolver un problema de decisión es elaborar el conjunto \mathcal{A} de las acciones a_i que podrían ser llevadas a cabo. La construcción del conjunto de acciones $\mathcal{A} = \{a_i, i \in I\}$ requiere una atención especial, porque el procedimiento que va a ser descrito se limitará a elegir uno de sus elementos. Formalmente, no es posible garantizar que se han incluido en \mathcal{A} todas las posibilidades interesantes; un buen decisor debe tener la creatividad y el conocimiento del tema suficientes para elaborar un conjunto de acciones *exhaustivo*, es decir que agote todas las posibilidades que, con la información disponible, pueda ser razonable tomar en consideración.

Es conveniente asimismo exigir que el conjunto de acciones esté formado por un conjunto de *alternativas*, de forma que la elección de uno de los elementos de \mathcal{A} excluya la elección de cualquier otro. Este planteamiento *no* supone pérdida de generalidad. Por ejemplo, la lista de medidas que pueden tomarse para atenuar la nube de partículas en suspensión producida por la descarga de un granel (tolvas especiales, irrigación,...) no es un conjunto de acciones adecuado, puesto que pueden ponerse simultáneamente en marcha varias de ellas, pero el conjunto de las partes de tal lista (el conjunto de todos sus subconjuntos) resulta serlo. De forma análoga, cualquier problema de decisión puede plantearse como el de la elección de un elemento, y uno sólo, de un conjunto apropiado de alternativas.

El principio, el conjunto \mathcal{A} puede contener infinitas alternativas. Por ejemplo, si debe decidirse cual es la presión más adecuada a la que debe mantenerse un tubo de descarga, el conjunto de acciones alternativas es el de todos los niveles de presión que soporta el sistema. Sin embargo, en la mayor parte de las aplicaciones \mathcal{A} es un conjunto finito, bien por serlo realmente, bien porque, con objeto de simplificar el problema, un conjunto de infinitas acciones alternativas es subdividido en un número finito de subconjuntos.

Sucesos inciertos relevantes. Determinar la mejor de un conjunto de alternativas es metodológicamente inmediato si se dispone de información completa sobre las consecuencias de cada una de ellas. El agente que debe reservar por anticipado determinados recursos no tendría problema si conociese de antemano la demanda a la que deberá atender. El médico que debe decidir un tratamiento ante un caso de alergia no dudaría si conociese de antemano con total exactitud las consecuencias que sufriría su paciente con cada uno de los tratamientos posibles. La dificultad principal que aparece al plantear un problema de decisión consiste en la falta de información precisa sobre lo que sucederá según se actúe de una u otra manera. El problema general de decisión se plantea así en *ambiente de incertidumbre*.

Existen situaciones en las que se tiene información completa y, sin embargo, es difícil tomar la decisión correcta, pero en estos casos la dificultad es de tipo técnico, no conceptual. Por ejemplo, a pesar de disponer de toda la información disponible, es difícil determinar la composición del pienso más barato que cumple

determinados requisitos de nutrición, pero se trata de un problema matemático bien definido (encontrar el mínimo condicionado de una función) sin que existan dudas sobre el *criterio de decisión* que debe adoptarse. En esta memoria no consideraremos tales dificultades técnicas: supondremos que en presencia de información completa siempre puede elegirse la mejor de un conjunto de acciones alternativas. Describiremos, en cambio, el *proceso lógico de decisión* en ambiente de incertidumbre, esto es, el *método* a seguir para tomar decisiones cuando *no* se dispone de toda la información que sería relevante.

Puesto que la dificultad esencial en un problema de decisión reside en los elementos inciertos presentes en la situación, es necesario considerar éstos con cuidado e introducirlos en la teoría. Así, para cada una de las posibles decisiones a_i , habrá que considerar el conjunto de los *sucesos inciertos relevantes* correspondientes

$$\Theta_i = \{\theta_{ij}, j \in J_i\}, \quad a_i \in \mathcal{A},$$

definido como el conjunto de aquellos sucesos θ_{ij} cuya ocurrencia permite determinar las eventuales *consecuencias* $c_{ij} = c(a_i, \theta_{ij})$ de optar por la acción a_i . Los sucesos inciertos deben constituir un conjunto *exhaustivo* de elementos *mutuamente excluyentes*, de forma que si se toma la decisión a_i , uno (y solamente uno) de los sucesos θ_{ij} debe tener lugar. Como en el caso de las decisiones, siempre puede conseguirse que los sucesos inciertos correspondientes a cada una de las acciones alternativas sean mutuamente excluyentes, pero la exhaustividad no es fácil de garantizar. La construcción de conjuntos de sucesos inciertos que realmente contemplen todas las eventualidades relevantes suele exigir un conocimiento importante del área de aplicación.

Si el conjunto de sucesos inciertos correspondiente a la alternativa a_i es un conjunto finito con m_i elementos, entonces tal alternativa a_i puede ser formalmente expresada como

$$a_i = \{c_{i1} | \theta_{i1}, c_{i2} | \theta_{i2}, \dots, c_{ij} | \theta_{ij}, \dots, c_{im_i} | \theta_{im_i}\}$$

de forma que a_i se *define* como la alternativa que da lugar a la consecuencia c_{i1} si sucede θ_{i1} , a la consecuencia c_{i2} si sucede θ_{i2} , y así sucesivamente. El caso más sencillo es el de las alternativas dicotómicas, de la forma $a_i = \{c_{i1} | \theta, c_{i2} | \theta^c\}$, donde θ^c representa el suceso complementario o negación de θ . La expresión general es de la forma $a_i = \{c_{i\theta} | \theta, \theta \in \Theta_i\}$, donde Θ_i es un espacio euclídeo de dimensión apropiada.

En el caso en el que tanto el conjunto de alternativas como los correspondientes conjuntos de sucesos inciertos sean finitos, el problema de decisión puede ser esquemáticamente descrito mediante un *árbol de decisión* como el de la Figura 1, en el que el cuadrado representa un nodo de decisión (donde debe elegirse una acción), y el círculo un nodo aleatorio, cuyo resultado no controla el decisor. Para cada

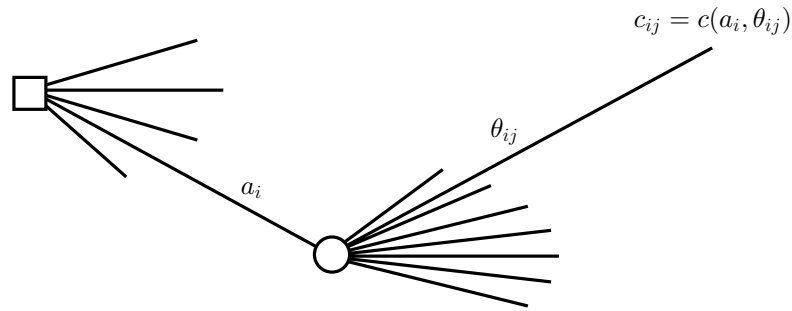


Figura 1. Árbol de decisión cualitativo

acción posible a_i , la ocurrencia de un suceso incierto θ_{ij} da lugar a una determinada consecuencia $c_{ij} = c(a_i, \theta_{ij})$.

Naturalmente, la mayor parte de los problemas reales involucran una sucesión de decisiones consecutivas, pero el análisis de tales problemas *secuenciales* puede ser reducido a un análisis repetido de estructuras simples como la descrita.

Solución intuitiva. Aunque los sucesos que componen cada uno de los conjuntos Θ_i son inciertos, en el sentido que no se sabe cuál de ellos tendrá lugar, los distintos sucesos posibles θ_{ij} no son (en general) igualmente verosímiles. Resulta intuitivamente obvia la necesidad de incorporar en el análisis la información de que se disponga sobre la relativa verosimilitud de los sucesos inciertos relevantes, lo que puede conseguirse determinando una *medida de probabilidad* que la describa. Por ejemplo, aunque no se disponga de información suficiente para determinar con precisión el contenido de partículas en suspensión que se observará en los alrededores del puerto como consecuencia de la descarga de un buque granelero, la lectura de los datos climáticos relevantes (intensidad y dirección del viento, humedad,...) y la medida de las características físicas del granel descargado permitirán determinar, haciendo uso de modelos adecuados, una distribución de probabilidad sobre el conjunto de sus valores posibles.

En el contexto de la teoría de la decisión, la *probabilidad* del suceso incierto θ_{ij} cuando se toma la decisión a_i en condiciones C , representada por $\Pr(\theta_{ij} | a_i, C)$, es una *medida* sobre una escala $(0, 1)$ de la verosimilitud de la ocurrencia de θ_{ij} en esas condiciones. Para una alternativa a_i con un número finito m_i de sucesos relevantes

$$0 \leq \Pr(\theta_{ij} | a_i, C) \leq 1, \quad \sum_{j=1}^{m_i} \Pr(\theta_{ij} | a_i, C) = 1,$$

de forma que la unidad de probabilidad se *distribuye* entre los m_i sucesos relevantes de los que se sabe que uno (y solamente uno) debe tener lugar.

Por otra parte, el decisor tendrá determinadas preferencias entre las distintas consecuencias y es también intuitivamente obvio que tales preferencias deben formar parte del análisis. En principio, las preferencias del decisor pueden ser cuantificadas asignando a cada una de las consecuencias c_{ij} un número $u(c_{ij})$ que mida la *utilidad* que cada una de ellas tenga para el decisor. Puede utilizarse cualquier escala que resulte conveniente, siempre que se utilice la misma escala para valorar todas las consecuencias posibles. Cualquiera que sea la escala elegida, $u(c_{ij}) = u(a_i, \theta_{ij})$ debe medir la *deseabilidad* de la consecuencia c_{ij} que se derivaría si se tomase la decisión a_i y sucediese θ_{ij} . Por ejemplo, podría asignársele una utilidad 100 a las consecuencias de una descarga granelera que no exija precauciones especiales y que no genere contaminación alguna, una utilidad 0 a las consecuencias de una descarga que genere niveles intolerables de contaminación en zonas habitadas, y valores intermedios a las consecuencias de métodos de descarga que utilicen determinados recursos para atenuar su efecto contaminante, dando lugar a niveles tolerables de contaminación atmosférica en zonas habitadas.

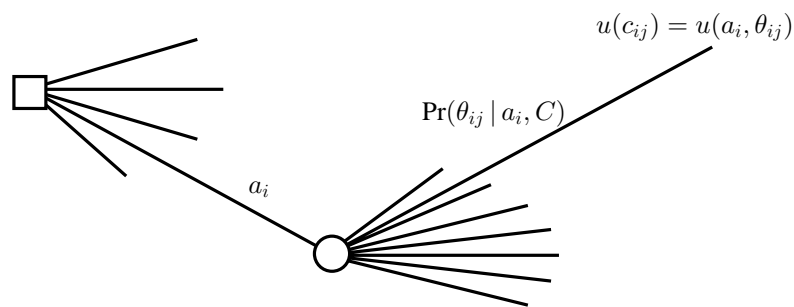


Figura 2. *Árbol de decisión cuantificado*

Una vez especificadas las *probabilidades* $\Pr(\theta_{ij} | a_i, C)$ que describen la relativa verosimilitud de los sucesos inciertos en las condiciones C en las que debe tomarse una decisión, y las *utilidades* $u(a_i, \theta_{ij})$ que describen las preferencias del decisor entre las posibles consecuencias, el problema de decisión plantado tiene una solución inmediata. En efecto, introduciendo en el árbol de decisión las probabilidades y las utilidades mencionadas (Figura 2) tendríamos una versión cuantitativa del problema de forma que la decisión a_i en las condiciones C daría lugar a una utilidad $u(a_i, \theta_{ij})$ con una probabilidad $\Pr(\theta_{ij} | a_i, C)$. Consecuentemente, en el caso finito, la utilidad *media* de la decisión a_i en las condiciones C , a la que llamaremos la *utilidad esperada* de a_i , vendrá dada por,

$$\bar{u}(a_i | C) = \sum_{j=1}^{m_i} u(a_i, \theta_{ij}) \Pr(\theta_{ij} | a_i, C),$$

que es simplemente la media de las utilidades que pueden alcanzarse si se toma la decisión a_i , ponderada con las probabilidades de que se alcancen tales valores. La acción óptima es aquella que *maximiza la utilidad esperada*.

Alternativa, pero equivalentemente, las preferencias del decisor puede ser descritas mediante una *función de pérdida* $l(a_i, \theta_{ij})$ que mida, en unidades que resulten adecuadas, la pérdida que sufriría el decisor si tomase la acción a_i y tuviese lugar el suceso incierto θ_{ij} . En este caso, la pérdida *media* de la decisión a_i en las condiciones C , a la que llamaremos la *pérdida esperada* de a_i , vendrá dada por,

$$\bar{l}(a_i | C) = \sum_{j=1}^{m_i} l(a_i, \theta_{ij}) \Pr(\theta_{ij} | a_i, C),$$

que es simplemente la media ponderada de las pérdidas que pueden alcanzarse si se toma la decisión a_i . La acción óptima es entonces aquella que *minimiza la pérdida esperada*.

El resultado más importante de la teoría de la decisión consiste en la demostración de que el procedimiento descrito, esto es:

- (i) cuantificar la incertidumbre mediante probabilidades,
- (ii) cuantificar las preferencias mediante utilidades,
- (iii) elegir aquella alternativa que maximice la utilidad esperada
(o minimice la pérdida esperada)

es, de hecho, el *único* procedimiento de toma de decisiones compatible con unos pocos principios, extremadamente intuitivos, de comportamiento racional.

Principios de coherencia. La teoría de la decisión es una teoría normativa. Partiendo de unos *axiomas* básicos o principios de coherencia que definen un *comportamiento racional*, demuestra la existencia de una *única* forma tomar decisiones (o, más precisamente, de ordenar las posibles alternativas) que es compatible con tales principios. Formalmente, el objeto de la teoría de la decisión es establecer las preferencias entre las acciones alternativas que *necesariamente* se deducen de comparaciones más sencillas realizadas entre los elementos básicos que forman parte del problema de decisión.

Es importante subrayar que los axiomas de comportamiento racional constituyen el fundamento de una teoría *normativa* de la decisión, no de una teoría *descriptiva*. No se trata de describir la forma en la que las personas toman habitualmente sus decisiones, sino de especificar la forma en que *deberían* tomarlas, si realmente pretenden evitar un comportamiento inconsistente.

La expresión formal de los principios de coherencia y las demostraciones matemáticas de los resultados a que dan lugar exceden ampliamente los límites de

esta memoria, pero describiremos brevemente su contenido intuitivo. Los lectores interesados en los detalles matemáticos pueden consultarlos, por ejemplo, en el texto de Bernardo & Smith (1994, Cap. 2).

- (i) *Comparabilidad.* El primer axioma afirma que existen al menos dos consecuencias c^* y c_* que no son igualmente deseables, y que siempre es posible elegir entre dos opciones dicotómicas, $\{c^* | \theta_1, c_* | \theta_1^c\}$ y $\{c^* | \theta_2, c_* | \theta_2^c\}$, basadas en ellas. La primera condición elimina los falsos problemas: si todas las consecuencias fuesen igualmente deseables, no habría problema de decisión en un sentido real. La segunda afirma que si se aspira a una elección racional entre opciones alternativas complicadas, entonces es necesario al menos expresar preferencias entre las posibilidades más sencillas (una forma racional de decisión exige saber lo que se quiere).
- (ii) *Transitividad.* El segundo axioma afirma la transitividad de las preferencias: se desea una teoría en la que si a_1 es preferible a a_2 , y a_2 es preferible a a_3 , entonces a_1 debe ser necesariamente preferible a a_3 .
- (iii) *Consistencia.* El tercer axioma afirma que las preferencias deben ser consistentes con la verosimilitud de los sucesos. Formalmente, si la consecuencia c^* es preferible a c_* y la opción $\{c^* | \theta_1, c_* | \theta_1^c\}$ es preferible a $\{c^* | \theta_2, c_* | \theta_2^c\}$, entonces θ_1 es juzgado por el decisor más verosímil que θ_2 .
- (iv) *Sucesos estándar.* El cuarto axioma afirma que es posible construir un conjunto de sucesos estándar con una medida verosimilitud conocida, y que tal medida satisface las propiedades inherentes a una medida de probabilidad. El ejemplo más sencillo es el proporcionado por la generación en ordenador de puntos aleatorios en el cuadrado unidad, de forma que la verosimilitud asociada al suceso de que el punto se sitúe en una determinada región del cuadrado es precisamente igual a su área. Los sucesos estándar actúan como unidades de medida para poder definir formalmente probabilidades y utilidades.
- (v) *Medida precisa.* El quinto y último axioma afirma la posibilidad de medir, por comparación con opciones construidas con sucesos estándar, la deseabilidad de cualquier consecuencia y la verosimilitud de cualquier suceso. Naturalmente, este axioma introduce un cierto grado de (inocua) idealización matemática, precisamente la misma que se utiliza en ingeniería al suponer que las medidas obtenidas son números reales, en lugar de números enteros en la unidad de medida más pequeña que permita apreciar el aparato de medida utilizado.

Existen muchas variaciones de los axiomas de comportamiento racional, que corresponden a distintos niveles de generalidad y a distintas formas de entender el concepto de “intuitivamente obvio”, pero sus consecuencias matemáticas son

siempre equivalentes. De los principios de comportamiento racional pueden deducirse resultados muy potentes. En particular, se demuestra que:

1. Es posible cuantificar la información de que se dispone sobre la verosimilitud relativa de los sucesos inciertos relevantes (también llamados *parámetros o variables de interés*, mediante una *medida de probabilidad* (una medida de incertidumbre que satisface las leyes aditiva y multiplicativa que definen el concepto matemático de probabilidad). En particular, pueden especificarse *distribuciones de probabilidad* que describen la verosimilitud relativa de los sucesos inciertos relevantes, correspondientes a cada una de las alternativas. Cada una de estas distribuciones de probabilidad se describe mediante el conjunto de probabilidades

$$\{0 \leq \Pr(\theta_{ij} | a_i, C) \leq 1, \quad j \in J_i\}, \quad \sum_{j \in J_i} \Pr(\theta_{ij} | a_i, C) = 1,$$

en el caso discreto, y mediante la densidad de probabilidad

$$p(\theta | a_i, C) \geq 0, \quad \theta \in \Theta_i, \quad \int_{\Theta_i} p(\theta | a_i, C) d\theta = 1,$$

en el caso continuo, de forma que la probabilidad de que θ se sitúe entre los valores a y b cuando se toma la decisión a_i en las condiciones C viene dada por

$$\Pr(a \leq \theta \leq b | a_i, C) = \int_a^b p(\theta | a_i, C) d\theta.$$

2. Es posible cuantificar las preferencias entre todas las consecuencias descritas en el problema mediante una *función de utilidad*,

$$\{u(a_i, \theta_{ij}), \quad i \in I, j \in J_i\}$$

en el caso discreto,

$$\{u(a_i, \theta), \quad i \in I, \theta \in \Theta_i\}$$

en el caso continuo, que asocia a cada consecuencia posible un número real que mide su deseabilidad.

3. La deseabilidad de cada una de las acciones alternativas, $a_i \in \mathcal{A}$, queda entonces descrita por su utilidad esperada $\bar{u}(a_i | C)$ que toma la forma

$$\bar{u}(a_i | C) = \sum_{j \in J_i} u(a_i, \theta_{ij}) \Pr(\theta_{ij} | a_i, C),$$

en el caso discreto y la forma

$$\bar{u}(a_i | C) = \int_{\Theta_i} u(a_i, \theta) p(\theta | a_i, C) d\theta$$

en el caso continuo. Consecuentemente, la decisión óptima en las condiciones C es aquella decisión $a^*(C)$ que maximiza la utilidad esperada en el conjunto de las alternativas consideradas, esto es,

$$a^*(C) = \arg \sup_{a_i \in \mathcal{A}} \bar{u}(a_i | C).$$

La metodología descrita proporciona la solución general a *cualquier* problema de decisión en términos de la función de utilidad que describe las preferencias del decisor y de las distribuciones de probabilidad que describen, para cada de las acciones alternativas a_i , la información disponible sobre los sucesos inciertos relevantes. La representación probabilística de esta información no es, sin embargo, inmediata, y constituye la aportación básica de los métodos estadísticos bayesianos a la solución de problemas de decisión.

2.2. DISCREPANCIA E INFORMACIÓN

El problema de predicción sobre el valor que acabará tomando una magnitud observable de interés $x \in \mathcal{X}$ puede ser descrito como un problema de decisión en el que el espacio de alternativas es el conjunto

$$\mathcal{A} = \left\{ p_x(\cdot), \quad p_x(x) > 0, \quad \int_{\mathcal{X}} p_x(x) dx = 1 \right\}$$

de las distribuciones de probabilidad sobre la magnitud a predecir, y el espacio de sucesos relevantes es el conjunto \mathcal{X} de los valores posibles de x . Para poder resolver este problema de decisión es necesario especificar la función $u\{p_x(\cdot), x\}$ que describe la utilidad que hubiera tenido la distribución predictiva $p_x(\cdot)$ si el valor de la magnitud de interés hubiese sido x .

Puede demostrarse (Savage, 1971; Bernardo, 1979a) que, bajo condiciones eminentemente razonables, la utilidad que puede esperarse de una distribución predictiva $\{p_1, \dots, p_k\}$ sobre el valor de una magnitud de interés x con un número finito de valores posibles $\mathcal{X} = \{x_1, \dots, x_k\}$, debe ser de la forma

$$u(\{p_1, \dots, p_k\}, x_j) = A \log[p_j] + B, \quad A > 0,$$

donde $p_j = \Pr(x = x_j)$, esto es debe ser una función lineal del logaritmo de la probabilidad asignada al verdadero valor. Consecuentemente, el incremento

de utilidad que puede esperarse de una predicción $\{p_1, \dots, p_k\}$ por encima de la predicción trivial $\{1/k, \dots, 1/k\}$ que da la misma probabilidad a todas las posibilidades resulta ser

$$\sum_{i=j}^k \left[u(\{p_1, \dots, p_k\}, x_j) - u(\{1/k, \dots, 1/k\}, x_j) \right] p_j = A \sum_{j=1}^k p_j \log \frac{p_j}{1/k},$$

lo que constituye una medida de la *cantidad de información* sobre el valor de la variable aleatoria discreta x que contiene la predicción $\{p_1, \dots, p_k\}$, tomando como origen la distribución uniforme $\{1/k, \dots, 1/k\}$.

En general, dado un vector aleatorio $\mathbf{x} \in \mathcal{X}$, la cantidad de información sobre \mathbf{x} contenida en una densidad de probabilidad p_x , tomando como origen otra densidad q_x , se define como

$$k\{q_x | p_x\} = \int_{\mathcal{X}} p_x(\mathbf{x}) \log \frac{p_x(\mathbf{x})}{q_x(\mathbf{x})} d\mathbf{x},$$

una cantidad no-negativa e invariante ante transformaciones biyectivas de \mathbf{x} , conocida como la *divergencia logarítmica dirigida* o *divergencia de Kullback-Leibler*, que separa q_x de p_x . La teoría matemática de la información proporciona una interpretación operativa de la divergencia logarítmica en términos de las unidades de información (*bits* si se utilizan logaritmos de base 2) necesarias para obtener p_x a partir de q_x , cuando p_x es la verdadera distribución de \mathbf{x} . Obsérvese que *no* se trata de una función simétrica de forma que, en general, $k\{q_x | p_x\} \neq k\{p_x | q_x\}$.

Una medida de discrepancia entre funciones es una función simétrica y no-negativa, que se anula si, y sólo si, las dos distribuciones son iguales casi por todas partes. La *discrepancia intrínseca* (Bernardo y Rueda, 2002), es una medida general de discrepancia entre distribuciones de probabilidad, que se define como el mínimo de sus dos posibles divergencias logarítmicas dirigidas:

$$\delta\{p_x, q_x\} = \min [k\{p_x | q_x\}, k\{q_x | p_x\}].$$

La función $\delta\{p_x, q_x\}$ es obviamente *simétrica*, y puede comprobarse que es *no-negativa*, y que se anula si, y solamente si, $p_x(x) = q_x(x)$ casi por todas partes. Se trata además de una función *aditiva* para variables aleatorias independientes, de forma que si $\mathbf{x} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, y $p_i(\mathbf{x}) = \prod_{l=1}^n q_l(\mathbf{y}_l)$, entonces

$$\delta\{p_1, p_2\} = n \delta\{q_1, q_2\}.$$

La discrepancia intrínseca es *invariante* ante transformaciones biyectivas, de forma que si $\mathbf{y} = \mathbf{y}(\mathbf{x})$ es una biyección y $q_i(\mathbf{y}) = p_i(\mathbf{x})/|J_{\mathbf{y}}|$, donde $J_{\mathbf{y}}$ el Jacobiano de la transformación, se verifica que $\delta\{p_1, p_2\} = \delta\{q_1, q_2\}$. Finalmente, la discrepancia

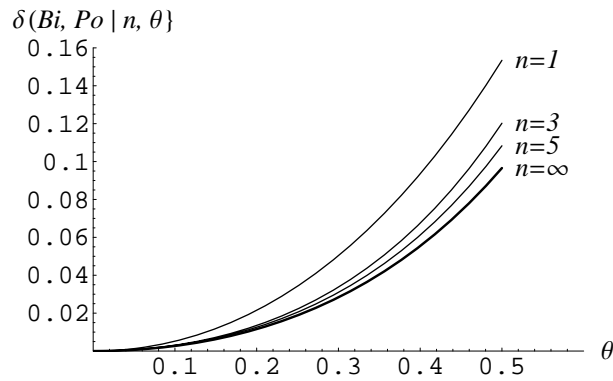


Figura 3 Discrepancia intrínseca asociada a la aproximación Poisson a una distribución binomial

intrínseca está bien definida incluso en el caso—muy frecuente en las aplicaciones— en el que las dos distribuciones tiene soportes estrictamente anidados, de forma que si $\mathcal{X}_i \subset \mathcal{X}_j$ (estrictamente) entonces $\delta\{p_i, p_j\} = \delta\{p_j, p_i\} = k\{p_j | p_i\}$.

Esta última propiedad permite estudiar aproximaciones definidas sobre distintos soportes. Por ejemplo, la Figura 3 representa en función de θ , y para distintos valores de n , la discrepancia intrínseca $\delta\{n, \theta\}$ entre una distribución binomial $\text{Bi}(x | \theta, n)$ y una Poisson $\text{Pn}(x | n\theta)$; la observación del comportamiento de $\delta\{n, \theta\}$ pone inmediatamente de manifiesto que, cualquiera que sea el valor de n , la condición necesaria y suficiente para que la aproximación funcione bien es que θ sea pequeño.

Si $\mathbf{x} \in \mathcal{X}$ es un vector aleatorio del que se sabe que su función de densidad de probabilidad es $p_1(\mathbf{x})$ o $p_2(\mathbf{x})$, entonces la discrepancia intrínseca $\delta\{p_1, p_2\}$ es el mínimo valor esperado del logaritmo del cociente de densidades $\log[p_i(\mathbf{x})/p_j(\mathbf{x})]$ en favor de la densidad verdadera. El particular, si $p_1(\mathbf{z})$ y $p_2(\mathbf{z})$ son modelos alternativos para un conjunto de datos $\mathbf{z} \in \mathcal{Z}$, y se supone que uno de los dos modelos es cierto, entonces $\delta\{p_1, p_2\}$ es el mínimo valor esperado del logaritmo del cociente de verosimilitudes en favor del modelo verdadero; esta propiedad es importante para el uso de la discrepancia intrínseca en el contexto de los problemas de contraste de hipótesis.

La discrepancia intrínseca permite definir un tipo de convergencia entre distribuciones de probabilidad especialmente útil en estadística matemática, la *convergencia intrínseca*. Se dice que una sucesión de densidades de probabilidad (funciones de probabilidad en el caso discreto) $\{p_i\}_{i=1}^{\infty}$ converge intrínsecamente a una densidad (función de probabilidad) p cuando la sucesión (de números reales

positivos) $\delta\{p_i, p\}_{i=1}^{\infty}$ converge a cero; formalmente,

$$\lim_{\text{int}} p_i = p \iff \lim_{i \rightarrow \infty} \delta\{p_i, p\} = 0.$$

Puede apreciarse que la discrepancia intrínseca $\delta\{p_{xy}, p_x p_y\}$ entre la densidad conjunta p_{xy} y el producto $p_x p_y$ de sus densidades marginales proporciona una medida general del *nivel de asociación* entre dos vectores aleatorios \mathbf{x} e \mathbf{y} , que se anula cuando \mathbf{x} e \mathbf{y} son independientes y tiende a infinito cuando la dependencia entre \mathbf{x} e \mathbf{y} tiende a una relación funcional, $\mathbf{y} = f(\mathbf{x})$. Obviamente, cualquier función biyectiva de la discrepancia intrínseca $\delta\{p_{xy}, p_x p_y\}$ es asimismo una medida general de la dependencia entre \mathbf{x} e \mathbf{y} .

El *coeficiente de asociación intrínseca* $\alpha(p_{xy})$ entre dos vectores aleatorios, $\mathbf{x} \in \mathcal{X}$, $\mathbf{y} \in \mathcal{Y}$, con densidad de probabilidad conjunta p_{xy} , definido por

$$\alpha(p_{xy}) = 1 - \exp(2\delta\{p_{xy}, p_x p_y\}),$$

es una medida general, en la escala $[0, 1]$ del grado de dependencia entre dos vectores aleatorios \mathbf{x} e \mathbf{y} . Si \mathbf{x} e \mathbf{y} son independientes, entonces $\alpha(p_{xy}) = 0$; si existe una dependencia funcional $\mathbf{y} = f(\mathbf{x})$, entonces $\alpha(p_{xy}) = 1$.

Si $\delta\{p_{xy}, p_x p_y\} = \log[10^k]$, entonces $\alpha(p_{xy}) = 1 - 10^{2k}$. En particular, si $\delta \approx \log[10]$, de forma que la densidad de probabilidad conjunta $p(\mathbf{x}, \mathbf{y})$ es, en valor medio, del orden de 10 veces mayor que el producto de las marginales $p(\mathbf{x})p(\mathbf{y})$, entonces $\alpha(p_{xy}) \approx 1 - 10^2 = 0.99$.

Si la distribución conjunta de dos variables aleatorias x e y es normal bivariable, su coeficiente de asociación intrínseca coincide con el coeficiente de determinación ρ^2 . En efecto, si la distribución conjunta de (x, y) es la normal bivariable

$$p(x, y | \mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = N_2 \left\{ \begin{pmatrix} x \\ y \end{pmatrix} \middle| \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \right\},$$

la discrepancia intrínseca entre $p(x, y | \mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ y el producto

$$p(x)p(y) = N(x | \mu_1, \sigma_1) N(y | \mu_2, \sigma_2)$$

de sus distribuciones marginales solamente depende de ρ , y resulta ser

$$\delta(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \delta(\rho) = -\frac{1}{2} \log(1 - \rho^2);$$

Consecuentemente, $\alpha(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = 1 - \exp[-2\delta(\rho)] = \rho^2$.

El frecuente (ab)uso del coeficiente de determinación como si se tratase una medida de asociación general entre dos variables aleatorias—cuando solo es

una medida de asociación apropiada bajo la (muy restrictiva) hipótesis de que su distribución conjunta sea normal bivalente—subraya la importancia de disponer de una medida general de asociación.

Si \mathbf{x} e \mathbf{y} son vectores aleatorios continuas y su distribución conjunta es regular $k\{p_{xy} | p_x p_y\}$ es siempre menor que $k\{p_x p_y | p_{xy}\}$ y, por lo tanto,

$$\delta\{p_{xy}, p_x p_y\} = \min\{k\{p_x p_y | p_{xy}\}, k\{p_{xy} | p_x p_y\}\} = k\{p_x p_y | p_{xy}\},$$

$$\alpha(p_{xy}) = 1 - \exp \left[2 \int_{\mathcal{X}\mathcal{Y}} p_{xy}(\mathbf{x}, \mathbf{y}) \log \frac{p_{xy}(\mathbf{x}, \mathbf{y})}{p_x(\mathbf{x}) p_y(\mathbf{y})} d\mathbf{x} d\mathbf{y} \right].$$

El coeficiente de asociación intrínseca entre dos vectores aleatorios \mathbf{x} e \mathbf{y} es una función matemática de su densidad de probabilidad conjunta $p(\mathbf{x}, \mathbf{y})$. Si solamente se dispone de una muestra aleatoria $\mathbf{z} = \{(\mathbf{x}_i, \mathbf{y}_i), i = 1, \dots, n\}$ de n pares de valores de tal distribución, su coeficiente de asociación intrínseca puede ser *estimado* por el método de Monte Carlo, de forma que

$$\hat{\alpha}(p_{xy}) = 1 - \exp \left[\frac{2}{n} \sum_{i=1}^n \log \frac{\hat{p}_{xy}(\mathbf{x}_i, \mathbf{y}_i)}{\hat{p}_x(\mathbf{x}_i) \hat{p}_y(\mathbf{y}_i)} \right],$$

donde \hat{p}_{xy} , \hat{p}_x y \hat{p}_y son estimaciones de las correspondientes densidades obtenidas a partir de los datos \mathbf{z} . En el caso continuo, y en ausencia de información sobre la posible estructura de la relación entre \mathbf{x} e \mathbf{y} , pueden utilizarse métodos no-paramétricos de estimación de densidades (ver, *e.g.*, Scott, 1992, Ch. 6).

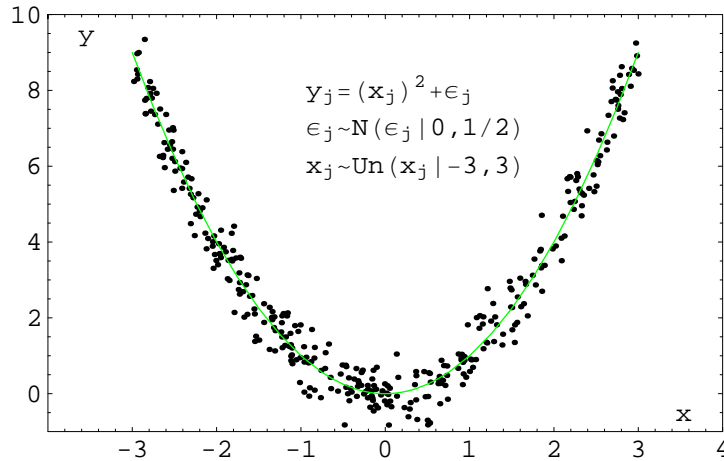


Figura 4 Conjunto de 350 observaciones aleatorias simuladas a partir de $N(y | x^2, 1)Un(x | -3, 3)$. Coeficiente de asociación intrínseca $\hat{\alpha}(z) = 0.986$; coeficiente de determinación $r^2 = 0.001$

En la Figura 4 se reproduce un conjunto de 350 pares de valores simulados a partir de la distribución $p(x, y) = N(y | x^2, 1)\text{Un}[x, -3, 3]$, que obviamente muestran un importante grado de dependencia, situándose alrededor de la función $y = x^2$. El verdadero valor de su coeficiente de asociación intrínseca (obtenido por integración numérica) es $\alpha(p_{xy}) = 0.986$, y un estimador por Monte Carlo (basado en una estimación no-paramétrica de las densidades p_{xy} , p_x y p_y) resultó ser $\hat{\alpha}(z) = 0.972$. Sin embargo, el coeficiente de determinación muestral es

$$\hat{\rho}^2(z) = r^2 = \frac{(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2} = 0.001,$$

perfectamente inútil para detectar la importante dependencia entre las dos variables aleatorias.

La discrepancia intrínseca entre dos *conjuntos* de distribuciones de probabilidad se define como la discrepancia intrínseca *mínima* entre elementos de ambos conjuntos. En particular, la discrepancia intrínseca $\delta(\mathcal{M}_1, \mathcal{M}_2)$ entre dos familias de modelos probabilísticos paramétricos definidos para $\mathbf{x} \in \mathcal{X}$,

$$\mathcal{M}_1 \equiv \{p(\mathbf{x} | \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}, \quad \mathcal{M}_2 \equiv \{q(\mathbf{x} | \boldsymbol{\phi}), \boldsymbol{\phi} \in \Phi\},$$

viene dada por

$$\delta\{\mathcal{M}_1, \mathcal{M}_2\} = \inf_{\boldsymbol{\theta} \in \Theta, \boldsymbol{\phi} \in \Phi} \delta\{p_x(\cdot | \boldsymbol{\theta}), q_x(\cdot | \boldsymbol{\phi})\},$$

lo que encuentra aplicaciones inmediatas en la formalización de los problemas convencionales de *estimación puntual* y de *contraste de hipótesis* (Bernardo y Rueda, 2002; Bernardo y Juárez, 2003).

El concepto de discrepancia intrínseca permite una definición general de la cantidad de información que puede esperarse de un experimento en función de la distribución inicial sobre sus parámetros.

Si el experimento consiste en la obtención de un conjunto de datos \mathbf{z} cuya densidad de probabilidad es un elemento de la familia

$$\mathcal{M} \equiv \{p(\mathbf{z} | \boldsymbol{\theta}), \mathbf{z} \in \mathcal{Z}, \boldsymbol{\theta} \in \Theta\},$$

y la distribución inicial de $\boldsymbol{\theta}$ es $p(\boldsymbol{\theta})$, la cantidad de información intrínseca que puede esperarse que proporcionen los datos \mathbf{z} sobre el valor de $\boldsymbol{\theta}$ es la discrepancia intrínseca $\delta\{p_{z\boldsymbol{\theta}}\}$ entre su densidad conjunta y el producto de sus densidades marginales de forma que

$$I\{p_{\boldsymbol{\theta}} | \mathcal{M}\} = \delta\{p(\mathbf{z}, \boldsymbol{\theta}), p(\mathbf{z})p(\boldsymbol{\theta})\}.$$

Bajo condiciones de regularidad, $\delta\{p(z, \theta), p(z)p(\theta)\} = k\{p(z, \theta) | p(z)p(\theta)\}$; en este caso, resulta

$$I\{p_\theta | \mathcal{M}\} = \int_{\mathcal{Z}} p(z) \int_{\Theta} p(\theta | z) \log\left[\frac{p(\theta | z)}{p(\theta)}\right] d\theta dz$$

y la información esperada intrínseca se reduce a la información esperada de Shannon.

La definición de información esperada de un experimento como función de la distribución inicial es crucial en la formulación de las distribuciones iniciales de referencia (Bernardo, 1979b, 2003, 2005), que constituyen la base de los métodos bayesianos objetivos.

2.3. MÉTODOS ESTADÍSTICOS BAYESIANOS

Sea Θ el conjunto de sucesos relevantes en el análisis de una determinada alternativa en un problema de decisión, de forma que $\theta \in \Theta$ es la variable o parámetro de interés. En el caso discreto, este conjunto será de la forma $\Theta = \{\theta_1, \theta_2, \dots\}$, y en el caso continuo $\Theta \subset \mathbb{R}^k$ será un subconjunto de un espacio euclídeo de dimensión finita k . Sea D un conjunto de *datos observados*, que presumiblemente aportan información relevante sobre el verdadero (y desconocido) valor de θ , y sean H las condiciones anteriores a la observación de los datos D .

Los resultados descritos arriba demuestran la *existencia* de una *distribución de probabilidad inicial*,

$$\{\Pr(\theta_1 | H), \Pr(\theta_2 | H), \dots\}, \quad \sum_j \Pr(\theta_j | H) = 1,$$

en el caso discreto, y

$$\{p(\theta | H), \theta \in \Theta\}, \quad \int_{\Theta} p(\theta | H) d\theta = 1,$$

en el caso continuo que describe la información de que se dispone sobre el valor de θ en las condiciones H anteriores a la observación de los datos D . Suponiendo que tal distribución ya ha sido determinada, se pretende encontrar la *distribución de probabilidad final*,

$$\{\Pr(\theta_1 | H, D), \Pr(\theta_2 | H, D), \dots\},$$

en el caso discreto, y

$$\{p(\theta | H, D), \quad \theta \in \Theta\}$$

en el caso continuo, que describe la información disponible sobre el valor de θ en el momento de tomar la decisión, esto es, en las condiciones $C = (H, D)$ en las que

se dispone tanto de la información inicial H como de la información proporcionada por los nuevos datos D .

El estudio estadístico de un conjunto D de *datos observados*, que presumiblemente aportan información relevante sobre el verdadero (y desconocido) valor de θ , suele empezar con un análisis descriptivo de su comportamiento, lo que permite sugerir un *modelo probabilístico* formal, $\{p(D|\theta, \theta \in \Theta)\}$, que describe (para el verdadero valor de θ) el mecanismo probabilístico que ha generado los datos observados D .

El *teorema de Bayes* (que da el nombre a los métodos bayesianos), permite obtener la distribución final buscada en términos de la distribución inicial y del modelo probabilístico. Específicamente,

$$\Pr(\theta_i | H, D) = \frac{p(D|\theta_i)\Pr(\theta_i|H)}{\sum_{j \in J} p(D|\theta_j)\Pr(\theta_j|H)}, \quad \theta_i \in \Theta,$$

en el caso discreto, y

$$p(\theta | H, D) = \frac{p(D|\theta)p(\theta|H)}{\int_{\Theta} p(D|\theta)p(\theta|H) d\theta}, \quad \theta \in \Theta,$$

en el caso continuo.

Ejemplo 1: Control de contaminación. Para verificar la posible contaminación por un agente químico de determinados productos hortícolas se dispone de un test del que se ha determinado en laboratorio que indica un resultado positivo en el 99% de los productos contaminados que se prueban (positivos correctos) y en el 2% de los no contaminados (falsos positivos). Si denotamos por θ_1 el suceso de que un producto esté contaminado y por θ_2 el suceso complementario de que no lo esté, la probabilidad final $\Pr(\theta_1 | H, +)$ de que un *determinado* producto este contaminado cuando el test ha dado positivo es

$$\begin{aligned} \Pr(\theta_1 | H, +) &= \frac{p(+|\theta_1)\Pr(\theta_1|H)}{p(+|\theta_1)\Pr(\theta_1|H) + p(+|\theta_2)\Pr(\theta_2|H)} \\ &= \frac{0.99p}{0.99p + 0.02(1-p)}, \quad p = \Pr(\theta_1|H), \end{aligned}$$

en función de la proporción p de productos contaminados entre los que forman parte del estudio. La Figura 5 muestra $\Pr(\theta_1 | H, +)$ en función de $p = \Pr(\theta_1 | H)$.

Como podría esperarse, la probabilidad final es igual a cero si (y solamente si) la probabilidad inicial es igual a cero, (de forma que *se sabe* que no hay ningún producto contaminado), y es igual a uno si (y solamente si) la probabilidad inicial es igual a uno, (de forma que *se sabe* que todos los productos están contaminados).

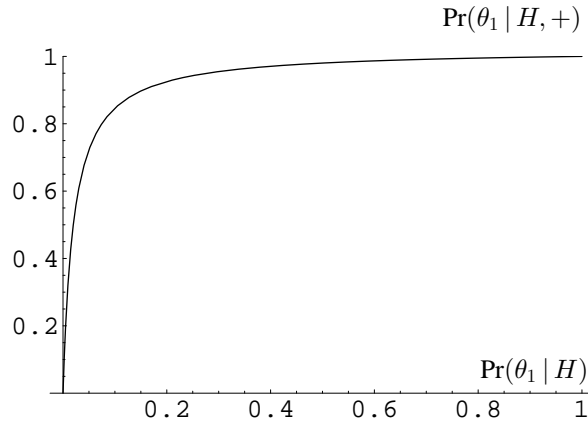


Figura 5. Probabilidad final de contaminación

Si un 50% de los productos están contaminados, de forma que $\Pr(\theta_1 | H) = 0.5$, entonces $\Pr(\theta_1 | H, +) = 0.980$ de forma que el 98% de los productos que den lugar a un test positivo estarán realmente contaminados. Obsérvese, sin embargo, que si la proporción $\Pr(\theta_1 | H)$ de productos contaminados entre los que son objeto de estudio es pequeña, entonces la probabilidad de que un producto escogido al azar esté contaminado será relativamente pequeña, *incluso* cuando el test haya dado positivo. Por ejemplo, si $\Pr(\theta_1 | H) = 0.005$, de forma que sólo el 0.5% están contaminados, resulta $\Pr(\theta_1 | H, +) = 0.199$, con lo que solamente el 19.9% de los productos que den lugar a un test positivo estarán realmente contaminados: la mayor parte de los resultados positivos serán *falsos* positivos.

Distribuciones iniciales de referencia. El teorema de Bayes permite la incorporación la información adicional sobre la variable de interés θ proporcionada por un conjunto de datos adicionales D en función del modelo que describe el comportamiento probabilístico de los datos, y de la distribución inicial de θ . Sin embargo, en muchos problemas no se dispone de información inicial sobre θ o esa información no es fácilmente objetivable y se quieren obtener conclusiones exclusivamente basadas en los datos observados D . En tales casos, es necesario especificar una *distribución inicial de referencia* $\pi(\theta)$ que describa matemáticamente la hipótesis de que *no* se dispone de información inicial sobre el valor de la cantidad de interés. La *teoría de la información* permite resolver este importante problema. En el caso particular en que θ sólo pueda tomar un número finito m de valores la solución es, como cabía esperar, la distribución uniforme,

$$\pi(\theta_i) = 1/m, \quad i = 1, \dots, m.$$

Cuando θ es una variable continua unidimensional y el modelo probabilístico suficientemente regular, la distribución de referencia viene dada por la fórmula de Jeffreys,

$$\pi(\theta) = \sqrt{i(\theta)}, \quad i(\theta) = - \int_D p(D|\theta) \frac{\partial^2}{\partial \theta^2} \log p(D|\theta) dD.$$

Para una introducción elemental a la metodología bayesiana objetiva, pueden consultarse Bernardo (2003) y Bernardo (2005). Para problemas más complicados ver, por ejemplo, Bernardo & Smith (1994, Cap. 5).

Ejemplo 2: Proporción de ciudadanos afectados. Para determinar la proporción de ciudadanos afectados por un episodio de contaminación atmosférica en un área determinada se analiza el estado de un conjunto de n ciudadanos residentes en esa área, aleatoriamente elegidos, y se encuentra que r de ellos han resultado afectados. Los datos observados D consisten pues en un conjunto de n observaciones Bernoulli con parámetro θ , la proporción de personas afectadas, de forma que

$$p(r|\theta, n) = \binom{n}{r} \theta^r (1-\theta)^{n-r}.$$

La fórmula de Jeffreys proporciona en este caso la distribución inicial de referencia

$$\pi(\theta) = \theta^{-1/2} (1-\theta)^{-1/2},$$

y utilizando el teorema de Bayes, la distribución final de referencia, que describe la información sobre θ obtenida exclusivamente a partir de los datos observados, resulta ser

$$\begin{aligned} \pi(\theta|r, n) &= \frac{p(r|\theta, n) \pi(\theta)}{\int_0^1 p(r|\theta, n) \pi(\theta) d\theta} \\ &= \frac{\Gamma(n+1)}{\Gamma(r+\frac{1}{2})\Gamma(n-r+\frac{1}{2})} \theta^{r-1/2} (1-\theta)^{n-r-1/2} \\ &= \text{Be}(\theta|r+\frac{1}{2}, n-r+\frac{1}{2}). \end{aligned}$$

donde $\Gamma(\cdot)$ es la función gamma, y $\text{Be}(\theta|a, b)$ denota una función beta de densidad de probabilidad con parámetros a y b . La distribución final de referencia es pues una distribución beta con parámetros $r+\frac{1}{2}$ y $n-r+\frac{1}{2}$, cuya media y desviación típica son, respectivamente,

$$\text{E}[\theta|r, n] = \frac{r+\frac{1}{2}}{n+1}, \quad \text{D}[\theta|r, n] = \sqrt{\frac{\text{E}[\theta|r, n](1-\text{E}[\theta|r, n])}{n+2}}.$$

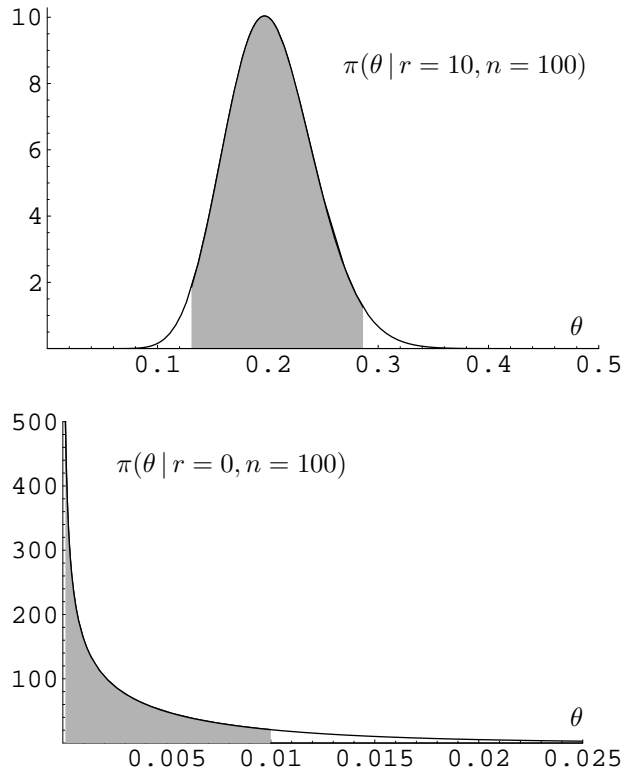


Figura 6. Distribuciones finales de la proporción de infectados

Si, por ejemplo, se observan $r = 20$ personas afectadas en una muestra de $n = 100$ residentes, la distribución final de referencia es $\text{Be}(\theta | 20.5, 80.5)$, representada en el panel superior de Figura 4. Consecuentemente, la media es $E[\theta | r, n] = 0.203$, y la desviación típica $D[\theta | r, n] = 0.040$; por tanto, la proporción θ de afectados debe situarse alrededor del $20.3 \pm 4.0\%$. Además, puede afirmarse, por ejemplo, que la proporción de afectados se sitúa entre el 13.1% y el 28.6% con probabilidad 0.95 (área sombreada en la figura).

De forma análoga, si no se observase *ninguna* persona afectada (de forma que $r = 0$) en una muestra del mismo tamaño ($n = 100$), la distribución final de referencia sería $\text{Be}(\theta | 0.5, 100.5)$, representada (en una escala muy distinta) en el panel inferior de la Figura 6. En particular, $\int_0^{0.01} \text{Be}(\theta | 0.5, 100.5) d\theta = 0.844$ y, por lo tanto, la probabilidad de que θ fuese menor del 1% sería 0.844 (área sombreada en la figura).

Parámetros marginales. En el desarrollo anterior, se ha supuesto que el modelo que describe el comportamiento probabilístico de los datos D dependía únicamente del parámetro de interés θ . En general, sin embargo, el modelo probabilístico es de la forma $p(D | \theta, \omega)$, que depende de θ y de un *parámetro marginal* $\omega \in \Omega$ (posiblemente multivariante). La solución en el caso general no presenta sin embargo nuevas dificultades metodológicas: basta hacer uso de la teoría de la probabilidad.

Como en el caso sencillo, hay que empezar por determinar la distribución de probabilidad que describe la información de que inicialmente se dispone sobre *todas* las variables desconocidas, ahora el conjunto $\{\theta, \omega\}$ (o la correspondiente distribución inicial de referencia si no se dispone de información inicial), y que será una *matriz* de probabilidades

$$\{0 \leq \Pr(\theta_j, \omega_k | H) \leq 1, \quad \theta_j \in \Theta, \omega_k \in \Omega\},$$

de suma total unidad, en el caso discreto, y una función de densidad de probabilidad multivariante

$$p(\theta, \omega | H) \geq 0, \quad \theta \in \Theta, \omega \in \Omega,$$

de integral unidad, en el caso continuo.

El teorema de Bayes permitirá entonces obtener la correspondiente distribución final *conjunta*, que será de la forma

$$\{0 \leq \Pr(\theta_j, \omega_k | H, D) \leq 1, \quad \theta_j \in \Theta, \omega_k \in \Omega\}$$

en el caso discreto, y de la forma

$$p(\theta, \omega, | H, D) \geq 0, \quad \theta \in \Theta, \omega \in \Omega$$

en el caso continuo. A partir de esta distribución final conjunta la distribución marginal de la magnitud de interés,

$$\{0 \leq \Pr(\theta_j | H, D) \leq 1, \quad \theta_j \in \Theta\}, \quad \sum_{\theta_j} \Pr(\theta_j | H, D) = 1$$

en el caso discreto, y

$$p(\theta | H, D) \geq 0, \quad \theta \in \Theta \quad \int_{\Theta} p(\theta | H, D) d\theta = 1$$

en el caso continuo, se obtiene inmediatamente sumando (en el caso discreto) o integrando (en el caso continuo) con respecto a los parámetros marginales.

Por ejemplo, en el caso continuo sin información inicial, la distribución final conjunta de referencia será

$$\pi(\theta, \omega | D) = \frac{p(D | \theta, \omega) \pi(\theta, \omega)}{\int_{\Theta} \int_{\Omega} p(D | \theta, \omega) \pi(\theta, \omega) d\theta d\omega}, \quad \theta \in \Theta, \omega \in \Omega,$$

y la correspondiente distribución final del parámetro de interés θ será

$$\pi(\theta | D) = \int_{\Omega} \pi(\theta, \omega | D) d\omega.$$

Las integrales necesarias sólo son analíticas en casos sencillos. Sin embargo, los métodos de integración numérica mediante cadenas markovianas de Monte Carlo (métodos MCMC) permiten obtener soluciones concretas en aplicaciones con miles de parámetros. El artículo de Gelfand y Smith (1990) y los libros de Gelman *et al.* (1995) y Gilks *et al.* (1996) constituyen una buena introducción a esta metodología.

Ejemplo 3: Densidad de partículas en suspensión. Para determinar la concentración media θ (en $\mu\text{g}/\text{m}^3$) de partículas en suspensión en el área portuaria en un determinado momento, se dispone de las medidas $D = \{x_1, \dots, x_n\}$ proporcionadas por n captadores apropiadamente situados. Bajo condiciones razonablemente generales, puede suponerse que tales observaciones constituyen una muestra aleatoria de tamaño n de una distribución normal de media desconocida θ (el parámetro de interés), y desviación típica desconocida σ (un parámetro marginal). Consecuentemente, el modelo probabilístico es de la forma $p(D | \theta, \sigma) = \prod_j N(x_j | \theta, \sigma)$. Puede demostrarse (ver, por ejemplo, Bernardo & Smith, 1994, p. 328), que la distribución inicial de referencia correspondiente (que describe matemáticamente una situación en la que no se dispone de información inicial relevante) es $\pi(\theta, \sigma) = \sigma^{-1}$, y que la correspondiente distribución final marginal de la variable de interés θ es la distribución de Student

$$\pi(\theta | D) = \pi(\theta | \bar{x}, s, n) = \text{St}(\theta | \bar{x}, \frac{s}{\sqrt{n-1}}, n-1),$$

con media $\bar{x} = n^{-1} \sum_j x_j$, parámetro de escala $(n-1)^{-1/2} s$ y $n-1$ grados de libertad, donde $s = n^{-1} \sum_j (x_j - \bar{x})^2$ es la desviación típica muestral.

En particular, si se dispone de $n = 4$ medidores, que dan los valores $\{30, 39, 52, 48\}$, (con media $\bar{x} = 42.45$, y desviación típica $s = 8.50$), la distribución final de θ , representada en la Figura 7, es $\text{St}(\theta | 42.25, 4.91, 3)$, de forma que puede afirmarse que la concentración media de partículas θ se sitúa alrededor de $42.25 \pm 4.91 \mu\text{g}/\text{m}^3$. Además, como puede comprobarse por integración directa (o utilizando las tablas de la función de distribución Student para la variable tipificada $\sqrt{(n-1)}(\theta - \bar{x})/s$),

$$\int_{26.64}^{57.86} \text{St}(\theta | 42.25, 4.91, 3) d\theta = 0.95,$$

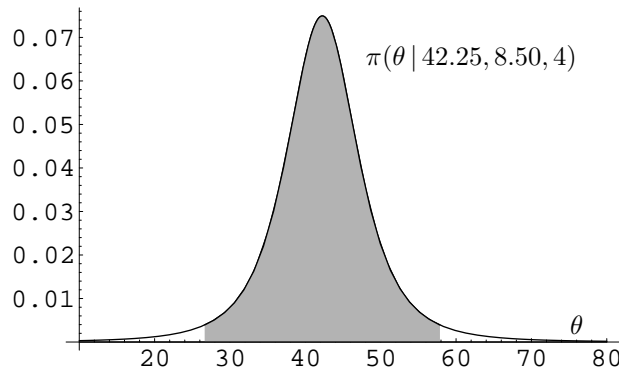


Figura 7. Distribución final de la concentración media

de forma que, dados los resultados observados (y haciendo uso únicamente de esa información) la probabilidad de que θ se sitúe entre 26.64 y $57.86 \mu\text{gr}/\text{m}^3$ es 0.95 (área sombreada en la figura).

Predicción. En numerosas ocasiones, las consecuencias de un problema de decisión no dependen de valores paramétricos (no observables) sino del valor de observaciones futuras. Por ejemplo, las consecuencias de una determinada forma de descarga de un buque granelero pueden depender del tiempo x que resulte necesario para descargarlo y, consecuentemente, deberá obtenerse una distribución de probabilidad $p(x | x_1, \dots, x_n)$ que permita *predecir* el valor de x dada la información proporcionada por un conjunto $D = \{x_1, \dots, x_n\}$ de observaciones anteriores realizadas en condiciones comparables.

En general, dado un conjunto de observaciones $D = \{x_1, \dots, x_n\}$, cuyo comportamiento probabilístico se supone que viene descrito por un determinado modelo probabilístico $\{p(D | \theta), \theta \in \Theta\}$, y dada una función de interés arbitraria $y = y(x_{n+1}, \dots, x_{n+m})$ de m observaciones futuras del mismo tipo, el problema es predecir el valor de y con la información proporcionada por D . Formalmente, el problema es de encontrar la *distribución predictiva* $p(y | x_1, \dots, x_n)$ que describe la información de que se dispone para predecir el valor de y en base a los datos $D = \{x_1, \dots, x_n\}$ y al modelo supuesto.

Si, como en el ejemplo mencionado, las observaciones x_i son *intercambiables* (condicionalmente independientes), de forma que las observaciones constituyen una muestra aleatoria de un modelo $\{p(x | \theta), \theta \in \Theta\}$, y la función a estimar es simplemente una nueva observación x , entonces el problema planteado tiene una solución sencilla. En efecto, en virtud del teorema de la probabilidad total,

$$p(x | x_1, \dots, x_n) = \sum_{\theta_i \in \Theta} p(x | \theta_i) \Pr(\theta_i | x_1, \dots, x_n)$$

en el caso discreto, y

$$p(x | x_1, \dots, x_n) = \int_{\Theta} p(x | \theta) p(\theta | x_1, \dots, x_n) d\theta,$$

en el caso continuo, donde $\Pr(\theta_i | x_1, \dots, x_n)$ y $p(\theta | x_1, \dots, x_n)$ son las dos formas posibles de la distribución final de θ , cuya obtención ya ha sido descrita.

En los problemas de predicción es frecuente contar con la información proporcionada por covariables relevantes que permiten obtener predicciones más precisas. En estos casos, los datos están constituidos por una colección de pares $D = \{(x_1, \mathbf{y}_1), \dots, (x_n, \mathbf{y}_n)\}$ (donde tanto las x_i como las \mathbf{y}_i pueden ser multivariantes). Dada una nueva observación, de la que se conocen los valores de las covariables x , se trata de predecir el valor correspondiente de \mathbf{y} utilizando la información proporcionada por los datos D . Formalmente, se trata de determinar la distribución predictiva $p(\mathbf{y} | x, D)$. El primer paso es establecer un modelo $\{p(\mathbf{y} | x, \theta), \theta \in \Theta\}$ que precise la relación probabilística entre las covariables x y la variable objeto de estudio \mathbf{y} .

El modelo más sencillo postula una relación lineal de una variable unidimensional y y un vector x de k covariables del tipo

$$p(y | x, \theta) = N(y | x^t \theta, \sigma),$$

pero es fácil proponer ejemplos que requieren modelos mucho más complejos.

En el caso lineal con estructura de errores normal, la distribución inicial de referencia es $\pi(\theta, \sigma) = \sigma^{-1}$, y la distribución predictiva resulta ser una distribución de Student con $n - k$ grados de libertad, donde k es la dimensión de θ , truncada y renormalizada al conjunto de valores admisibles de y . Específicamente, si los datos consisten en el conjunto $\{(y_i, \mathbf{x}_i), i = 1, \dots, n\}$, la distribución predictiva del valor de y que corresponde a un nuevo vector x , dados el vector $\mathbf{y} = \{y_1, \dots, y_n\}^t$ y la matriz \mathbf{X} de tamaño $n \times k$ cuyas filas son los vectores $\mathbf{x}_j = \{x_{j1}, \dots, x_{jk}\}$, es

$$p(y | x, \mathbf{y}, \mathbf{X}) = \text{St}(y | x \hat{\theta}, s \sqrt{\frac{nf}{n-k}}, n-k),$$

$$\hat{\theta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}, \quad f = 1 + x (\mathbf{X}^t \mathbf{X})^{-1} x^t,$$

con $ns^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i \hat{\theta})^t (y_i - \mathbf{x}_i \hat{\theta})$.

Si sólo se considera una covariable y se incluye un término independiente, de forma que $k = 2$, $p(y | x, \theta, \sigma) = N(y | \theta_1 + \theta_2 x, \sigma)$, y $\mathbf{x}_i = (1, x_i)$, las fórmulas anteriores se reducen a

$$p(y | x, \mathbf{y}, \mathbf{X}) = \text{St}(y | \hat{\theta}_1 + \hat{\theta}_2 x, s \sqrt{\frac{nf}{n-2}}, n-2),$$

$$\hat{\theta}_1 = \bar{y} - \hat{\theta}_2 \bar{x}, \quad \hat{\theta}_2 = \frac{s_{xy}}{s_x^2}, \quad f = 1 + \frac{1}{n} \frac{(x - \bar{x})^2 + s_x^2}{s_x^2}$$

donde $ns^2 = \sum_{i=1}^n (y_i - \hat{\theta}_1 - \hat{\theta}_2 x_i)^2$, y, con una notación convencional, $n\bar{x} = \sum_i x_i$, $n\bar{y} = \sum_i y_i$, $ns_{xy} = \sum_i (x_i - \bar{x})(y_i - \bar{y})$, y $ns_x^2 = \sum_i (x_i - \bar{x})^2$.

Ejemplo 6: Contaminación en función del viento. Supongamos que en una ciudad portuaria se desean predecir los niveles de contaminación atmosférica en función de la intensidad del viento.

x	3.7	0.5	0.1	1.0	0.3	2.6	3.9	4.6	2.4	0.8	0.1	3.0
y	1062	188	78	280	163	708	1122	1216	524	347	79	694

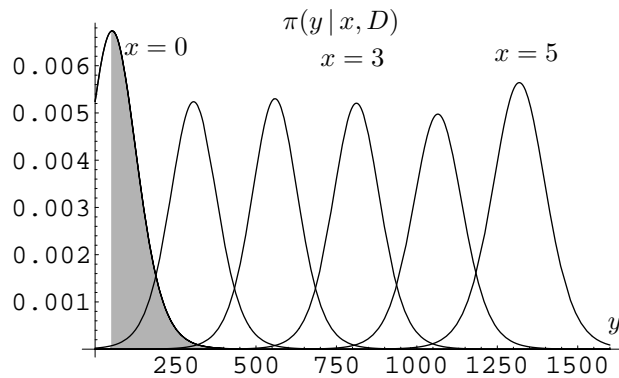
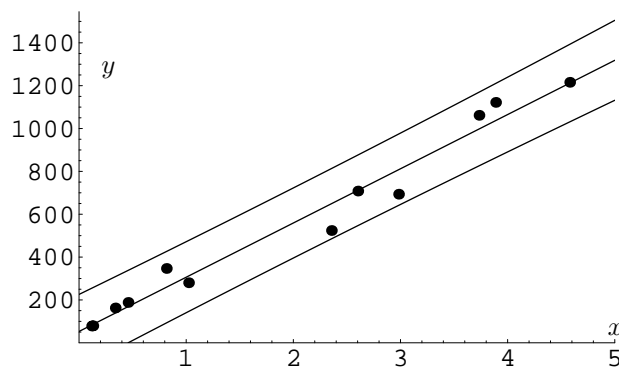


Figura 8. Concentración de partículas en función del viento

Los datos D disponibles sobre la cantidad de partículas en suspensión en el punto de la ciudad más cercano a su puerto (y_i , en $\mu\text{g}/\text{m}^3$), y la velocidad del viento procedente de la zona de descarga del puerto (x_i , módulo de la velocidad en m/s por el coseno del ángulo de separación), en un total de $n = 12$ episodios de descarga granelera son los reproducidos y representados en la Figura 6.

Utilizando la metodología antes descrita, puede determinarse la distribución predictiva de la cantidad de partículas en suspensión correspondiente a *cualquier* velocidad del viento. En el panel superior de la Figura 8 se han superpuesto la *recta de regresión*, $y = \hat{\theta}_1 + \hat{\theta}_2 x$, y las líneas (ligeramente divergentes) correspondientes a los intervalos predictivos con contenido probabilístico 0.95. En el panel inferior de la Figura 8 se representan las distribuciones predictivas correspondientes a 0, 1, 2, 3, 4, y 5 m/s. La distribución correspondiente al caso en que no hay nada de viento procedente del puerto ($x = 0$, con lo que se mide en nivel subyacente de partículas en suspensión debido a otras causas) merece especial atención. Se trata de una distribución Student truncada y renormalizada a valores no-negativos, con parámetro de localización $53\mu\text{g}/\text{m}^3$, parámetro de escala $78\mu\text{g}/\text{m}^3$, y con 18 grados de libertad; la probabilidad de que en ese caso se superen los $50\mu\text{g}/\text{m}^3$ (un nivel generalmente considerado medio-alto) es en ese caso 0.69 (área sombreada bajo la curva correspondiente).

Un conjunto de problemas de predicción especialmente interesantes son los planteados por conjuntos de datos temporales, en los que se trata de predecir un valor futuro y_{n+k} de una serie temporal, conocidos sus últimos n valores $D = \{y_1, \dots, y_n\}$. En este caso las observaciones y_i ya *no* son intercambiables y la determinación de la correspondiente distribución predictiva $p(y_{n+k} | y_1, \dots, y_n)$ es bastante más compleja. La descripción técnica de este tipo de problemas, el *análisis bayesiano de series temporales* excede ampliamente los límites de esta sencilla introducción. El lector interesado puede consultar la excelente monografía de West & Harrison (1989).

2.4. BIBLIOGRAFÍA COMENTADA

- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. Berlin: Springer.
[Una descripción muy completa de los métodos bayesianos que enfatiza su componente normativa en términos de teoría de la decisión]
- Bernardo, J. M. (1979a). Expected information as expected utility. *Ann. Statist.* **7**, 686–690.
[Demuestra que la inferencia estadística tiene la estructura matemática de un problema de decisión en el que la función de utilidad es una medida de información]

- Bernardo, J. M. (1979b). Reference posterior distributions for Bayesian inference. *J. Roy. Statist. Soc. B* **41**, 113–147 (con discusión).
Reimpreso en *Bayesian Inference* (N. G. Polson and G. C. Tiao, eds.), Brookfield, VT: Edward Elgar, (1995), 229–263.
[Artículo en el que se introdujeron las distribuciones iniciales de referencia, base de la mayor parte de los métodos bayesianos objetivos actuales]
- Bernardo, J. M. (1997). Noninformative priors do not exist. *J. Statist. Planning and Inference* **65**, 159–189 (con discusión).
[Una descripción elemental de la polémica académica en torno a los métodos bayesianos objetivos]
- Bernardo, J. M. (1999). Nested hypothesis testing: The Bayesian reference criterion. *Bayesian Statistics 6* (J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.). Oxford: University Press, 101–130 (con discusión).
[Una aproximación al contraste de hipótesis precisas basada en la teoría de la información y en la teoría bayesiana de la decisión]
- Bernardo, J. M. (2003). Bayesian Statistics. *Encyclopedia of Life Support Systems (EOLSS). Probability and Statistics* (R. Viertl, ed). London, UK: UNESCO.
[Un largo artículo, en uno de los tomos de matemáticas de la enciclopedia de la UNESCO, que proporciona una introducción elemental a los métodos estadísticos bayesianos. El apéndice a esta memoria contiene un resumen actualizado del contenido de este trabajo.]
- Bernardo, J. M. (2005). Reference analysis. *Handbook of Statistics* **25** (D. Dey & C. R. Rao, eds). Amsterdam: North Holland (en prensa).
<www.uv.es/~bernardo/RefAna.pdf>
[Un largo artículo, en el volumen dedicado a la metodología bayesiana de una enciclopedia de Estadística, que proporciona una introducción a los métodos bayesianos objetivos, con especial énfasis en la definición y construcción de distribuciones iniciales de referencia.]
- Bernardo, J. M. and Juárez, M. (2003). Intrinsic Estimation. *Bayesian Statistics 7* (J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith and M. West, eds.). Oxford: University Press, 465–476. [Una nueva formulación del problema de estimación puntual, que resulta en estimadores bayesianos objetivos invariantes]
- Bernardo, J. M. & Ramón, J. M. (1998). An introduction to Bayesian reference analysis: inference on the ratio of multinomial parameters. *The Statistician* **47**, 1–35.
[Una introducción elemental a los métodos bayesianos objetivos]

- Bernardo, J. M. and Rueda, R. (2002). Bayesian hypothesis testing: A reference approach. *International Statistical Review* **70**, 351-372. [Una nueva formulación del problema del contrast de hipótesis que utiliza funciones de pérdida basadas en la teoría de la información]
- Bernardo, J. M. & Smith, A. F. M. (1994). *Bayesian Theory*, Chichester: Wiley.
[Un tratado detallado, a nivel de posgrado, de los conceptos y resultados más relevantes de la estadística bayesiana, con una bibliografía muy extensa]
- Bernardo, J. M., Berger, J. O., Dawid, A. P. & Smith, A. F. M. (Eds.) (1999). *Bayesian Statistics 6*, Oxford: University Press.
[*Proceedings* del 6th Valencia International Meeting on Bayesian Statistics; los Valencia Meetings, que tienen lugar cada cuatro años, recogen los últimos avances en investigación y aplicaciones, aparecidos dentro del paradigma Bayesiano]
- Berry, D. A. (1996). *Statistics, a Bayesian Perspective*. Belmont, CA: Wadsworth.
[Excelente introducción a la estadística bayesiana desde un planteamiento subjetivista]
- Box, G. E. P. & Tiao, G. C. (1973). *Bayesian Inference in Statistical Analysis*. Reading, MA: Addison-Wesley.
[Un excelente clásico, todavía útil, de la metodología bayesiana objetiva]
- DeGroot, M. H. (1970). *Optimal Statistical Decisions*, New York: McGraw-Hill.
[Uno de los mejores tratados de teoría bayesiana de la decisión de todos los tiempos, que incluye un riguroso tratamiento matemático de sus fundamentos]
- DeGroot, M. H. & Schervish, M. J. (2002). *Probability and Statistics*, (3rd ed). Reading, MA: Addison-Wesley.
[Un texto ecléctico de introducción a la estadística matemática que describe tanto métodos frecuentistas como bayesianos. La tercera edición de un clásico]
- Efron, B. (1986). Why isn't everyone a Bayesian? *Amer. Statist.* **40**, 1-11 (con discusión).
[Un buen ejemplo de la polémica académica entre las aproximaciones bayesiana y frecuentista a la estadística matemática]
- de Finetti, B. (1970). *Teoria delle Probabilità*, Turin: Einaudi. Traducido como *Theory of Probability* in 1975, Chichester: Wiley.
[Un extraordinario texto en probabilidad y estadística desde una perspectiva subjetiva]
- Geisser, S. (1993). *Predictive Inference: An Introduction*. London: Chapman and Hall.
[Un análisis comparativo de los métodos frecuentistas y bayesianos de abordar los problemas de predicción]

- Gelfand, A. E. & Smith, A. F. M. (1990). Sampling based approaches to calculating marginal densities. *J. Amer. Statist. Assoc.* **85**, 398-409.
[Una excelente introducción a las técnicas de integración numérica por simulación en el contexto de la estadística bayesiana]
- Gelman, A., Carlin, J. B., Stern, H. & Rubin, D. B. (2004). *Bayesian Data Analysis* (2nd ed). London: Chapman and Hall.
[Un tratamiento muy completo del análisis de datos desde una perspectiva bayesiana, con énfasis en los procedimientos de simulación]
- Gilks, W. R., Richardson, S. & Spiegelhalter, D. J. (1996). *Markov Chain Monte Carlo in Practice*. London: Chapman and Hall.
[Excelente introducción a los métodos de integración por Monte Carlo y a sus aplicaciones]
- Kass, R. E. & Raftery, A. E. (1995). Bayes factors. *J. Amer. Statist. Assoc.* **90**, 773-795.
[Una buena revisión de los métodos de contraste de hipótesis basados en factores de Bayes]
- Lindley, D. V. (1972). *Bayesian Statistics, a Review*. Philadelphia, PA: SIAM.
[Una brillante y concisa síntesis de la bibliografía bayesiana hasta los años 70, enfatizando su consistencia interna]
- Lindley, D. V. (1985). *Making Decisions*. (2nd ed.) Chichester: Wiley.
[La mejor introducción elemental a la teoría bayesiana de la decisión]
- Lindley, D. V. (1990). The present position in Bayesian statistics. *Statist. Sci.* **5**, 44-89 (con discusión).
[Una interesante descripción del paradigma Bayesiano y de su relación con otras escuelas de inferencia estadística]
- Lindley, D. V. (2000). The philosophy of statistics. *The Statistician* **49**, 293-337 (con discusión).
[Una descripción reciente del paradigma bayesiano desde una perspectiva subjetivista]
- O'Hagan, A. (1994). *Bayesian Inference*. London: Edward Arnold.
[Una buena descripción de la inferencia bayesiana integrada en la colección Kendall de Estadística]
- Press, S. J. (1972). *Applied Multivariate Analysis: using Bayesian and Frequentist Methods of Inference*. Melbourne, FL: Krieger.
[Un tratamiento muy completo de los métodos estadísticos multivariantes, en el que se comparan las soluciones bayesianas con las soluciones frecuentistas]
- Robert, C. P. (2001). *The Bayesian Choice* (2nd ed). Berlin: Springer.
[Una buena introducción a la metodología estadística bayesiana que enfatiza sus fundamentos en la teoría de la decisión]

- Savage, L. J. (1971). Elicitation of personal probabilities and expectations. *J. Amer. Statist. Assoc.* **66**, 781–801. Reimpreso en *Studies in Bayesian Econometrics and Statistics: in Honor of Leonard J. Savage* (S. E. Fienberg and A. Zellner, eds.). Amsterdam: North-Holland, 111–156 (1974).
[Artículo póstumo de L. J. Savage en que establece condiciones para la unicidad de la función de evaluación logarítmica en el caso discreto finito]
- Scott, D. W. (1992). *Multivariate Density Estimation*. Chichester: Wiley.
[Una buena introducción a los métodos convencionales de estimación no-paramétrica de densidades de probabilidad]
- West, M. & Harrison, P. J. (1989). *Bayesian Forecasting and Dynamic Models*. Berlin: Springer.
[Una excelente descripción de los métodos bayesianos de análisis de series temporales]
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics*. New York: Wiley. Reprinted in 1987, Melbourne, FL: Krieger.
[Una detallada descripción del análisis bayesiano objetivo de modelos estadísticos lineales y de sus aplicaciones en econometría]

Capítulo 3.

Algunos Problemas Medioambientales Portuarios

En este capítulo se indica, muy sucintamente, la formulación bayesiana de algunos de los problemas de decisión planteados por las autoridades de los puertos que forman parte de este proyecto. Como el lector podrá comprobar, la solución efectiva de cualquiera de ellos exige consagrarles tiempo y recursos de cierta envergadura, por lo que su ejecución probablemente requerirá un proyecto específico en cada caso.

3.1. RED DE CONTROL AMBIENTAL EN A CORUÑA

Por su ubicación en el corazón de la ciudad (Figura 9), la carga y descarga en el puerto de La Coruña, especialmente en los muelles del Centenario y San Diego producen nubes de carbón y de soja que, en determinadas condiciones meteorológicas pueden desplazarse hacia la ciudad. Como consecuencia, las autoridades portuarias, con el apoyo técnico de la Universidad de Santiago de Compostela, tienen previsto establecer una red predicción de inmisión de partículas que permita alertar sobre las condiciones en las que podrían generarse niveles de contaminación que superen los límites establecidos por la legislación vigente. Esto plantea el problema de decidir el número y la posible localización de las estaciones de medida de partículas en suspensión de las que deberá depender el proceso de alerta.

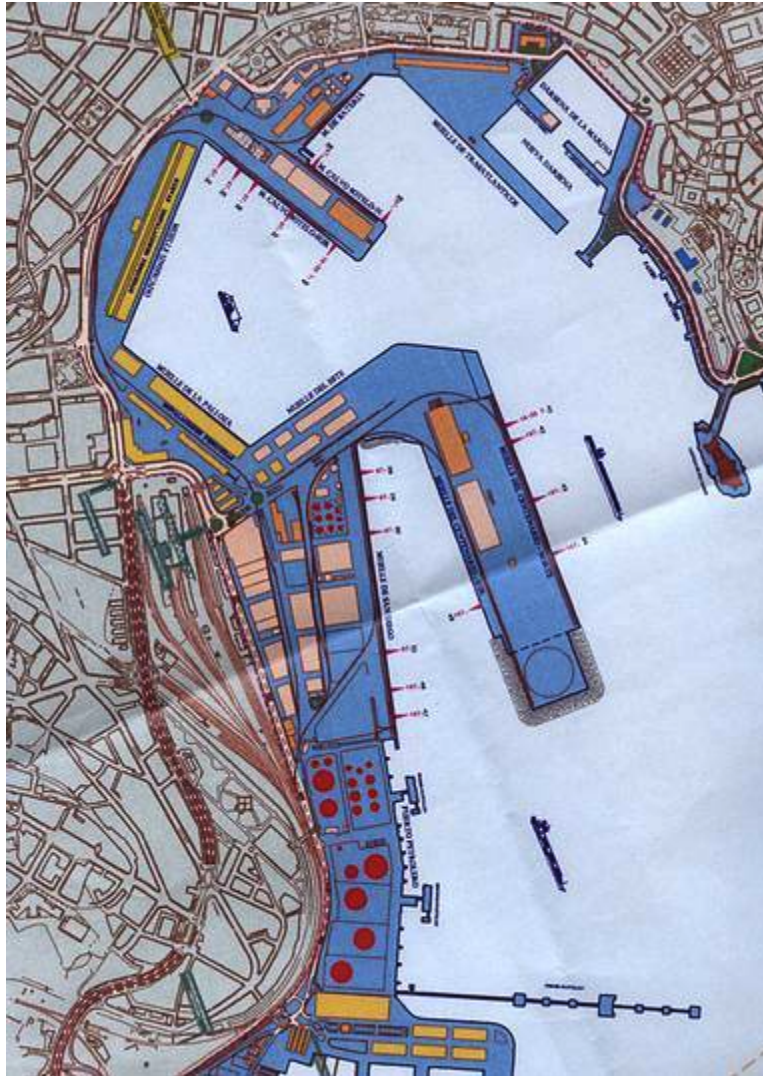


Figura 9. Estructura del puerto interior de A Coruña.

Formalmente, se trata de un problema de decisión cuyo espacio de alternativas es de la forma $\mathcal{A} = \{a_i, i \in N\}$, donde $a_i = a_i(\alpha_i, \beta_i)$ es la alternativa $a_i \equiv \{(\alpha_{ij}, \beta_{ij}), j = 1, \dots, i\}$ que consiste en instalar un total de i estaciones de medida, y situarlas precisamente en las coordenadas $(\alpha_{i1}, \beta_{i1}), \dots, (\alpha_{ii}, \beta_{ii})$ de una determinada malla topográfica de la zona.

El objetivo declarado de la red de control es el de prever las consecuencias de una posible operación de carga o descarga con tiempo suficiente como para emitir una alerta que permita suspenderla preventivamente para evitar niveles de contaminación que violarían la legislación en vigor. Consecuentemente, la magnitud de interés es precisamente la probabilidad

$$\theta_i(\alpha_i, \beta_i) = \Pr(A | a_i, C, B)$$

de que, en las condiciones C en las que una operación rutinaria de carga o descarga generaría niveles de contaminación inaceptables, se hubiese disparado una alarma con tiempo suficiente (suceso A), si se hubiese dispuesto i medidores situados en los lugares (α_i, β_i) , y de una base de datos locales B .

Ignorando de momento el coste de instalación, la función de utilidad será de tipo cero-uno, con el valor uno si la probabilidad de producir una alarma necesaria es suficientemente alta y cero en caso contrario. Formalmente,

$$u(a_i, \theta_i) = \begin{cases} 1, & \theta_i \geq 1 - \epsilon \\ 0, & \theta_i < 1 - \epsilon \end{cases}$$

donde $\epsilon > 0$ es el margen de error (la probabilidad de no prever un episodio de contaminación de niveles intolerables) que el decisor está dispuesto a asumir.

La combinación de un estudio de dispersión de partículas en función de las variables meteorológicas relevantes, como el realizado por el Grupo de Física No Lineal de la Universidad de Santiago, con un modelo de predicción meteorológica a corto plazo basado en un banco de datos locales B apropiado permitiría determinar las distribuciones de probabilidad $p(\theta_i | a_i, B)$ de las magnitudes de interés θ_i . Para cada posible alternativa, la utilidad esperada será

$$\bar{u}[a_i(\alpha_i, \beta_i)] = \int_{\theta_i > 1 - \epsilon} p(\theta_i | a_i(\alpha_i, \beta_i), B) d\theta_i$$

y la localización óptima de i estaciones medidoras será aquella localización (α_i^*, β_i^*) que maximice la utilidad esperada

$$\bar{u}[a_i(\alpha_i, \beta_i)]$$

entre todas las localizaciones posibles. Como el coste de instalación es una función creciente del número de estaciones medidoras instaladas, el número apropiado i^* de estaciones a instalar será el del número mínimo de estaciones que, óptimamente distribuidas, garantizarían una utilidad esperada suficiente, esto es

$$i^* = \arg \min_{i \in \mathcal{N}} \{i; \bar{u}[a_i(\alpha_i^*, \beta_i^*)] > 1 - \epsilon\}.$$

3.2. DESCARGA DE SOJA EN BARCELONA

La presencia en el puerto de Barcelona de una instalación de descarga de soja en una localización muy próxima a la ciudad (señalada con un círculo en la Figura 10) genera la posibilidad, en determinadas condiciones meteorológicas, de un tipo de contaminación atmosférica que puede dar lugar a importantes reacciones alérgicas en un sector de la población especialmente vulnerable. Consecuentemente, las autoridades portuarias han instalado una serie de medidores que permiten emitir una alerta cuando tal tipo de contaminación es previsible.

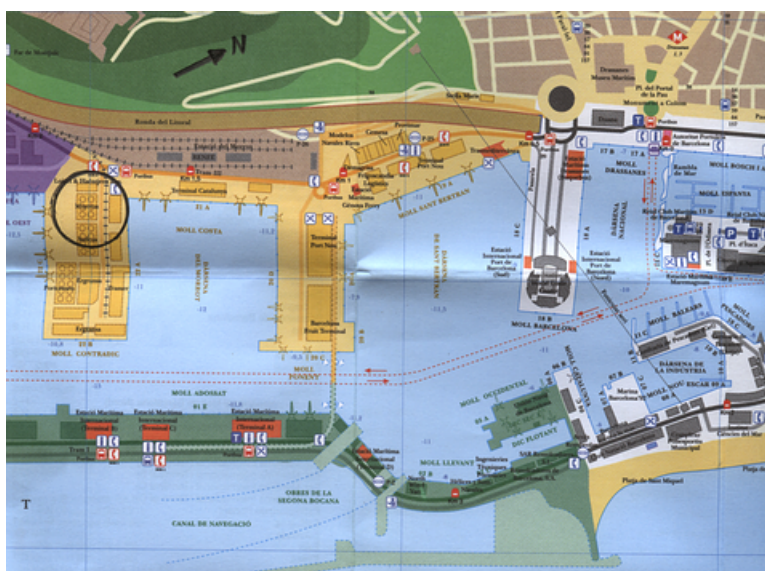


Figura 10. Punto de descarga de soja en el puerto de Barcelona.

El problema de decisión que se plantea aquí, muy frecuente en sus numerosas variantes en la gestión medioambiental de las actividades portuarias, es el de definir cual debe ser la acción a adoptar $a = a(D)$ en función de los datos D proporcionados por las estaciones de medida. Formalmente, se trata de establecer una *regla de decisión* que permita implementar, *en tiempo real*, medidas correctoras que tengan simultáneamente en cuenta de manera satisfactoria tanto los riesgos sanitarios para la población como las necesidades comerciales del puerto.

Hay que distinguir dos horizontes temporales. A corto plazo, se trata de decidir si, en una situación concreta, es necesario detener las actividades relacionadas con la soja, o es posible mantenerlas. A medio y largo plazo, se trata de decidir si pueden mejorarse las instalaciones, de forma que pueda permitirse su actividad

en un conjunto más amplio de condiciones atmosféricas, o si convendría trasladar las instalaciones a una nueva ubicación. En esta sección nos limitaremos a discutir el primer problema.

A corto plazo se trata por tanto problema de decisión con sólo dos alternativas $\mathcal{A} = \{a_0, a_1\}$: mantener las actividades (a_0) o proceder a su suspensión temporal (a_1). Las consecuencias de mantener las actividades dependen de que tenga o no tenga lugar el suceso S que se intenta evitar, esto es que se produzcan concentraciones de polvo de soja sobre la ciudad capaces de generar reacciones alérgicas en algunos de sus habitantes. En el caso de suspender temporalmente las actividades, no hay riesgo alguno de contaminación, pero existen pérdidas comerciales para la empresa concesionaria y para el puerto.

Como en todo problema de decisión, hay que cuantificar las pérdidas asociadas a cada posible combinación de decisión y resultado; en un problema con estructura tan sencilla como este, es fácil expresarlo en forma de tabla:

	S	\bar{S}
a_0	$l_{01} = l(a_0, S)$	$l_{02} = l(a_0, \bar{S})$
a_1	$l(a_1)$	

La pérdidas pueden ser medidas en una escala cualquiera. Si, por ejemplo, utilizamos una escala [0-1], la pérdida asociada a la mejor consecuencia posible (no polución, ningún comercial) sería $l_{02} = 0$, mientras que la pérdida asociada a la peor consecuencia posible (polución sobre la ciudad) sería $l_{02} = 1$, y solamente es necesario determinar, en una escala [0-1] la pérdida asociada a detener el proceso, $l(a_1)$.

Consecuentemente, la pérdida esperada de la decisión a_0 de mantener las actividades cuando los medidores proporcionan los datos D será

$$\bar{l}(a_0) = l_{01} \Pr(S | D, B) + l_{02} \Pr(\bar{S} | D, B) = \Pr(S | D, B),$$

donde $\Pr(S | D, B)$ es la probabilidad de que vayan a producirse concentraciones peligrosas de polvo de soja sobre la ciudad dada la información proporcionada por los datos D y la base de datos históricos B . Por otra parte, la pérdida esperada de detener las actividades es simplemente $\bar{l}(a_1) = l(a_1)$. Consecuentemente, la decisión óptima será detener las actividades si, y solamente si, la pérdida esperada de mantener la actividad $\bar{l}(a_0)$ es mayor que la pérdida esperada de suspenderla $\bar{l}(a_1)$. La regla de decisión debe por tanto de la forma:

$$\text{Suspender actividades siempre que } \Pr(S | D, B) > l(a_1).$$

Este sencillo análisis ha permitido identificar la forma que, *necesariamente*, debe tener la regla de decisión buscada. Es necesario plantear un modelo probabilístico que permita determinar en tiempo real, la probabilidad $\Pr(S | D, B)$ de que

vaya a producirse un episodio contaminante en función de los datos D proporcionados por las estaciones medidoras y de la base de datos disponible B . Una vez validado e implementado este modelo, las autoridades portuarias pueden limitarse a observar la evolución de $\Pr(S | D, B)$ y ordenar la suspensión de actividades siempre que $\Pr(S | D, B)$ supere un determinado nivel p^* .

Obsérvese que la especificación de p^* es una *decisión política insoslayable*. Una forma de ayudar a los responsables de tomar tal decisión es analizar sus implicaciones. Así, si quienes tengan autoridad para hacerlo fijan p^* en el valor $p^* = 0.001$ (ordenando detener las actividades si la probabilidad de un episodio contaminante es mayor de 0.001), entonces la probabilidad de que un día cualquiera se produzca un episodio de contaminación será del orden de 0.999, la de que pase una año sin problemas será del orden de $0.999^{365} \approx 0.694$ y la de que pasen 10 años sin problemas es solamente $0.999^{3650} \approx 0.026$. Inversamente, si se quiere una probabilidad q de que no haya problemas en n años, el valor umbral debe tomarse $p^* = 1 - q^{1/(365n)}$; para $q = 0.90$ y $n = 10$, resulta $p^* \approx 0.0003$ de forma que, para tener una probabilidad razonable de que no habrá problemas de contaminación si se estima en unos 10 años la vida esperada de las instalaciones actuales, el umbral de suspensión de actividades debería ser del orden de $p^* = 0.0003$.

Obsérvese que, debido a la importancia de los valores extremos de la probabilidad $\Pr(S | D, B)$ en la determinación de la solución, es *imperativo* disponer de un modelo probabilístico satisfactorio en ese dominio, lo que puede requerir un análisis muy sofisticado de la base de datos disponible B .

3.3. POSIBLE IMPACTO AMBIENTAL EN HUELVA

El puerto de Huelva, situado sobre la ría que forma el río Odiel, está separado del Atlántico por un paraje natural, la Isla de Saltés, de reconocido valor ecológico. La isla está situada inmediatamente al Oeste del muelle Juan Gonzalo, punto de descarga de graneles sólidos, y se trata de determinar si es o no necesario modificar los procedimientos vigentes de carga y descarga, y los de almacenamiento al aire libre, para garantizar a la isla (y a la población residencial de Punta Umbría, situada a continuación) una protección efectiva.

Efectivamente, se trata de determinar si los vientos del Este son capaces de arrastrar hasta la isla partículas contaminantes procedentes de las parvas acumuladas en los muelles o de las nubes producidas en los procedimientos de carga y descarga.

Se trata en este caso de un problema decisión de “contraste de hipótesis”, con dos únicas alternativas: los datos observados son compatibles con la hipótesis de que los usos actuales no tienen una incidencia apreciable sobre el parque (a_0) o, alternativamente, puede demostrarse una contaminación apreciable y, consecuentemente, hay que modificar los procedimientos en vigor (a_1).

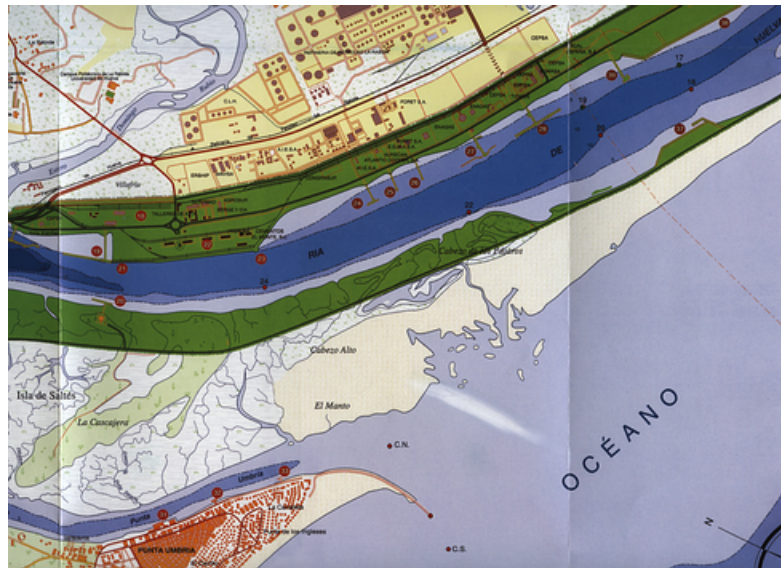


Figura 11. Ría de Huelva, Isla de Saltes y Punta Umbría.

La magnitud de interés en este problema es la cantidad θ de partículas por m^2 que se depositan sobre la isla durante un periodo de viento procedente del este, que debe ser estimada a partir de las medidas $\mathbf{x} = \{x_1, \dots, x_n\}$ obtenidas por n captadores instalados al efecto. Si $p(\mathbf{x} | \theta)$ es el modelo probabilístico que describe la relación entre las medidas realizables y la magnitud de interés, se trata de decidir si los datos observados son o no compatibles con la hipótesis de que $\theta \leq \theta_0$, donde θ_0 es el máximo valor tolerable. En este caso, la función de pérdida es una medida $\delta(\theta, \theta_0)$ de la discrepancia entre el modelo verdadero $p(\mathbf{x} | \theta)$ y el la subfamilia $\{p(\mathbf{x} | \theta), \theta \leq \theta_0\}$, que se mide en unidades de información, y la hipótesis de incidencia no apreciable (a_0) debe ser desestimada si, y solamente si, el valor esperado de esa discrepancia es suficientemente grande, *i.e.*, si

$$\int_{\Theta} \delta(\theta, \theta_0) p(\theta | \mathbf{x}) d\theta > d^*$$

donde $p(\theta | \mathbf{x})$ es la distribución de probabilidad que describe la información proporcionada sobre los niveles de polución θ por las muestras tomadas \mathbf{x} y donde, convencionalmente, se toma $d^* = 5$ unidades de información (equivalente a tres desviaciones típicas en el lenguaje convencional en ingeniería).

En un problema como el que nos ocupa es crucial realizar un análisis preciso de las muestras obtenidas \mathbf{x} para garantizar que provienen de la fuente considerada

objeto del problema. Por ejemplo, en el caso descrito, sería necesario asegurar que las muestras recogidas provienen de los muelles del puerto y no de las industrias instaladas en el polígono industrial situado al Este del muelle Juan Gonzalo.

3.4. DESCARGA DE GRANELES SÓLIDOS EN SANTANDER

En condiciones de viento adversas, la polución atmosférica sobre ciudades portuarias como consecuencia de partículas en suspensión levantadas desde parvas situadas al aire libre, o como consecuencia de operaciones de carga o descarga de graneles sólidos, es un problema medioambiental al que frecuentemente deben enfrentarse las autoridades portuarias.

En el caso particular del puerto de Santander, las parvas de carbón situadas sobre el espigón central de Raos así como las operaciones de carga y descarga en ese muelle (frecuente, pero no exclusivamente de carbón) dan lugar, con vientos fuertes del sur, a importantes episodios de contaminación atmosférica sobre el barrio de la ciudad situado al norte de ese espigón (cuadrante superior izquierda de la Figura 12).



Figura 12. Puerto de Santander.

El problema planteado por las parvas acumuladas es técnicamente muy diferente del que plantean las operaciones de carga y descarga. En el primer caso es necesario plantearse opciones tales como la instalación o mejora de sistemas de riego

sobre de las parvas o su cambio de ubicación. Este tipo de problema de decisión será analizado en la sección siguiente al describir un problema de estas características en el puerto de Tarragona.

En su versión más sencilla, el problema planteado por la contaminación atmosférica que las operaciones previstas de carga y descarga de graneles sólidos pueden llegar a producir admite tres acciones alternativas, $\mathcal{A} = \{a_0, a_1, a_2\}$: autorizar la operación en la forma habitual (a_0), autorizarla en forma restringida (a_1) (operando solamente con tolvas especiales, disminuyendo la altura de volcado,...), o posponer la operación hasta que mejore las situación meteorológica prevista (a_3) (Figura 13).

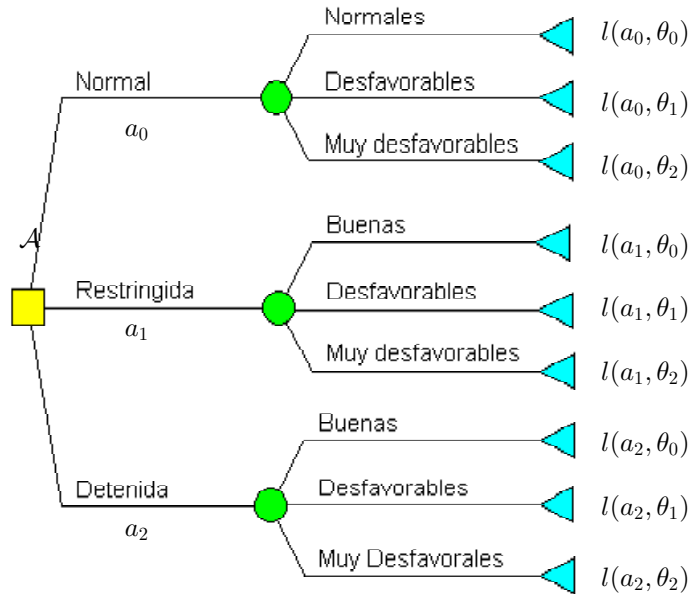


Figura 13. Árbol de decisión básico en operaciones de carga y descarga

Las consecuencias de cualquiera de estas acciones dependen de *forma probabilística* de las condiciones meteorológicas (fundamentalmente de la dirección e intensidad del viento y, en menor medida, de la humedad relativa). En una primera aproximación, tales condiciones pueden ser clasificadas en tres grandes grupos: buenas (θ_0), desfavorables (θ_1) y muy desfavorables o malas (θ_2). Para cada combinación de acción posible a_i y condiciones meteorológicas θ_j es necesario evaluar, en una escala adecuada, las pérdidas $l_{ij} = l(a_i, \theta_j)$, que se generarían, dando lugar a una tabla de la forma:

	θ_0	θ_1	θ_2
a_0	$l_{00} = l(a_0, \theta_0)$	$l_{01} = l(a_0, \theta_1)$	$l_{02} = l(a_0, \theta_2)$
a_1	$l_{10} = l(a_1, \theta_0)$	$l_{11} = l(a_1, \theta_1)$	$l_{12} = l(a_1, \theta_2)$
a_2	$l_{20} = l(a_2, \theta_0)$	$l_{21} = l(a_2, \theta_1)$	$l_{22} = l(a_2, \theta_2)$

Además, desde la perspectiva que nos ocupa, la consecuencia básica de la acción a_i si se presentan las condiciones θ_j es una *variable aleatoria*, la cantidad ω_{ij} de partículas recogidas por m^2 , cuya distribución, obviamente, depende de a_i y de θ_j . Formalmente,

$$l_{ij} = l(a_i, \theta_j) = \int_0^\infty l(c_i, \omega_{ij}) p(\omega_{ij} | a_i, \theta_j, B) d\omega_{ij},$$

de forma que la pérdida sufrida si se toma la acción a_i en las condiciones θ_j es una función del coste c_i de implementar a_i y de la distribución de probabilidad $p(\omega_{ij} | a_i, \theta_j, B)$ del nivel ω_{ij} de polución que puede esperarse en esas circunstancias en base a la información proporcionada por una base de datos histórica B . Obsérvese que si se detienen las operaciones, las actividades del puerto no pueden contribuir a la polución de la ciudad; consecuentemente, w_{20} , w_{21} y w_{22} describen, respectivamente, los niveles de polución del sector de la ciudad estudiado (en las distintas condiciones meteorológicas estudiadas) cuando *no* existen emisiones procedentes del puerto.

La función $l(c, \omega)$ debe describir en una escala conveniente (por ejemplo en una escala [0-1]) la pérdida asociada a una situación en la que con un gasto c se detecta una nivel de polución ω . Determinar esa función es una tarea difícil (pero inevitable) que debe ser el resultado de algún tipo de acuerdo entre todas las partes implicadas

Finalmente, una vez determinada la función de pérdida es necesario analizar la base de datos histórica B para calcular las probabilidades $\Pr(\theta_j | M, B)$ asociadas a cada uno de los tipos de condiciones meteorológicas consideradas, $\{\theta_0, \theta_1, \theta_2\}$, dados los datos meteorológicos M en el momento de tomar la decisión.

La utilidad esperada de cada una de las alternativas considerada será entonces

$$\bar{l}(a_i | M, B) = \sum_{j=0}^2 l(a_i, \theta_j) \Pr(\theta_j | M, B)$$

donde los $l(a_i, \theta_j)$ son los valores esperados de las pérdidas $l(c_i, \omega_{ij})$ asociadas a los niveles de polución ω_{ij} . La decisión óptima a^* en condiciones meteorológicas M es entonces aquella que minimiza la pérdida esperada,

$$a^* = a^*(M, B) = \arg \min \bar{l}(a_i | M, B).$$

3.5. DESCARGA DE CARBÓN EN TARRAGONA

La descarga de carbón en el puerto de Tarragona tiene lugar a cielo abierto, en el muelle de Cataluña, situado en el extremo sur de puerto (Figura 14). El dique rompeolas que separa el muelle del mar abierto soporta una estrecha carretera, a unos metros sobre el nivel del mar, que permite el acceso al faro de la Banya, y que constituye uno de los paseos deportivos de la ciudad. Además sobre esa carretera, inmediatamente detrás del muelle, se sitúa un centro de buceo, debido a que esa zona de mar presenta características adecuadas para esa práctica deportiva. Cuando sopla viento del norte, las partículas de carbón de las parvas acumuladas en el muelle barren la carretera, haciendo muy desagradable pasear sobre ella, y contaminan de carbón al zona colindante de buceo, deteriorándola de forma progresiva.

En este caso que trata de un problema continuo, no limitado a los momentos de descarga de buques, y cuya solución exigiría obras de envergadura. Frente a la situación actual a_0 , un primer análisis pone de manifiesto tres alternativas:

- a_1 : mejorar el sistema de riego de las parvas, tal vez con aditivos apropiados, para limitar el desplazamiento de las partículas situadas en su superficie.
- a_2 : Diseñar y construir entre el muelle y la carretera una barrera física que no sea superada por las partículas de carbón con vientos fuertes del norte.
- a_3 : Buscar una ubicación alternativa para la descarga del carbón.

Las consecuencias finales que se derivan de cualquiera de estas acciones alternativas son función de dos variables: la cantidad c de carbón que se depositaría al año por m^2 de la zona de buceo y el coste e en euros de implementar esa medida. La primera tarea necesaria es por tanto cuantificar, en una escala apropiada, las preferencias asociadas a cada solución $\{c, e\}$ (carbón, euros) posible. Una forma eficiente de hacerlo sería ajustar una función de pérdida $l[c, e]$ que combine las pérdidas asociadas a un depósito de carbón c con las de un coste e en una escala $[0 - 100]$. La solución ideal, obviamente inexistente, sería polución $c = 0$ con coste $e = 0$, de forma que $l[0, 0] = 0$. Cualquier solución con un coste $e > e_1$ no asumible tendría pérdida máxima, de forma que si $e > e_1$, $l[e, c] = 100$. Análogamente, cualquier solución con una polución $c > c_1$ excesiva (por ilegal o por socialmente inaceptable) tendría pérdida máxima, de forma que si $c > c_1$, $l[e, c] = 100$. Para cualquier otra solución $\{c, e\}$ habría que obtener un valor numérico $l(c, e)$ entre 0 y 100, fruto de un pacto social, de algún tipo de acuerdo, entre todas las instituciones implicadas. Este acuerdo podría ser facilitado por el diseño y análisis de una encuesta que determinase la opinión al respecto de los habitantes de la ciudad. La determinación de la función $l[e, c]$ que describe las preferencias no es sencilla, pero es importante subrayar que es una condición *necesaria* para una solución *racional* al problema planteado.

Para cada una de las alternativas consideradas es necesario precisar la distribución de probabilidad de los niveles de polución c y los costes e a que daría lugar.

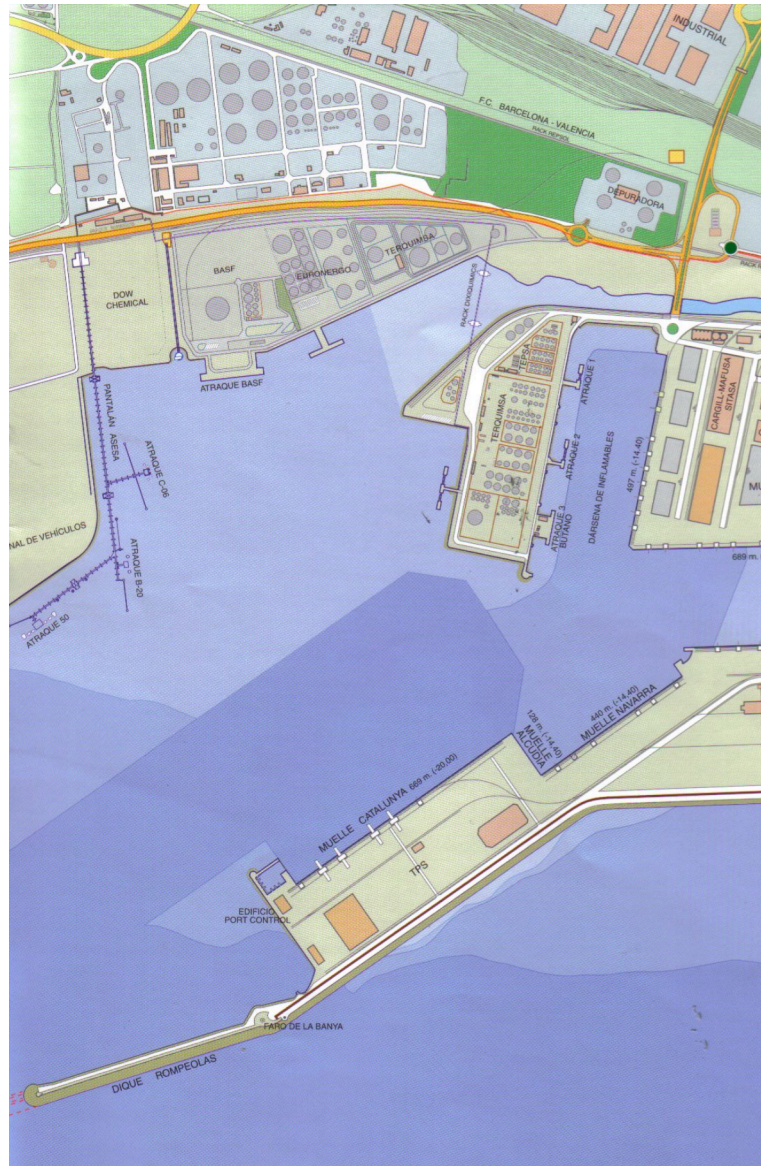


Figura 14. Muelle de Catalunya, Puerto de Tarragona.

Por otra parte, sería necesario disponer de una base de datos meteorológicos M que permitan estimar la distribución de la intensidad y dirección del viento sobre el muelle *a lo largo de todo un año* (para integrar previsible efectos estacionales).

En el caso a_1 de los sistemas de riego, habría que estudiar, posiblemente mediante experimentación *in situ*, el comportamiento de cada posible sistema S_j para, en función de los datos meteorológicos M , determinar la distribución de probabilidad $p(c | S_j, M)$ sobre la polución c a que daría lugar. Si el coste de instalación y mantenimiento del sistema S_j es $e(S_j)$, la pérdida esperada de la solución a_1 con el sistema S_j sería

$$\bar{l}(S_j) = \int_C l[c, e(S_j)] p(c | S_j, M) dc$$

y el mejor sistema el que minimice $\bar{l}(S_j)$ entre todos los posibles, de forma que la mejor solución entre las de riego tiene una pérdida esperada

$$\bar{l}(a_1) = \min_j \bar{l}(S_j)$$

De forma totalmente análoga, en el caso a_2 de las barreras físicas, habría que estudiar, posiblemente mediante métodos de simulación, el comportamiento de cada posible barrera B_j para, en función de los datos meteorológicos M , determinar la distribución de probabilidad $p(c | B_j, M)$ sobre la polución c a que daría lugar. Si el coste de construcción de la barrera B_j es $e(B_j)$, la pérdida esperada de la solución a_1 con el diseño B_j sería

$$\bar{l}(B_j) = \int_C l[c, e(B_j)] p(c | B_j, M) dc$$

y el mejor diseño el que minimice $\bar{l}(B_j)$ entre todos los posibles, de forma que la mejor solución entre las de barrera tiene una pérdida esperada $\bar{l}(a_2) = \min_j \bar{l}(B_j)$.

En el caso a_3 de traslado de las instalaciones, la polución sobre la zona de buceo sería obviamente nula. Si el coste de traslado a una ubicación U_j es $c(U_j)$, la pérdida esperada de la solución a_3 con la ubicación U_j sería $\bar{l}(a_3) = \min_j l[0, c(U_j)]$.

Finalmente, la pérdida esperada de la alternativa “nula” a_0 , esto es de mantener el statu quo, que no tiene coste monetario directo, sería

$$\bar{l}(a_0) = \int_C l[c, 0] p(c | M) dc,$$

donde $p(c | M)$ describe la distribución de probabilidad de la cantidad c de carbón por m^2 que se depositaría en un año sobre la zona de buceo si no se efectuase acción correctora alguna.

La decisión óptima a^* , será aquella que minimice la pérdida esperada,

$$a^* = \arg \min \bar{l}(a_i), \quad i = 0, 1, 2, 3.$$

El resultado dependerá obviamente de la función de pérdida utilizada. Para disponer de resultados generales, se puede plantear una familia paramétrica de funciones de pérdida suficientemente amplia, y proporcionar una tabla que recoja la solución óptima en función de los valores de los parámetros utilizados.

Apéndice

BAYESIAN STATISTICS

Este apéndice contiene una versión abreviada y actualizada del artículo preparado por el autor de este informe para la enciclopedia científica de la UNESCO:

Bernardo, J. M. (2003). Bayesian Statistics. *Encyclopedia of Life Support Systems (EOLSS). Probability and Statistics* (R. Viertl, ed). London, UK: UNESCO.

Contiene una exposición más detallada, y de mayor nivel matemático, de los conceptos introducidos en el capítulo 2 de este informe.

Bayesian Statistics

José M. Bernardo

Summary

Statistics is the study of uncertainty. The field of statistics is based on two major paradigms: conventional and Bayesian. Bayesian methods provide a *complete* paradigm for both statistical inference and decision making under uncertainty. Bayesian methods may be derived from an axiomatic system and provide a *coherent* methodology which makes it possible to incorporate relevant initial information, and which alleviates many of the difficulties faced by conventional statistical methods. The Bayesian paradigm is based on an interpretation of probability as a *conditional measure of uncertainty* which closely matches the sense of the word 'probability' in ordinary language. Statistical inference about a quantity of interest is described as the modification of the uncertainty about its value in the light of evidence, and Bayes' theorem specifies how this modification should be made. Bayesian methods may be

applied to highly structured complex problems, which have been often intractable by traditional statistical methods. The special situation, often met in scientific reporting and public decision making, where the only acceptable information is that which may be deduced from available documented data, is addressed as an important particular case.

1. Introduction

Scientific experimental or observational results generally consist of (possibly many) sets of data of the general form $D = \{x_1, \dots, x_n\}$, where the x_i 's are somewhat "homogeneous" (possibly multidimensional) observations x_i . Statistical methods are then typically used to derive conclusions on both the nature of the process which has produced those observations, and on the expected behaviour at future instances of the same process. A central element of *any* statistical analysis is the specification of a *probability model* which is assumed to describe the mechanism which has generated the observed data D as a function of a (possibly multidimensional) parameter (vector) $\omega \in \Omega$, sometimes referred to as the *state of nature*, about whose value only limited information (if any) is available. All derived statistical conclusions are obviously conditional on the assumed probability model.

Unlike most other branches of mathematics, conventional methods of statistical inference suffer from the lack of an axiomatic basis; as a consequence, their proposed desiderata are often mutually incompatible, and the analysis of the same data may well lead to incompatible results when different, apparently intuitive procedures are tried (see the 1970's monographs by Lindley and by Jaynes for many instructive examples). In marked contrast, the Bayesian approach to statistical inference is firmly based on axiomatic foundations which provide a unifying logical structure, and guarantee the mutual consistency of the methods proposed. Bayesian methods constitute a *complete* paradigm to statistical inference, a scientific revolution in Kuhn's sense.

Bayesian statistics only require the *mathematics* of probability theory and the *interpretation* of probability which most closely corresponds to the standard use of this word in everyday language: it is no accident that some of the more important seminal books on Bayesian statistics, such as the works of de Laplace, de Finetti or Jeffreys, are actually entitled "Probability Theory". The practical consequences of adopting the Bayesian paradigm are far reaching. Indeed, Bayesian methods (i) reduce statistical inference to problems in probability theory, thereby minimizing the need for completely new concepts, and (ii) serve to discriminate among conventional statistical techniques, by either providing a logical justification to some (and making explicit the conditions under which they are valid), or proving the logical inconsistency of others.

The main consequence of these foundations is the mathematical *need* to describe by means of probability distributions all uncertainties present in the problem. In particular, unknown parameters in probability models *must* have a joint probability distribution which describes the available information about their values; this is often regarded as *the* characteristic element of a Bayesian approach. Notice that (in sharp contrast to conventional statistics) *parameters are treated as random variables* within the Bayesian paradigm. This is not a description of their variability (parameters are typically *fixed unknown* quantities) but a description of the *uncertainty* about their true values.

An important particular case arises when either no relevant prior information is readily available, or that information is subjective and an “objective” analysis is desired, one that is exclusively based on accepted model assumptions and well-documented data. This is addressed by *reference analysis* which uses information-theoretic concepts to derive appropriate reference posterior distributions, defined to encapsulate inferential conclusions on the quantities of interest solely based on the supposed model and the observed data.

In this article it is assumed that probability distributions may be described through their probability density functions, and no distinction is made between a random quantity and the particular values that it may take. Bold italic roman fonts are used for *observable* random vectors (typically data) and bold italic greek fonts are used for unobservable random vectors (typically parameters); lower case is used for variables and upper case for their dominion sets. Moreover, the standard mathematical convention of referring to *functions*, say f and g of $\mathbf{x} \in X$, respectively by $f(\mathbf{x})$ and $g(\mathbf{x})$, will be used throughout. Thus, $p(\boldsymbol{\theta} | C)$ and $p(\mathbf{x} | C)$ respectively represent general *probability densities* of the random vectors $\boldsymbol{\theta} \in \Theta$ and $\mathbf{x} \in X$ under conditions C , so that $p(\boldsymbol{\theta} | C) \geq 0$, $\int_{\Theta} p(\boldsymbol{\theta} | C) d\boldsymbol{\theta} = 1$, and $p(\mathbf{x} | C) \geq 0$, $\int_X p(\mathbf{x} | C) d\mathbf{x} = 1$. This admittedly imprecise notation will greatly simplify the exposition. If the random vectors are discrete, these functions naturally become probability mass functions, and integrals over their values become sums.

Density functions of specific distributions are denoted by appropriate names. Thus, if x is a random quantity with a normal distribution of mean μ and standard deviation σ , its probability density function will be denoted $N(x | \mu, \sigma)$. Table 1 contains definitions of other distributions used in this article.

Bayesian methods make frequent use of the the concept of logarithmic divergence, a very general measure of the goodness of the approximation of a probability density $p(\mathbf{x})$ by another density $\hat{p}(\mathbf{x})$. The Kullback-Leibler, or *logarithmic divergence* of a probability density $\hat{p}(\mathbf{x})$ of the random vector $\mathbf{x} \in X$ from its true probability density $p(\mathbf{x})$, is defined as

$$k\{\hat{p}(\mathbf{x}) | p(\mathbf{x})\} = \int_X p(\mathbf{x}) \log\{p(\mathbf{x})/\hat{p}(\mathbf{x})\} d\mathbf{x}.$$

Table 1. Notation for common probability distributions

Name	Probability Mass (or Density) Function
Beta	$\text{Be}(x \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}, x \in (0, 1)$
Binomial	$\text{Bi}(x n, \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}, x \in \{0, \dots, n\}$
Exponential	$\text{Ex}(x \theta) = \theta e^{-\theta x}, x > 0$
ExpGamma	$\text{Eg}(x \alpha, \beta) = \frac{\alpha\beta^\alpha}{(x+\beta)^{\alpha+1}}, x > 0$
Gamma	$\text{Ga}(x \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, x > 0$
NegBinomial	$\text{Nb}(x r, \theta) = \theta^r \binom{r+x-1}{r-1} (1-\theta)^x, x \in \{0, 1, \dots\}$
Normal	$\text{N}_k(\mathbf{x} \boldsymbol{\mu}, \Sigma) = \frac{ \Sigma ^{-1/2}}{(2\pi)^{k/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right], \mathbf{x} \in \mathfrak{R}^k$
Poisson	$\text{Pn}(x \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, x \in \{0, 1, \dots\}$
Student	$\text{St}(x \mu, \sigma, \alpha) = \frac{\Gamma(\frac{\alpha+1}{2})}{\Gamma(\frac{\alpha}{2})} \frac{1}{\sigma\sqrt{\alpha\pi}} \left[1 + \frac{1}{\alpha} \left(\frac{x-\mu}{\sigma}\right)^2\right]^{-(\alpha+1)/2}, x \in \mathfrak{R}$

It may be shown that (i) the logarithmic divergence is non-negative (and it is zero if, and only if, $\hat{p}(\mathbf{x}) = p(\mathbf{x})$ almost everywhere), and (ii) that $\delta\{\hat{p}(\mathbf{x}) | p(\mathbf{x})\}$ is invariant under one-to-one transformations of \mathbf{x} .

This article contains a brief summary of the mathematical foundations of Bayesian statistical methods (Section 2), an overview of the paradigm (Section 3), a description of useful inference summaries, including estimation and hypothesis testing (Section 4), an explicit discussion of objective Bayesian methods (Section 5), the detailed analysis of a simplified case study (Section 6), and a final discussion which includes pointers to further issues not addressed here (Section 7).

2. Foundations

A central element of the Bayesian paradigm is the use of probability distributions to describe all relevant unknown quantities, interpreting the probability of an event as a conditional measure of uncertainty, on a $[0, 1]$ scale, about the occurrence of the event in some specific conditions. The limiting extreme values 0 and 1, which are typically inaccessible in applications, respectively describe impossibility and certainty of the

occurrence of the event. This interpretation of probability includes and extends all other probability interpretations. There are two independent arguments which prove the mathematical inevitability of the use of probability distributions to describe uncertainties; these are summarized later in this section.

2.1. Probability as a Measure of Conditional Uncertainty

Bayesian statistics uses the word *probability* in precisely the same sense in which this word is used in everyday language, as a *conditional measure of uncertainty* associated with the occurrence of a particular event, given the available information and the accepted assumptions. Thus, $\Pr(E | C)$ is a measure of (presumably rational) belief in the occurrence of the *event* E under *conditions* C . It is important to stress that probability is *always* a function of two arguments, the event E whose uncertainty is being measured, and the conditions C under which the measurement takes place; “absolute” probabilities do not exist. In typical applications, one is interested in the probability of some event E given the available *data* D , the set of *assumptions* A which one is prepared to make about the mechanism which has generated the data, and the relevant contextual *knowledge* K which might be available. Thus, $\Pr(E | D, A, K)$ is to be interpreted as a measure of (presumably rational) belief in the occurrence of the *event* E , given data D , assumptions A and any other available knowledge K , as a measure of how “likely” is the occurrence of E in these conditions. Sometimes, but certainly not always, the probability of an event under given conditions may be associated with the relative frequency of “similar” events in “similar” conditions. The following examples are intended to illustrate the use of probability as a conditional measure of uncertainty.

Probabilistic diagnosis. A human population is known to contain 0.2% of people infected by a particular virus. A person, *randomly selected* from that population, is subject to a test which is from laboratory data known to yield positive results in 98% of infected people and in 1% of non-infected, so that, if V denotes the event that a person carries the virus and $+$ denotes a positive result, $\Pr(+ | V) = 0.98$ and $\Pr(+ | \bar{V}) = 0.01$. Suppose that the result of the test turns out to be positive. Clearly, one is then interested in $\Pr(V | +, A, K)$, the *probability* that the person carries the virus, given the positive result, the assumptions A about the probability mechanism generating the test results, and the available knowledge K of the prevalence of the infection in the population under study (described here by $\Pr(V | K) = 0.002$). An elementary exercise in probability algebra, which involves Bayes’ theorem in its simplest form (see Section 3), yields $\Pr(V | +, A, K) = 0.164$. Notice that the four probabilities involved in the problem have *the same interpretation*: they are all conditional measures of uncertainty. Besides, $\Pr(V | +, A, K)$ is *both* a measure of the uncertainty associated with the event that the particular person who tested positive is actually infected, *and* an *estimate* of the proportion of people in that

population (about 16.4%) that would eventually prove to be infected among those which yielded a positive test.

Estimation of a proportion. A survey is conducted to estimate the proportion θ of individuals in a population who share a given property. A random sample of n elements is analyzed, r of which are found to possess that property. One is then typically interested in using the results from the sample to establish regions of $[0, 1]$ where the unknown value of θ may plausibly be expected to lie; this information is provided by *probabilities* of the form $\Pr(a < \theta < b \mid r, n, A, K)$, a conditional measure of the uncertainty about the event that θ belongs to (a, b) given the information provided by the data (r, n) , the assumptions A made on the behaviour of the mechanism which has generated the data (a random sample of n Bernoulli trials), and any relevant knowledge K on the values of θ which might be available. For example, after a political survey in which 720 citizens out of a random sample of 1500 have declared their support to a particular political measure, one may conclude that $\Pr(\theta < 0.5 \mid 720, 1500, A, K) = 0.933$, indicating a probability of about 93% that a referendum of that issue would be lost. Similarly, after a screening test for an infection where 100 people have been tested, none of which has turned out to be infected, one may conclude that $\Pr(\theta < 0.01 \mid 0, 100, A, K) = 0.844$, or a probability of about 84% that the proportion of infected people is smaller than 1%.

Measurement of a physical constant. A team of scientists, intending to establish the unknown value of a physical constant μ , obtain data $D = \{x_1, \dots, x_n\}$ which are considered to be measurements of μ subject to error. The probabilities of interest are then typically of the form $\Pr(a < \mu < b \mid x_1, \dots, x_n, A, K)$, the *probability* that the unknown value of μ (fixed in nature, but unknown to the scientists) lies within an interval (a, b) given the information provided by the data D , the assumptions A made on the behaviour of the measurement mechanism, and whatever knowledge K might be available on the value of the constant μ . Again, those probabilities are conditional measures of uncertainty which describe the (necessarily probabilistic) conclusions of the scientists on the true value of μ , given available information and accepted assumptions. For example, after a classroom experiment to measure the gravitational field with a pendulum, a student may report (in m/sec^2) something like $\Pr(9.788 < g < 9.829 \mid D, A, K) = 0.95$, meaning that, under accepted knowledge K and assumptions A , the *observed* data D indicate that the true value of g lies within 9.788 and 9.829 with probability 0.95, a conditional uncertainty measure on a $[0, 1]$ scale. This is naturally compatible with the fact that the value of the gravitational field at the laboratory may well be known with high precision from available literature or from precise previous experiments, but the student may have been instructed *not* to use that information as part of the accepted knowledge K . Under some conditions, it is also true that if the same *procedure* were actually used by many other students with similarly obtained data sets, their reported intervals would

actually cover the true value of g in approximately 95% of the cases, thus providing some form of *calibration* for the student's probability statement (see Section 5.2).

Prediction. An experiment is made to count the number r of times that an event E takes place in each of n replications of a well defined situation; it is observed that E does take place r_i times in replication i , and it is desired to forecast the number of times r that E will take place in a future, similar situation. This is a *prediction* problem on the value of an *observable* (discrete) quantity r , given the information provided by data D , accepted assumptions A on the probability mechanism which generates the r_i 's, and any relevant available knowledge K . Hence, simply the computation of the probabilities $\{\Pr(r | r_1, \dots, r_n, A, K)\}$, for $r = 0, 1, \dots$, is required. For example, the quality assurance engineer of a firm which produces automobile restraint systems may report something like $\Pr(r = 0 | r_1 = \dots = r_{10} = 0, A, K) = 0.953$, after observing that the entire production of airbags in each of $n = 10$ consecutive months has yielded no complaints from their clients. This should be regarded as a measure, on a $[0, 1]$ scale, of the conditional uncertainty, given observed data, accepted assumptions and contextual knowledge, associated with the event that no airbag complaint will come from next month's production and, if conditions remain constant, this is also an estimate of the proportion of months expected to share this desirable property.

A similar problem may naturally be posed with continuous observables. For instance, after measuring some continuous magnitude in each of n randomly chosen elements within a population, it may be desired to forecast the proportion of items in the whole population whose magnitude satisfies some precise specifications. As an example, after measuring the breaking strengths $\{x_1, \dots, x_{10}\}$ of 10 randomly chosen safety belt webbings to verify whether or not they satisfy the requirement of remaining above 26 kN, the quality assurance engineer may report something like $\Pr(x > 26 | x_1, \dots, x_{10}, A, K) = 0.9987$. This should be regarded as a measure, on a $[0, 1]$ scale, of the conditional uncertainty (given observed data, accepted assumptions and contextual knowledge) associated with the event that a randomly chosen safety belt webbing will support no less than 26 kN. If production conditions remain constant, it will also be an estimate of the proportion of safety belts which will conform to this particular specification.

Often, additional information of future observations is provided by related covariates. For instance, after observing the outputs $\{y_1, \dots, y_n\}$ which correspond to a sequence $\{x_1, \dots, x_n\}$ of different production conditions, it may be desired to forecast the output y which would correspond to a particular set x of production conditions. For instance, the viscosity of commercial condensed milk is required to be within specified values a and b ; after measuring the viscosities $\{y_1, \dots, y_n\}$ which correspond to samples of condensed milk produced under different physical conditions $\{x_1, \dots, x_n\}$, production engineers will require probabilities of the form

$\Pr(a < y < b \mid \mathbf{x}, (y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n), A, K)$. This is a conditional measure of the uncertainty (always given observed data, accepted assumptions and contextual knowledge) associated with the event that condensed milk produced under conditions \mathbf{x} will actually satisfy the required viscosity specifications.

2.2. Statistical Inference and Decision Theory

Decision theory not only provides a precise methodology to deal with decision problems under uncertainty, but its solid axiomatic basis also provides a powerful reinforcement to the logical force of the Bayesian approach. We now summarize the basic argument.

A decision problem exists whenever there are two or more possible courses of action; let \mathcal{A} be the class of possible actions. Moreover, for each $a \in \mathcal{A}$, let Θ_a be the set of *relevant events* which may affect the result of choosing a , and let $c(a, \boldsymbol{\theta}) \in \mathcal{C}_a$, $\boldsymbol{\theta} \in \Theta_a$, be the *consequence* of having chosen action a when event $\boldsymbol{\theta}$ takes place. The class of pairs $\{(\Theta_a, \mathcal{C}_a), a \in \mathcal{A}\}$ describes the *structure* of the decision problem. Without loss of generality, it may be assumed that the possible actions are mutually exclusive, for otherwise one would work with the appropriate Cartesian product.

Different sets of principles have been proposed to capture a minimum collection of logical rules that could sensibly be required for “rational” decision-making. These all consist of axioms with a strong intuitive appeal; examples include the *transitivity* of preferences (if $a_1 > a_2$ given C , and $a_2 > a_3$ given C , then $a_1 > a_3$ given C), and the *sure-thing principle* (if $a_1 > a_2$ given C and E , and $a_1 > a_2$ given C and \bar{E} , then $a_1 > a_2$ given C). Notice that these rules are not intended as a description of actual human decision-making, but as a *normative* set of principles to be followed by someone who aspires to achieve coherent decision-making.

There are naturally different options for the set of acceptable principles, but all of them lead basically to the same conclusions, namely:

(i) Preferences among consequences should be measured with a real-valued bounded *utility* function $U(c) = U(a, \boldsymbol{\theta})$ which specifies, on some numerical scale, their desirability.

(ii) The uncertainty of relevant events should be measured with a set of *probability* distributions $\{p(\boldsymbol{\theta} \mid C, a), \boldsymbol{\theta} \in \Theta_a, a \in \mathcal{A}\}$ describing their plausibility given the conditions C under which the decision must be taken.

(iii) The desirability of the available actions is measured by their corresponding *expected utility*

$$\bar{U}(a \mid C) = \int_{\Theta_a} U(a, \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid C, a) d\boldsymbol{\theta}, \quad a \in \mathcal{A}. \quad (1)$$

It is often convenient to work in terms of the non-negative *loss* function defined by

$$L(a, \boldsymbol{\theta}) = \sup_{a \in \mathcal{A}} \{U(a, \boldsymbol{\theta})\} - U(a, \boldsymbol{\theta}), \quad (2)$$

which directly measures, as a function of $\boldsymbol{\theta}$, the “penalty” for choosing a wrong action. The relative undesirability of available actions $a \in \mathcal{A}$ is then measured by their *expected loss*

$$\bar{L}(a | C) = \int_{\Theta_a} L(a, \boldsymbol{\theta}) p(\boldsymbol{\theta} | C, a) d\boldsymbol{\theta}, \quad a \in \mathcal{A}. \quad (3)$$

Notice that, in particular, the argument described above establishes the need to quantify the uncertainty about all relevant unknown quantities (the actual values of the $\boldsymbol{\theta}$'s), and specifies that this quantification *must* have the mathematical structure of probability distributions. These probabilities are conditional on the circumstances C under which the decision is to be taken, which typically, but not necessarily, include the results D of some relevant experimental or observational data.

It has been argued that the development described above (which is not questioned when decisions have to be made) does not apply to problems of statistical inference, where no specific decision making is envisaged. However, there are two powerful counterarguments to this. Indeed, (i) a problem of statistical inference is typically considered worth analysing because it *may* eventually help make sensible decisions (as Ramsey put it in the 1930's, a lump of arsenic is poisonous because it *may* kill someone, not because it has actually killed someone), and (ii) it has been shown (by Bernardo in the 1970's) that statistical inference on $\boldsymbol{\theta}$ actually *has* the mathematical structure of a decision problem, where the class of alternatives is the functional space

$$\mathcal{A} = \left\{ p(\boldsymbol{\theta} | D); \quad p(\boldsymbol{\theta} | D) > 0, \int_{\Theta} p(\boldsymbol{\theta} | D) d\boldsymbol{\theta} = 1 \right\} \quad (4)$$

of the conditional probability distributions of $\boldsymbol{\theta}$ given the data, and the utility function is a measure of the amount of information about $\boldsymbol{\theta}$ which the data may be expected to provide.

2.3. Exchangeability and Representation Theorem

Available data often take the form of a set $\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_n\}$ of “homogeneous” (possibly multidimensional) observations, in the precise sense that only their *values* matter and not the *order* in which they appear. Formally, this is captured by the notion of *exchangeability*. The set of random vectors $\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_n\}$ is exchangeable if their joint distribution is invariant under permutations. An infinite sequence $\{\boldsymbol{x}_j\}$ of random vectors is exchangeable if all its finite subsequences are exchangeable. Notice that, in particular, any random sample from any model is exchangeable

in this sense. The concept of exchangeability, introduced by de Finetti in the 1930's, is central to modern statistical thinking. Indeed, the general *representation theorem* implies that if a set of observations is assumed to be a subset of an exchangeable sequence, then it constitutes a *random sample* from some probability model $\{p(\mathbf{x}|\boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$, $\mathbf{x} \in X$, described in terms of (labelled by) some *parameter vector* $\boldsymbol{\omega}$; furthermore this parameter $\boldsymbol{\omega}$ is *defined* as the limit (as $n \rightarrow \infty$) of some function of the observations. Available information about the value of $\boldsymbol{\omega}$ in prevailing conditions C is *necessarily* described by *some* probability distribution $p(\boldsymbol{\omega} | C)$.

For example, in the case of a sequence $\{x_1, x_2, \dots\}$ of dichotomous exchangeable random quantities $x_j \in \{0, 1\}$, de Finetti's representation theorem establishes that the joint distribution of (x_1, \dots, x_n) has the *integral representation*

$$p(x_1, \dots, x_n | C) = \int_0^1 \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} p(\theta | C) d\theta \quad (5)$$

where $\theta = \lim_{n \rightarrow \infty} (r/n)$ and $r = \sum x_j$ is the number of positive trials. This is nothing but the joint distribution of a set of (conditionally) independent Bernoulli trials with parameter θ , over which some probability distribution $p(\theta | C)$ is therefore proven to exist. More generally, for sequences of arbitrary random quantities $\{x_1, x_2, \dots\}$, exchangeability leads to integral representations of the form

$$p(x_1, \dots, x_n | C) = \int_{\Omega} \prod_{i=1}^n p(x_i | \boldsymbol{\omega}) p(\boldsymbol{\omega} | C) d\boldsymbol{\omega}, \quad (6)$$

where $\{p(\mathbf{x} | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$ denotes some probability *model*, $\boldsymbol{\omega}$ is the limit as $n \rightarrow \infty$ of some function $f(x_1, \dots, x_n)$ of the observations, and $p(\boldsymbol{\omega} | C)$ is some probability distribution over Ω . This formulation includes "nonparametric" (distribution free) modelling, where $\boldsymbol{\omega}$ may index, for instance, all continuous probability distributions on X . Notice that $p(\boldsymbol{\omega} | C)$ does *not* describe a possible variability of $\boldsymbol{\omega}$ (since $\boldsymbol{\omega}$ will typically be a fixed *unknown* vector), but a description on the uncertainty associated with its actual value.

Under appropriate conditioning, exchangeability is a very general assumption, a powerful extension of the traditional concept of a *random sample*. Indeed, many statistical analyses directly assume data (or subsets of the data) to be a random sample of conditionally independent observations from some probability model, so that $p(x_1, \dots, x_n | \boldsymbol{\omega}) = \prod_{i=1}^n p(x_i | \boldsymbol{\omega})$; but *any* random sample is exchangeable, since $\prod_{i=1}^n p(x_i | \boldsymbol{\omega})$ is obviously invariant under permutations. Notice that the observations in a random sample are only independent *conditional* on the parameter value $\boldsymbol{\omega}$; as nicely put by Lindley, the mantra that the observations $\{x_1, \dots, x_n\}$ in a random sample are independent is ridiculous when they are used to infer x_{n+1} .

Notice also that, under exchangeability, the general representation theorem provides an *existence theorem* for a probability distribution $p(\omega | C)$ on the parameter space Ω , and that this is an argument which only depends on mathematical probability theory.

Another important consequence of exchangeability is that it provides a formal *definition* of the parameter ω which labels the model as the limit, as $n \rightarrow \infty$, of *some* function $f(x_1, \dots, x_n)$ of the observations; the function f obviously depends both on the assumed model and the chosen parametrization. For instance, in the case of a sequence of Bernoulli trials, the parameter θ is *defined* as the limit, as $n \rightarrow \infty$, of the relative frequency r/n . It follows that, under exchangeability, the sentence “the true value of ω ” has a well-defined meaning, if only asymptotically verifiable. Moreover, if two different models have parameters which are functionally related by their definition, then the corresponding posterior distributions may be meaningfully compared, for they refer to functionally related quantities. For instance, if a finite subset $\{x_1, \dots, x_n\}$ of an exchangeable sequence of integer observations is assumed to be a random sample from a Poisson distribution $\text{Po}(x | \lambda)$, so that $E[x | \lambda] = \lambda$, then λ is *defined* as $\lim_{n \rightarrow \infty} \{\bar{x}_n\}$, where $\bar{x}_n = \sum_j x_j/n$; similarly, if for some fixed non-zero integer r , the same data are assumed to be a random sample for a negative binomial $\text{Nb}(x | r, \theta)$, so that $E[x | \theta, r] = r(1 - \theta)/\theta$, then θ is *defined* as $\lim_{n \rightarrow \infty} \{r/(\bar{x}_n + r)\}$. It follows that $\theta \equiv r/(\lambda + r)$ and, hence, θ and $r/(\lambda + r)$ may be treated as the *same* (unknown) quantity whenever this might be needed as, for example, when comparing the relative merits of these alternative probability models.

3. The Bayesian Paradigm

The statistical analysis of some observed data D typically begins with some informal *descriptive* evaluation, which is used to suggest a tentative, formal *probability model* $\{p(D | \omega), \omega \in \Omega\}$ assumed to represent, for some (unknown) value of ω , the probabilistic mechanism which has generated the observed data D . The arguments outlined in Section 2 establish the logical need to assess a *prior* probability distribution $p(\omega | K)$ over the parameter space Ω , describing the available knowledge K about the value of ω prior to the data being observed. It then follows from standard probability theory that, if the probability model is correct, all available information about the value of ω after the data D have been observed is contained in the corresponding *posterior* distribution whose probability density, $p(\omega | D, A, K)$, is immediately obtained from Bayes’ theorem,

$$p(\omega | D, A, K) = \frac{p(D | \omega) p(\omega | K)}{\int_{\Omega} p(D | \omega) p(\omega | K) d\omega}, \quad (7)$$

where A stands for the assumptions made on the probability model. It is this systematic use of Bayes’ theorem to incorporate the information provided by the

data that justifies the adjective *Bayesian* by which the paradigm is usually known. It is obvious from Bayes' theorem that any value of ω with zero prior density will have zero posterior density. Thus, it is typically assumed (by appropriate restriction, if necessary, of the *parameter space* Ω) that prior distributions are *strictly positive* (as Savage put it, keep the mind open, or at least ajar). To simplify the presentation, the accepted assumptions A and the available knowledge K are often omitted from the notation, but the fact that *all* statements about ω given D are *also* conditional to A and K should always be kept in mind.

Example 1. (Bayesian inference with a finite parameter space). Let $p(D|\theta)$, $\theta \in \{\theta_1, \dots, \theta_m\}$, be the probability mechanism which is assumed to have generated the observed data D , so that θ may only take a *finite* number of values. Using the finite form of Bayes' theorem, and omitting the prevailing conditions from the notation, the posterior probability of θ_i after data D have been observed is

$$\Pr(\theta_i | D) = \frac{p(D|\theta_i) \Pr(\theta_i)}{\sum_{j=1}^m p(D|\theta_j) \Pr(\theta_j)}, \quad i = 1, \dots, m. \quad (8)$$

For any prior distribution $p(\theta) = \{\Pr(\theta_1), \dots, \Pr(\theta_m)\}$ describing available knowledge about the value of θ , $\Pr(\theta_i | D)$ measures how likely should θ_i be judged, given both the initial knowledge described by the prior distribution, and the information provided by the data D .

An important, frequent application of this simple technique is provided by probabilistic diagnosis. For example, consider the simple situation where a particular test designed to detect a virus is known from laboratory research to give a positive result in 98% of infected people and in 1% of non-infected. Then, the posterior probability that a person who tested positive is infected is given by $\Pr(V|+) = (0.98p)/\{0.98p + 0.01(1-p)\}$ as a function of $p = \Pr(V)$, the prior probability of a person being infected (the *prevalence* of the infection in the population under study). Figure 1 shows $\Pr(V|+)$ as a function of $\Pr(V)$.

As one would expect, the posterior probability is only zero if the prior probability is zero (so that it is *known* that the population is free of infection) and it is only one if the prior probability is one (so that it is *known* that the population is universally infected). Notice that if the infection is rare, then the posterior probability of a randomly chosen person being infected will be relatively low even if the test is positive. Indeed, for say $\Pr(V) = 0.002$, one finds $\Pr(V|+) = 0.164$, so that in a population where only 0.2% of individuals are infected, only 16.4% of those testing positive within a random sample will actually prove to be infected: most positives would actually be *false* positives.

In this section, we describe in some detail the learning process described by Bayes' theorem, discuss its implementation in the presence of nuisance parameters,

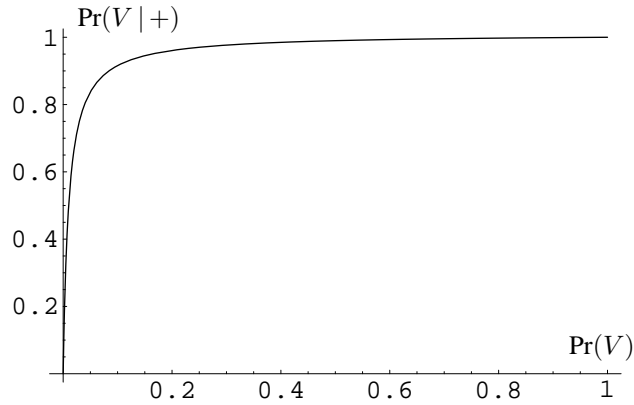


Figure 1. Posterior probability of infection $\Pr(V | +)$ given a positive test, as a function of the prior probability of infection $\Pr(V)$.

show how it can be used to forecast the value of future observations, and analyse its large sample behaviour.

3.1. The Learning Process

In the Bayesian paradigm, the process of learning from the data is systematically implemented by making use of Bayes' theorem to combine the available prior information with the information provided by the data to produce the required posterior distribution. Computation of posterior densities is often facilitated by noting that Bayes' theorem may be simply expressed as

$$p(\omega | D) \propto p(D | \omega) p(\omega), \quad (9)$$

(where \propto stands for 'proportional to' and where, for simplicity, the accepted assumptions A and the available knowledge K have been omitted from the notation), since the missing proportionality constant $[\int_{\Omega} p(D | \omega) p(\omega) d\omega]^{-1}$ may always be deduced from the fact that $p(\omega | D)$, a probability density, must integrate to one. Hence, to identify the form of a posterior distribution it suffices to identify a *kernel* of the corresponding probability density, that is a function $k(\omega)$ such that

$$p(\omega | D) = c(D) k(\omega)$$

for some $c(D)$ which does not involve ω . In the examples which follow, this technique will often be used.

An *improper prior function* is defined as a positive function $\pi(\omega)$ such that $\int_{\Omega} \pi(\omega) d\omega$ is not finite. Equation (9), the formal expression of Bayes' theorem, remains technically valid if $p(\omega)$ is replaced by an improper prior function $\pi(\omega)$ provided the proportionality constant exists, thus leading to a well defined *proper* posterior density $\pi(\omega | D) \propto p(D | \omega)\pi(\omega)$. It will later be established (Section 5) that Bayes' theorem also remains philosophically valid if $p(\omega)$ is replaced by an appropriately chosen reference “noninformative” (typically improper) prior function $\pi(\omega)$.

Considered as a function of ω , $l(\omega, D) = p(D | \omega)$ is often referred to as the *likelihood function*. Thus, Bayes' theorem is simply expressed in words by the statement that *the posterior is proportional to the likelihood times the prior*. It follows from equation (9) that, provided the *same* prior $p(\omega)$ is used, two different data sets D_1 and D_2 , with possibly different probability models $p_1(D_1 | \omega)$ and $p_2(D_2 | \omega)$ but yielding *proportional* likelihood functions, will produce identical posterior distributions for ω . This immediate consequence of Bayes theorem has been proposed as a principle on its own, the *likelihood principle*, and it is seen by many as an obvious requirement for reasonable statistical inference. In particular, for any given prior $p(\omega)$, the posterior distribution does not depend on the set of possible data values, or the *outcome space*. Notice, however, that the likelihood principle only applies to inferences about the parameter vector ω once the data have been obtained. Consideration of the outcome space is essential, for instance, in model criticism, in the design of experiments, in the derivation of predictive distributions, or (see Section 5) in the construction of objective Bayesian procedures.

Naturally, the terms prior and posterior are only *relative* to a particular set of data. As one would expect from the coherence induced by probability theory, if data $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ are sequentially presented, the final result will be the same whether data are globally or sequentially processed. Indeed, $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_{i+1}) \propto p(\mathbf{x}_{i+1} | \omega) p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_i)$, for $i = 1, \dots, n - 1$, so that the “posterior” at a given stage becomes the “prior” at the next.

In most situations, the posterior distribution is “sharper” than the prior so that, in most cases, the density $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_{i+1})$ will be more concentrated around the true value of ω than $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_i)$. However, this is not always the case: occasionally, a “surprising” observation will increase, rather than decrease, the uncertainty about the value of ω . For instance, in probabilistic diagnosis, a sharp posterior probability distribution (over the possible causes $\{\omega_1, \dots, \omega_k\}$ of a syndrome) describing, a “clear” diagnosis of disease ω_i (that is, a posterior with a large probability for ω_i) would typically update to a less concentrated posterior probability distribution over $\{\omega_1, \dots, \omega_k\}$ if a new clinical analysis yielded data which were unlikely under ω_i .

For a given probability model, one may find that some particular function of the data $\mathbf{t} = \mathbf{t}(D)$ is a *sufficient* statistic in the sense that, given the model, $\mathbf{t}(D)$

contains all information about ω which is available in D . Formally, $\mathbf{t} = \mathbf{t}(D)$ is sufficient if (and only if) there exist nonnegative functions f and g such that the likelihood function may be factorized in the form $p(D | \omega) = f(\omega, \mathbf{t})g(D)$. A sufficient statistic always exists, for $\mathbf{t}(D) = D$ is obviously sufficient; however, a much simpler sufficient statistic, with a fixed dimensionality which is independent of the sample size, often exists. In fact this is known to be the case whenever the probability model belongs to the *generalized exponential family*, which includes many of the more frequently used probability models. It is easily established that if \mathbf{t} is sufficient, the posterior distribution of ω only depends on the data D through $\mathbf{t}(D)$, and may be directly computed in terms of $p(\mathbf{t} | \omega)$, so that,

$$p(\omega | D) = p(\omega | \mathbf{t}) \propto p(\mathbf{t} | \omega) p(\omega).$$

Naturally, for fixed data and model assumptions, different priors lead to different posteriors. Indeed, Bayes' theorem may be described as a data-driven probability transformation machine which maps prior distributions (describing prior knowledge) into posterior distributions (representing combined prior and data knowledge). It is important to analyse whether or not sensible changes in the prior would induce noticeable changes in the posterior. Posterior distributions based on reference "noninformative" priors play a central role in this *sensitivity analysis* context. Investigation of the sensitivity of the posterior to changes in the prior is an important ingredient of the comprehensive analysis of the sensitivity of the final results to *all* accepted assumptions which any responsible statistical study should contain.

Example 2. (Inference on a binomial parameter). If the data D consist of n Bernoulli observations with parameter θ which contain r positive trials, then

$$p(D | \theta, n) = \theta^r (1 - \theta)^{n-r},$$

so that $\mathbf{t}(D) = \{r, n\}$ is sufficient. Suppose that prior knowledge about θ is described by a Beta distribution $\text{Be}(\theta | \alpha, \beta)$, so that

$$p(\theta | \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}.$$

Using Bayes' theorem, the posterior density of θ is

$$p(\theta | r, n, \alpha, \beta) \propto \theta^r (1 - \theta)^{n-r} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \propto \theta^{r+\alpha-1} (1 - \theta)^{n-r+\beta-1},$$

which is the Beta distribution $\text{Be}(\theta | r + \alpha, n - r + \beta)$.

Suppose, for example, that in the light of precedent surveys, available information on the proportion θ of citizens who would vote for a particular political measure in a referendum is described by a Beta distribution $\text{Be}(\theta | 50, 50)$, so that

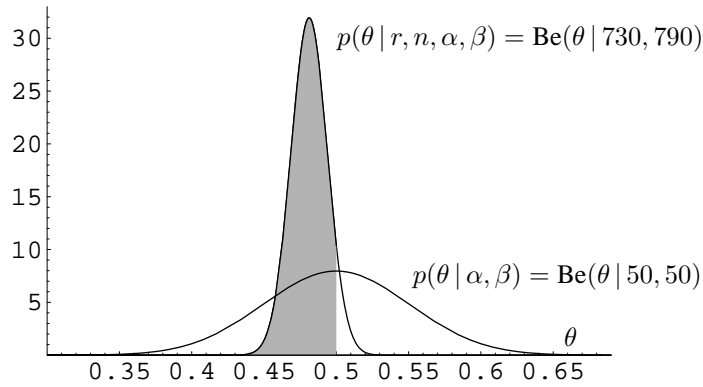


Figure 2. Prior and posterior densities of the proportion θ of citizens that would vote in favour of a referendum.

it is judged to be equally likely that the referendum would be won or lost, and it is judged that the probability that either side wins less than 60% of the vote is 0.95.

A random survey of size 1500 is then conducted, where only 720 citizens declare to be in favour of the proposed measure. Using the results above, the corresponding posterior distribution is then $\text{Be}(\theta | 770, 830)$. These prior and the posterior densities are plotted in Figure 2; it may be appreciated that, as one would expect, the effect of the data is to drastically reduce the initial uncertainty on the value of θ and, hence, on the referendum outcome. More precisely,

$$\Pr(\theta < 0.5 | 720, 1500, H, K) = 0.933$$

(shaded region in Figure 2) so that, after the information from the survey has been included, the probability that the referendum will be lost should be judged to be about 93%.

The general situation where the vector of interest is not the whole parameter vector ω , but some function $\theta = \theta(\omega)$ of possibly lower dimension than ω , will now be considered. Let D be some observed data, let $\{p(D | \omega), \omega \in \Omega\}$ be a probability model assumed to describe the probability mechanism which has generated D , let $p(\omega)$ be a probability distribution describing any available information on the value of ω , and let $\theta = \theta(\omega) \in \Theta$ be a function of the original parameters over whose value inferences based on the data D are required. Any valid conclusion on the value of the *vector of interest* θ will then be contained in its posterior probability distribution $p(\theta | D)$ which is conditional on the observed data D and will naturally also depend, although not explicitly shown in the notation, on the assumed model

$\{p(D|\omega), \omega \in \Omega\}$, and on the available prior information encapsulated by $p(\omega)$. The required posterior distribution $p(\theta|D)$ is found by standard use of probability calculus. Indeed, by Bayes' theorem, $p(\omega|D) \propto p(D|\omega)p(\omega)$. Moreover, let $\lambda = \lambda(\omega) \in \Lambda$ be some other function of the original parameters such that $\psi = \{\theta, \lambda\}$ is a one-to-one transformation of ω , and let $J(\omega) = (\partial\psi/\partial\omega)$ be the corresponding Jacobian matrix. Naturally, the introduction of λ is not necessary if $\theta(\omega)$ is a one-to-one transformation of ω . Using standard change-of-variable probability techniques, the posterior density of ψ is

$$p(\psi|D) = p(\theta, \lambda|D) = \left[\frac{p(\omega|D)}{|J(\omega)|} \right]_{\omega=\omega(\psi)} \quad (10)$$

and the required posterior of θ is the appropriate *marginal* density, obtained by integration over the *nuisance parameter* λ ,

$$p(\theta|D) = \int_{\Lambda} p(\theta, \lambda|D) d\lambda. \quad (11)$$

Notice that elimination of unwanted nuisance parameters, a simple integration within the Bayesian paradigm is, however, a difficult (often polemic) problem for conventional statistics.

Sometimes, the range of possible values of ω is effectively restricted by contextual considerations. If ω is known to belong to $\Omega_c \subset \Omega$, the prior distribution is only positive in Ω_c and, using Bayes' theorem, it is immediately found that the restricted posterior is

$$p(\omega|D, \omega \in \Omega_c) = \frac{p(\omega|D)}{\int_{\Omega_c} p(\omega|D)}, \quad \omega \in \Omega_c, \quad (12)$$

and obviously vanishes if $\omega \notin \Omega_c$. Thus, to incorporate a restriction on the possible values of the parameters, it suffices to *renormalize* the unrestricted posterior distribution to the set $\Omega_c \subset \Omega$ of parameter values which satisfy the required condition. Incorporation of known constraints on the parameter values, a simple renormalization within the Bayesian paradigm, is another very difficult problem for conventional statistics.

Example 3. (Inference on normal parameters). Let $D = \{x_1, \dots, x_n\}$ be a random sample from a normal distribution $N(x|\mu, \sigma)$. The corresponding likelihood function is immediately found to be

$$p(D|\mu, \sigma) \propto \sigma^{-n} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)],$$

with $n\bar{x} = \sum_i x_i$, and $ns^2 = \sum_i (x_i - \bar{x})^2$. It may be shown (see Section 5) that absence of initial information on the value of both μ and σ may formally be

described by a joint prior function which is uniform in both μ and $\log(\sigma)$, that is, by the (improper) prior function $p(\mu, \sigma) = \sigma^{-1}$. Using Bayes' theorem, the corresponding joint posterior is

$$p(\mu, \sigma | D) \propto \sigma^{-(n+1)} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)]. \quad (13)$$

Thus, using the Gamma integral in terms of $\lambda = \sigma^{-2}$ to integrate out σ ,

$$p(\mu | D) \propto \int_0^\infty \sigma^{-(n+1)} \exp\left[-\frac{n}{2\sigma^2}[s^2 + (\bar{x} - \mu)^2]\right] d\sigma \propto [s^2 + (\bar{x} - \mu)^2]^{-n/2}, \quad (14)$$

which is recognized as a kernel of the Student density $\text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$. Similarly, integrating out μ ,

$$p(\sigma | D) \propto \int_{-\infty}^\infty \sigma^{-(n+1)} \exp\left[-\frac{n}{2\sigma^2}[s^2 + (\bar{x} - \mu)^2]\right] d\mu \propto \sigma^{-n} \exp\left[-\frac{ns^2}{2\sigma^2}\right]. \quad (15)$$

Changing variables to the precision $\lambda = \sigma^{-2}$ results in $p(\lambda | D) \propto \lambda^{(n-3)/2} e^{ns^2\lambda/2}$, a kernel of the Gamma density $\text{Ga}(\lambda | (n-1)/2, ns^2/2)$. In terms of the standard deviation σ this becomes $p(\sigma | D) = p(\lambda | D) |\partial\lambda/\partial\sigma| = 2\sigma^{-3} \text{Ga}(\sigma^{-2} | (n-1)/2, ns^2/2)$, a square-root inverted gamma density.

A frequent example of this scenario is provided by laboratory measurements made in conditions where central limit conditions apply, so that (assuming no experimental bias) those measurements may be treated as a random sample from a normal distribution centred at the quantity μ which is being measured, and with some (unknown) standard deviation σ . Suppose, for example, that in an elementary physics classroom experiment to measure the gravitational field g with a pendulum, a student has obtained $n = 20$ measurements of g yielding (in m/sec²) a mean $\bar{x} = 9.8087$, and a standard deviation $s = 0.0428$. Using no other information, the corresponding posterior distribution is $p(g | D) = \text{St}(g | 9.8087, 0.0098, 19)$ represented in Figure 3(a). In particular, $\text{Pr}(9.788 < g < 9.829 | D) = 0.95$, so that, with the information provided by this experiment, the gravitational field at the location of the laboratory may be expected to lie between 9.788 and 9.829 with probability 0.95.

Formally, the posterior distribution of g should be restricted to $g > 0$; however, as immediately obvious from Figure 3a, this would not have any appreciable effect, due to the fact that the likelihood function is actually concentrated on positive g values.

Suppose now that the student is further instructed to incorporate into the analysis the fact that the value of the gravitational field g at the laboratory is known

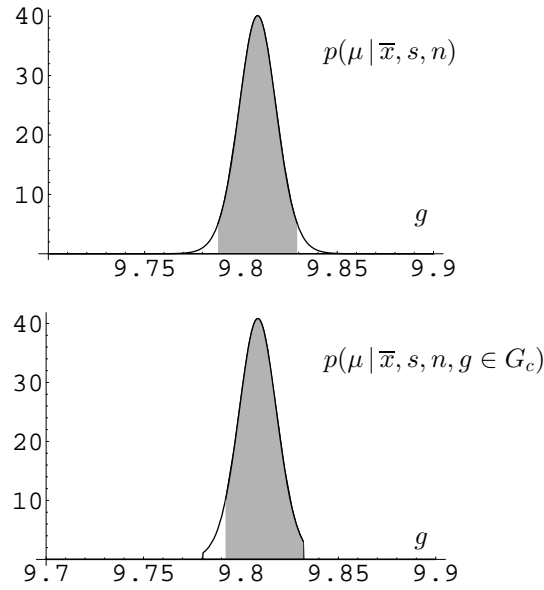


Figure 3. Posterior density $p(g | m, s, n)$ of the value g of the gravitational field, given $n = 20$ normal measurements with mean $m = 9.8087$ and standard deviation $s = 0.0428$, with no additional information (upper panel), and with the value of g restricted to region $G_c = \{g; 9.7803 < g < 9.8322\}$ (lower panel). Shaded areas represent 95%-credible regions of g .

to lie between 9.7803 m/sec² (average value at the Equator) and 9.8322 m/sec² (average value at the poles). The updated posterior distribution will be

$$p(g | D, g \in G_c) = \frac{\text{St}(g | m, s/\sqrt{n-1}, n)}{\int_{g \in G_c} \text{St}(g | m, s/\sqrt{n-1}, n)}, \quad g \in G_c, \quad (16)$$

represented in Figure 3(b), where $G_c = \{g; 9.7803 < g < 9.8322\}$. One-dimensional numerical integration may be used to verify that $\Pr(g > 9.792 | D, g \in G_c) = 0.95$. Moreover, if inferences about the standard deviation σ of the measurement procedure are also requested, the corresponding posterior distribution is found to be $p(\sigma | D) = 2\sigma^{-3}\text{Ga}(\sigma^{-2} | 9.5, 0.0183)$. This has a mean $E[\sigma | D] = 0.0458$ and yields $\Pr(0.0334 < \sigma < 0.0642 | D) = 0.95$.

3.2. Predictive Distributions

Let $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\mathbf{x}_i \in X$, be a set of exchangeable observations, and consider now a situation where it is desired to predict the value of a future observation $\mathbf{x} \in X$ generated by the same random mechanism that has generated the data D . It follows from the foundations arguments discussed in Section 2 that the solution to this prediction problem is simply encapsulated by the *predictive* distribution $p(\mathbf{x} | D)$ describing the uncertainty on the value that \mathbf{x} will take, given the information provided by D and any other available knowledge. Suppose that contextual information suggests the assumption that data D may be considered to be a random sample from a distribution in the family $\{p(\mathbf{x} | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$, and let $p(\boldsymbol{\omega})$ be a prior distribution describing available information on the value of $\boldsymbol{\omega}$. Since $p(\mathbf{x} | \boldsymbol{\omega}, D) = p(\mathbf{x} | \boldsymbol{\omega})$, it then follows from standard probability theory that

$$p(\mathbf{x} | D) = \int_{\Omega} p(\mathbf{x} | \boldsymbol{\omega}) p(\boldsymbol{\omega} | D) d\boldsymbol{\omega},$$

which is an average of the probability distributions of \mathbf{x} conditional on the (unknown) value of $\boldsymbol{\omega}$, weighted with the posterior distribution of $\boldsymbol{\omega}$ given D .

If the assumptions on the probability model are correct, the posterior predictive distribution $p(\mathbf{x} | D)$ will converge, as the sample size increases, to the distribution $p(\mathbf{x} | \boldsymbol{\omega})$ which has generated the data. Indeed, the best technique to assess the quality of the inferences about $\boldsymbol{\omega}$ encapsulated in $p(\boldsymbol{\omega} | D)$ is to check against the observed data the predictive distribution $p(\mathbf{x} | D)$ generated by $p(\boldsymbol{\omega} | D)$.

Example 4. (Prediction in a Poisson process). Let $D = \{r_1, \dots, r_n\}$ be a random sample from a Poisson distribution $\text{Pn}(r | \lambda)$ with parameter λ , so that $p(D | \lambda) \propto \lambda^t e^{-\lambda n}$, where $t = \sum r_i$. It may be shown (see Section 5) that absence of initial information on the value of λ may be formally described by the (improper) prior function $p(\lambda) = \lambda^{-1/2}$. Using Bayes' theorem, the corresponding posterior is

$$p(\lambda | D) \propto \lambda^t e^{-\lambda n} \lambda^{-1/2} \propto \lambda^{t-1/2} e^{-\lambda n}, \quad (17)$$

the kernel of a Gamma density $\text{Ga}(\lambda |, t + 1/2, n)$, with mean $(t + 1/2)/n$. The corresponding predictive distribution is the Poisson-Gamma mixture

$$p(r | D) = \int_0^{\infty} \text{Pn}(r | \lambda) \text{Ga}(\lambda |, t + \frac{1}{2}, n) d\lambda = \frac{n^{t+1/2}}{\Gamma(t + 1/2)} \frac{1}{r!} \frac{\Gamma(r + t + 1/2)}{(1 + n)^{r+t+1/2}}. \quad (18)$$

Suppose, for example, that in a firm producing automobile restraint systems, the entire production in each of 10 consecutive months has yielded no complaint from their clients. With no additional information on the average number λ of

complaints per month, the quality assurance department of the firm may report that the probabilities that r complaints will be received in the next month of production are given by equation (18), with $t = 0$ and $n = 10$. In particular, $p(r = 0 | D) = 0.953$, $p(r = 1 | D) = 0.043$, and $p(r = 2 | D) = 0.003$. Many other situations may be described with the same model. For instance, if meteorological conditions remain similar in a given area, $p(r = 0 | D) = 0.953$ would describe the chances of no flash flood next year, given 10 years without flash floods in the area.

Example 5. (Prediction in a Normal process). Consider now prediction of a continuous variable. Let $D = \{x_1, \dots, x_n\}$ be a random sample from a normal distribution $N(x | \mu, \sigma)$. As mentioned in Example 3, absence of initial information on the values of both μ and σ is formally described by the *improper* prior function $p(\mu, \sigma) = \sigma^{-1}$, and this leads to the joint posterior density (13). The corresponding (posterior) predictive distribution is

$$p(x | D) = \int_0^\infty \int_{-\infty}^\infty N(x | \mu, \sigma) p(\mu, \sigma | D) d\mu d\sigma$$

leading to

$$p(x | D) = \text{St}\left(x \mid \bar{x}, s\sqrt{\frac{n+1}{n-1}}, n-1\right). \quad (19)$$

If μ is known to be positive, the appropriate prior function will be the restricted prior function

$$p(\mu, \sigma) = \begin{cases} \sigma^{-1} & \text{if } \mu > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (20)$$

However, the result in equation (19) will still hold, provided the likelihood function $p(D | \mu, \sigma)$ is concentrated on positive μ values. Suppose, for example, that in the firm producing automobile restraint systems, the observed breaking strengths of $n = 10$ randomly chosen safety belt webbings have mean $\bar{x} = 28.011$ kN and standard deviation $s = 0.443$ kN, and that the relevant engineering specification requires breaking strengths to be larger than 26 kN. If data may truly be assumed to be a random sample from a normal distribution, the likelihood function is only appreciable for positive μ values, and only the information provided by this small sample is to be used, then the quality engineer may claim that the probability that a safety belt randomly chosen from the same batch as the sample tested would satisfy the required specification is $\Pr(x > 26 | D) = 0.9987$. Besides, if production conditions remain constant, 99.87% of the safety belt webbings may be expected to have acceptable breaking strengths.

3.3. Asymptotic Behaviour

The behaviour of posterior distributions when the sample size is large is now considered. This is important for, at least, two different reasons: (i) asymptotic results provide useful first-order approximations when actual samples are relatively large, and (ii) objective Bayesian methods typically depend on the asymptotic properties of the assumed model. Let $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\mathbf{x} \in X$, be a random sample of size n from $\{p(\mathbf{x} | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$. It may be shown that, as $n \rightarrow \infty$, the posterior distribution $p(\boldsymbol{\omega} | D)$ of a *discrete* parameter $\boldsymbol{\omega}$ typically converges to a degenerate distribution which gives probability one to the true value of $\boldsymbol{\omega}$, and that the posterior distribution of a *continuous* parameter $\boldsymbol{\omega}$ typically converges to a normal distribution centred at its *maximum likelihood estimate* $\hat{\boldsymbol{\omega}}$ (MLE), with a variance matrix which decreases with n as $1/n$.

Consider first the situation where $\Omega = \{\boldsymbol{\omega}_1, \boldsymbol{\omega}_2, \dots\}$ consists of a *countable* (possibly infinite) set of values, such that the probability model which corresponds to the true parameter value $\boldsymbol{\omega}_t$ is *distinguishable* from the others in the sense that the logarithmic divergence $\delta\{p(\mathbf{x} | \boldsymbol{\omega}_i) | p(\mathbf{x} | \boldsymbol{\omega}_t)\}$ of each of the $p(\mathbf{x} | \boldsymbol{\omega}_i)$ from $p(\mathbf{x} | \boldsymbol{\omega}_t)$ is strictly positive. Taking logarithms in Bayes' theorem, defining $z_j = \log[p(\mathbf{x}_j | \boldsymbol{\omega}_i) / p(\mathbf{x}_j | \boldsymbol{\omega}_t)]$, $j = 1, \dots, n$, and using the strong law of large numbers on the n conditionally independent and identically distributed random quantities z_1, \dots, z_n , it may be shown that

$$\lim_{n \rightarrow \infty} p(\boldsymbol{\omega}_t | \mathbf{x}_1, \dots, \mathbf{x}_n) = 1, \quad \lim_{n \rightarrow \infty} p(\boldsymbol{\omega}_i | \mathbf{x}_1, \dots, \mathbf{x}_n) = 0, \quad i \neq t. \quad (21)$$

Thus, under appropriate regularity conditions, the posterior probability of the true parameter value converges to one as the sample size grows.

Consider now the situation where $\boldsymbol{\omega}$ is a k -dimensional *continuous* parameter. Expressing Bayes' theorem as

$$p(\boldsymbol{\omega} | \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \exp\{\log[p(\boldsymbol{\omega})] + \sum_{j=1}^n \log[p(\mathbf{x}_j | \boldsymbol{\omega})]\},$$

expanding $\sum_j \log[p(\mathbf{x}_j | \boldsymbol{\omega})]$ about its maximum (the MLE $\hat{\boldsymbol{\omega}}$), and assuming regularity conditions (to ensure that terms of order higher than quadratic may be ignored and that the sum of the terms from the likelihood will dominate the term from the prior) it is found that the posterior density of $\boldsymbol{\omega}$ is the approximate k -variate normal

$$p(\boldsymbol{\omega} | \mathbf{x}_1, \dots, \mathbf{x}_n) \approx N_k\{\hat{\boldsymbol{\omega}}, \mathbf{S}(D, \hat{\boldsymbol{\omega}})\}, \quad (22)$$

where

$$\mathbf{S}^{-1}(D, \boldsymbol{\omega}) = \left(- \sum_{l=1}^n \frac{\partial^2 \log[p(\mathbf{x}_l | \boldsymbol{\omega})]}{\partial \omega_i \partial \omega_j} \right).$$

A simpler, but somewhat poorer, approximation may be obtained by using the strong law of large numbers on the sums in (22) to establish that

$$\mathbf{S}^{-1}(D, \hat{\boldsymbol{\omega}}) \approx n \mathbf{F}(\hat{\boldsymbol{\omega}}),$$

where $\mathbf{F}(\boldsymbol{\omega})$ is Fisher's information matrix, with general element

$$\mathbf{F}_{ij}(\boldsymbol{\omega}) = - \int_X p(\mathbf{x} | \boldsymbol{\omega}) \frac{\partial^2 \log[p(\mathbf{x} | \boldsymbol{\omega})]}{\partial \omega_i \partial \omega_j} d\mathbf{x}, \quad (23)$$

so that

$$p(\boldsymbol{\omega} | \mathbf{x}_1, \dots, \mathbf{x}_n) \approx N_k(\boldsymbol{\omega} | \hat{\boldsymbol{\omega}}, n^{-1} \mathbf{F}^{-1}(\hat{\boldsymbol{\omega}})). \quad (24)$$

Thus, under appropriate regularity conditions, the posterior probability density of the parameter vector $\boldsymbol{\omega}$ approaches, as the sample size grows, a multivariate normal density centred at the MLE $\hat{\boldsymbol{\omega}}$, with a dispersion matrix which decreases with n as n^{-1} .

Example 2. (Inference on a binomial parameter, continued). Let $D = (x_1, \dots, x_n)$ consist of n independent Bernoulli trials with parameter θ , so that $p(D | \theta, n) = \theta^r (1 - \theta)^{n-r}$. This likelihood function is maximized at $\hat{\theta} = r/n$, and Fisher's information function is $F(\theta) = \theta^{-1}(1 - \theta)^{-1}$. Thus, using the results above, the posterior distribution of θ will be the approximate normal,

$$p(\theta | r, n) \approx N(\theta | \hat{\theta}, s(\hat{\theta})/\sqrt{n}), \quad s(\theta) = \{\theta(1 - \theta)\}^{1/2} \quad (25)$$

with mean $\hat{\theta} = r/n$ and variance $\hat{\theta}(1 - \hat{\theta})/n$. This will provide a reasonable approximation to the exact posterior if (i) the prior $p(\theta)$ is relatively "flat" in the region where the likelihood function matters, and (ii) both r and n are moderately large. If, say, $n = 1500$ and $r = 720$, this leads to $p(\theta | D) \approx N(\theta | 0.480, 0.013)$, and to $\Pr(\theta > 0.5 | D) \approx 0.940$, which may be compared with the exact value $\Pr(\theta > 0.5 | D) = 0.933$ obtained from the posterior distribution which corresponds to the prior $\text{Be}(\theta | 50, 50)$.

It follows from the *joint* posterior asymptotic behaviour of $\boldsymbol{\omega}$ and from the properties of the multivariate normal distribution that, if the parameter vector is decomposed into $\boldsymbol{\omega} = (\boldsymbol{\theta}, \boldsymbol{\lambda})$, and Fisher's information matrix is correspondingly partitioned, so that

$$\mathbf{F}(\boldsymbol{\omega}) = \mathbf{F}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \begin{pmatrix} \mathbf{F}_{\theta\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{F}_{\theta\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \\ \mathbf{F}_{\lambda\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{F}_{\lambda\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \end{pmatrix} \quad (26)$$

and

$$\mathbf{S}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbf{F}^{-1}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \begin{pmatrix} \mathbf{S}_{\theta\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{S}_{\theta\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \\ \mathbf{S}_{\lambda\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{S}_{\lambda\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \end{pmatrix}, \quad (27)$$

then the *marginal* posterior distribution of $\boldsymbol{\theta}$ will be

$$p(\boldsymbol{\theta} | D) \approx \mathbf{N}\{\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}, n^{-1} \mathbf{S}_{\theta\theta}(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\lambda}})\}, \quad (28)$$

while the *conditional* posterior distribution of $\boldsymbol{\lambda}$ given $\boldsymbol{\theta}$ will be

$$p(\boldsymbol{\lambda} | \boldsymbol{\theta}, D) \approx \mathbf{N}\{\boldsymbol{\lambda} | \hat{\boldsymbol{\lambda}} - \mathbf{F}_{\lambda\lambda}^{-1}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}}) \mathbf{F}_{\lambda\theta}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}}) (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}), n^{-1} \mathbf{F}_{\lambda\lambda}^{-1}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}})\}. \quad (29)$$

Notice that $\mathbf{F}_{\lambda\lambda}^{-1} = \mathbf{S}_{\lambda\lambda}$ if (and only if) \mathbf{F} is block diagonal, *i.e.*, if (and only if) $\boldsymbol{\theta}$ and $\boldsymbol{\lambda}$ are asymptotically independent.

Example 3. (Inference on normal parameters, continued). Let $D = (x_1, \dots, x_n)$ be a random sample from a normal distribution $\mathbf{N}(x | \mu, \sigma)$. The corresponding likelihood function $p(D | \mu, \sigma)$ is maximized at $(\hat{\mu}, \hat{\sigma}) = (\bar{x}, s)$, and Fisher's matrix is diagonal, with $F_{\mu\mu} = \sigma^{-2}$. Hence, the posterior distribution of μ is *approximately* $\mathbf{N}(\mu | \bar{x}, s/\sqrt{n})$. This may usefully be compared with the *exact* result $p(\mu | D) = \text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$, obtained previously under the assumption of no prior knowledge.

4. Inference Summaries

From a Bayesian viewpoint, the final outcome of a problem of inference about *any* unknown quantity is nothing but the corresponding posterior distribution. Thus, given some data D and conditions C , *all* that can be said about any function $\boldsymbol{\omega}$ of the parameters which govern the model is contained in the posterior distribution $p(\boldsymbol{\omega} | D, C)$, and *all* that can be said about some function \mathbf{y} of future observations from the same model is contained in its posterior predictive distribution $p(\mathbf{y} | D, C)$. As mentioned before, Bayesian inference may technically be described as a decision problem where the space of available actions is the class of those posterior probability distributions of the quantity of interest which are compatible with accepted assumptions.

However, to make it easier for the user to assimilate the appropriate conclusions, it is often convenient to *summarize* the information contained in the posterior distribution by (i) providing values of the quantity of interest which, in the light of the data, are likely to be "close" to its true value and by (ii) measuring the compatibility of the results with hypothetical values of the quantity of interest which might have been suggested in the context of the investigation. In this section, those Bayesian counterparts of traditional *estimation* and *hypothesis testing* problems are briefly considered.

4.1. Estimation

In one or two dimensions, a graph of the posterior probability density of the quantity of interest (or the probability mass function in the discrete case) immediately conveys an intuitive, “impressionist” summary of the main conclusions which may possibly be drawn on its value. Indeed, this is greatly appreciated by users, and may be quoted as an important asset of Bayesian methods. From a plot of its posterior density, the region where (given the data) a univariate quantity of interest is likely to lie is easily distinguished. For instance, all important conclusions about the value of the gravitational field in Example 3 are qualitatively available from Figure 3. However, this does not easily extend to more than two dimensions and, besides, *quantitative* conclusions (in a simpler form than that provided by the mathematical expression of the posterior distribution) are often required.

Point Estimation. Let D be the available data, which are assumed to have been generated by a probability model $\{p(D | \omega), \omega \in \Omega\}$, and let $\theta = \theta(\omega) \in \Theta$ be the quantity of interest. A *point estimator* of θ is some function of the data $\tilde{\theta} = \tilde{\theta}(D)$ which could be regarded as an appropriate proxy for the actual, unknown value of θ . Formally, to choose a point estimate for θ is a *decision problem*, where the action space is the class Θ of possible θ values. From a decision-theoretic perspective, to choose a point estimate $\tilde{\theta}$ of some quantity θ is a *decision* to act as though $\tilde{\theta}$ were θ , not to assert something about the value of θ (although desire to assert something simple may well be the reason to obtain an estimate). As prescribed by the foundations of decision theory (Section 2), to solve this decision problem it is necessary to specify a *loss function* $L(\tilde{\theta}, \theta)$ measuring the consequences of acting *as if* the true value of the quantity of interest were $\tilde{\theta}$, when it is actually θ . The expected posterior loss if $\tilde{\theta}$ were used is

$$\bar{L}[\tilde{\theta} | D] = \int_{\Theta} L(\tilde{\theta}, \theta) p(\theta | D) d\theta, \quad (30)$$

and the corresponding *Bayes estimator* θ^* is that function of the data, $\theta^* = \theta^*(D)$, which minimizes this expectation.

Example 6. (Conventional Bayes estimators). For any given model and data, the Bayes estimator obviously depends on the chosen loss function. The loss function is context specific, and should be chosen in terms of the anticipated uses of the estimate; however, a number of conventional loss functions have been suggested for those situations where no particular uses are envisaged. These loss functions produce estimates which may be regarded as simple descriptions of the *location* of the posterior distribution. For example, if the loss function is quadratic, so that $L(\tilde{\theta}, \theta) = (\tilde{\theta} - \theta)^t(\tilde{\theta} - \theta)$, then the Bayes estimator is the *posterior mean* $\theta^* = E[\theta | D]$, assuming that the mean exists. Similarly, if the loss function is a zero-one function, so that $L(\tilde{\theta}, \theta) = 0$ if $\tilde{\theta}$ belongs to a ball or radius ϵ centred

in θ and $L(\tilde{\theta}, \theta) = 1$ otherwise, then the Bayes estimator θ^* tends to the *posterior mode* as the ball radius ϵ tends to zero, assuming that a unique mode exists. If θ is univariate and the loss function is linear, so that $L(\tilde{\theta}, \theta) = c_1(\tilde{\theta} - \theta)$ if $\tilde{\theta} \geq \theta$, and $L(\tilde{\theta}, \theta) = c_2(\theta - \tilde{\theta})$ otherwise, then the Bayes estimator is the *posterior quantile* of order $c_2/(c_1 + c_2)$, so that $\Pr[\theta < \theta^*] = c_2/(c_1 + c_2)$. In particular, if $c_1 = c_2$, the Bayes estimator is the *posterior median*. The results derived for linear loss functions clearly illustrate the fact that *any* possible parameter value may turn out be the Bayes estimator: it all depends on the loss function describing the consequences of the anticipated uses of the estimate.

Example 7. (Intrinsic estimation). Conventional loss functions are typically non-invariant under reparametrization, so that the Bayes estimator ϕ^* of a one-to-one transformation $\phi = \phi(\theta)$ of the original parameter θ is not necessarily $\phi(\theta^*)$ (the *univariate* posterior median, which *is* invariant, is an interesting exception). Moreover, conventional loss functions focus on the “distance” between the estimate $\tilde{\theta}$ and the true value θ , rather than on the “distance” between the probability models they label. Intrinsic losses directly focus on how different the probability *model* $p(D | \theta, \lambda)$ is from its closest approximation within the family

$$\{p(D | \tilde{\theta}, \lambda_i), \lambda_i \in \Lambda\},$$

and typically produce invariant solutions. An attractive example is the *intrinsic discrepancy*, $\delta(\tilde{\theta}, \theta)$ defined as the minimum logarithmic divergence between a probability model labelled by θ and a probability model labelled by $\tilde{\theta}$. When there are no nuisance parameters, this is given by

$$\delta(\tilde{\theta}, \theta) = \min\{k(\tilde{\theta} | \theta), k(\theta | \tilde{\theta})\} \quad (31)$$

where

$$k(\theta_i | \theta_j) = \int_T p(\mathbf{t} | \theta_j) \log \frac{p(\mathbf{t} | \theta_j)}{p(\mathbf{t} | \theta_i)} d\mathbf{t},$$

and $\mathbf{t} = \mathbf{t}(D) \in T$ is *any* sufficient statistic (which may well be the whole data set D). The definition is easily extended to problems with nuisance parameters; in this case,

$$\delta(\tilde{\theta}, \theta, \lambda) = \min_{\lambda_i \in \Lambda} \delta(\tilde{\theta}, \lambda_i, \theta, \lambda) \quad (32)$$

measures the logarithmic divergence from $p(\mathbf{t} | \theta, \lambda)$ of its closest approximation with $\theta = \tilde{\theta}$, and the loss function now depends on the complete parameter vector (θ, λ) . Although not explicitly shown in the notation, the intrinsic discrepancy function typically depends on the sample size n ; indeed, when the data consist of a random sample $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ from some model $p(\mathbf{x} | \theta)$ then

$$k(\theta_i | \theta) = n \int_X p(\mathbf{x} | \theta) \log \frac{p(\mathbf{x} | \theta)}{p(\mathbf{x} | \theta_i)} d\mathbf{x}, \quad (33)$$

so that the intrinsic discrepancy associated with the full model is simply n times the intrinsic discrepancy which corresponds to a single observation. The intrinsic discrepancy is a symmetric, non-negative loss function with a direct interpretation in information-theoretic terms as the minimum amount of information which is expected to be necessary to distinguish between the model $p(D | \theta, \lambda)$ and its closest approximation within the class $\{p(D | \tilde{\theta}, \lambda_i), \lambda_i \in \Lambda\}$. Moreover, it is invariant under one-to-one reparametrization of the parameter of interest θ , and does not depend on the choice of the nuisance parameter λ . The *intrinsic estimator* is naturally obtained by minimizing the posterior expected intrinsic discrepancy

$$d(\tilde{\theta} | D) = \int_{\Lambda} \int_{\Theta} \delta(\tilde{\theta}, \theta, \lambda) p(\theta, \lambda | D) d\theta d\lambda. \quad (34)$$

Since the intrinsic discrepancy is invariant under reparametrization, minimizing its posterior expectation produces invariant estimators.

Example 2. (Inference on a binomial parameter, continued). In the estimation of a binomial proportion θ , given data $D = (n, r)$ and a Beta prior $\text{Be}(\theta | \alpha, \beta)$, the Bayes estimator associated with the quadratic loss (the corresponding posterior mean) is $E[\theta | D] = (r + \alpha)/(n + \alpha + \beta)$, while the quadratic loss based estimator of, say, the log-odds $\phi(\theta) = \log[\theta/(1 - \theta)]$, is $E[\phi | D] = \psi(r + \alpha) - \psi(n - r + \beta)$ (where $\psi(x) = d \log[\Gamma(x)]/dx$ is the *digamma* function), which is *not* equal to $\phi(E[\theta | D])$. The intrinsic loss function in this problem is

$$\delta(\tilde{\theta}, \theta) = n \min\{k(\tilde{\theta} | \theta), k(\theta | \tilde{\theta})\} \quad (35)$$

where

$$k(\theta_i | \theta_j) = \theta_j \log \frac{\theta_j}{\theta_i} + (1 - \theta_j) \log \frac{1 - \theta_j}{1 - \theta_i},$$

and the corresponding intrinsic estimator θ^* is obtained by minimizing the expected posterior loss $d(\tilde{\theta} | D) = \int \delta(\tilde{\theta}, \theta) p(\theta | D) d\theta$. The exact value of θ^* may be obtained by numerical minimization, but a very good approximation is given by

$$\theta^* \approx \frac{r + \alpha/2}{n + \alpha}. \quad (36)$$

In particular, with no prior information, the reference prior is $\text{Be}(\theta | \frac{1}{2}, \frac{1}{2})$ and, hence, the objective intrinsic estimator is approximately $(r + 1/4)/(n + 1/2)$.

Since intrinsic estimation is an invariant procedure, the intrinsic estimator of the log-odds will simply be the log-odds of the intrinsic estimator of θ . As one

would expect, when $r + \alpha$ and $n - r + \beta$ are both large, all Bayes estimators of any well-behaved function $\phi(\theta)$ will cluster around $\phi(E[\theta | D])$.

Interval Estimation. To describe the inferential content of the posterior distribution of the quantity of interest $p(\theta | D)$ it is often convenient to quote regions $R \subset \Theta$ of given probability under $p(\theta | D)$. For example, the identification of regions containing 50%, 90%, 95%, or 99% of the probability under the posterior may be sufficient to convey the general quantitative messages implicit in $p(\theta | D)$; indeed, this is the intuitive basis of graphical representations of univariate distributions like those provided by boxplots. Any region $R \subset \Theta$ such that $\int_R p(\theta | D) d\theta = q$ (so that, given data D , the true value of θ belongs to R with probability q), is said to be a posterior q -credible region of θ . Notice that this provides immediately a direct intuitive statement about the unknown quantity of interest θ in probability terms, in marked contrast to the circumlocutory statements provided by frequentist confidence intervals. Clearly, for any given q there are generally infinitely many credible regions. A credible region is invariant under reparametrization; thus, for any q -credible region R of θ , $\phi(R)$ is a q -credible region of $\phi = \phi(\theta)$. Sometimes, credible regions are selected to have minimum size (length, area, volume), resulting in *highest probability density* (HPD) regions, where all points in the region have larger probability density than all points outside. However, HPD regions are *not* invariant under reparametrization: the image $\phi(R)$ of an HPD region R will be a credible region for ϕ , but will not generally be HPD; indeed, there is no compelling reason to restrict attention to HPD credible regions. Posterior quantiles are often used to derive credible regions. Thus, if $\theta_q = \theta_q(D)$ is the $100q\%$ posterior quantile of θ , then $R = \{\theta; \theta \leq \theta_q\}$ is a one-sided, typically unique q -credible region, and it is invariant under reparametrization. Indeed, *probability centred* q -credible regions of the form $R = \{\theta; \theta_{(1-q)/2} \leq \theta \leq \theta_{(1+q)/2}\}$ are easier to compute, and are often quoted in preference to HPD regions.

A better alternative is to use a loss function $L(\tilde{\theta}, \theta)$. A Lowest Posterior Loss (LPL) q -credible region R is then defined a q -credible region such that all the elements in the region have smaller posterior loss than those outside R , *i.e.*, such that

$$\int_R p(\theta | D) d\theta = q, \quad \forall \theta_1 \in R, \forall \theta_2 \notin R, \bar{L}(\theta_1 | D) \leq \bar{L}(\theta_2 | D),$$

where

$$\bar{L}(\theta_i | D) = \int_{\Theta} L(\theta, \theta_i) p(\theta | D) d\theta.$$

Intrinsic, invariant loss functions will produce *invariant* LPL credible regions. The intrinsic discrepancy loss (35) produces very attractive results.

Example 3. (Inference on normal parameters, continued). In the numerical example about the value of the gravitational field described in the upper panel

of Figure 3, the interval [9.788, 9.829] in the unrestricted posterior density of g is a HPD, 95%-credible region for g and it is also an intrinsic LDL region. The interval [9.7803, 9.8322] in the lower panel of Figure 3 is also an intrinsic LDL 95%-credible region for g , but it is not HPD.

The concept of a credible region for a function $\theta = \theta(\omega)$ of the parameter vector is trivially extended to prediction problems. Thus, a posterior q -credible region for $x \in X$ is a subset R of the outcome space X with posterior predictive probability q , so that $\int_R p(x | D) dx = q$.

4.2. Hypothesis Testing

The posterior distribution $p(\theta | D)$ of the quantity of interest θ conveys immediate intuitive information on those values of θ which, given the assumed model, may be taken to be *compatible* with the observed data D , namely, those with a relatively high probability density. Sometimes, a *restriction* $\theta \in \Theta_0 \subset \Theta$ of the possible values of the quantity of interest (where Θ_0 may possibly consists of a single value θ_0) is suggested in the course of the investigation as deserving special consideration, either because restricting θ to Θ_0 would greatly simplify the model, or because there are additional, context specific arguments suggesting that $\theta \in \Theta_0$. Intuitively, the *hypothesis* $H_0 \equiv \{\theta \in \Theta_0\}$ should be judged to be *compatible* with the observed data D if there are elements in Θ_0 with a relatively high posterior density. However, a more precise conclusion is often required and, once again, this is made possible by adopting a decision-oriented approach. Formally, testing the hypothesis $H_0 \equiv \{\theta \in \Theta_0\}$ is a *decision problem* where the action space has only two elements, namely to accept (a_0) or to reject (a_1) the proposed restriction. To solve this decision problem, it is necessary to specify an appropriate loss function, $L(a_i, \theta)$, measuring the consequences of accepting or rejecting H_0 as a function of the actual value θ of the vector of interest. Notice that this requires the statement of an *alternative* a_1 to accepting H_0 ; this is only to be expected, for an action is taken not because it is good, but because it is better than anything else that has been imagined.

Given data D , the optimal action will be to reject H_0 if (and only if) the expected posterior loss of accepting, $\int_{\Theta} L(a_0, \theta) p(\theta | D) d\theta$, is larger than the expected posterior loss of rejecting, $\int_{\Theta} L(a_1, \theta) p(\theta | D) d\theta$, that is, if (and only if)

$$\int_{\Theta} [L(a_0, \theta) - L(a_1, \theta)] p(\theta | D) d\theta = \int_{\Theta} \Delta L(\theta) p(\theta | D) d\theta > 0. \quad (37)$$

Therefore, only the loss difference $\Delta L(\theta) = L(a_0, \theta) - L(a_1, \theta)$, which measures the *advantage* of rejecting H_0 as a function of θ , has to be specified. Thus, as common sense dictates, the hypothesis H_0 should be rejected whenever the expected advantage of rejecting H_0 is positive.

A crucial element in the specification of the loss function is a description of what is actually meant by rejecting H_0 . By assumption a_0 means to act as if H_0 were true, *i.e.*, as if $\theta \in \Theta_0$, but there are at least two obvious options for the alternative action a_1 . This may either mean (i) the *negation* of H_0 , that is to act as if $\theta \notin \Theta_0$ or, alternatively, it may rather mean (ii) to reject the simplification implied by H_0 and to keep the unrestricted model, $\theta \in \Theta$, which is true by assumption. Both options have been analyzed in the literature, although it may be argued that the problems of scientific data analysis where hypothesis testing procedures are typically used are better described by the second alternative. Indeed, an established model, identified by $H_0 \equiv \{\theta \in \Theta_0\}$, is often embedded into a more general model, $\{\theta \in \Theta, \Theta_0 \subset \Theta\}$, constructed to include possibly promising departures from H_0 , and it is required to verify whether presently available data D are still compatible with $\theta \in \Theta_0$, or whether the extension to $\theta \in \Theta$ is really required.

Example 8. (Conventional hypothesis testing). Let $p(\theta|D)$, $\theta \in \Theta$, be the posterior distribution of the quantity of interest, let a_0 be the decision to work under the restriction $\theta \in \Theta_0$ and let a_1 be the decision to work under the complementary restriction $\theta \notin \Theta_0$. Suppose, moreover, that the loss structure has the simple, zero-one form given by $\{L(a_0, \theta) = 0, L(a_1, \theta) = 1\}$ if $\theta \in \Theta_0$ and, similarly, $\{L(a_0, \theta) = 1, L(a_1, \theta) = 0\}$ if $\theta \notin \Theta_0$, so that the *advantage* $\Delta L(\theta)$ of rejecting H_0 is 1 if $\theta \notin \Theta_0$ and it is -1 otherwise. With this loss function it is immediately found that the optimal action is to reject H_0 if (and only if) $\Pr(\theta \notin \Theta_0 | D) > \Pr(\theta \in \Theta_0 | D)$. Notice that this formulation requires that $\Pr(\theta \in \Theta_0) > 0$, that is, that the hypothesis H_0 has a strictly positive prior probability. If θ is a continuous parameter and Θ_0 has zero measure (for instance if H_0 consists of a single point θ_0), this requires the use of a non-regular “sharp” prior concentrating a positive probability mass on Θ_0 .

Example 9. (Intrinsic hypothesis testing). Again, let $p(\theta|D)$, $\theta \in \Theta$, be the posterior distribution of the quantity of interest, and let a_0 be the decision to work under the restriction $\theta \in \Theta_0$, but let a_1 now be the decision to keep the general, unrestricted model $\theta \in \Theta$. In this case, the advantage $\Delta L(\theta)$ of rejecting H_0 as a function of θ may safely be assumed to have the form $\Delta L(\theta) = d(\Theta_0, \theta) - d^*$, for some $d^* > 0$, where (i) $d(\Theta_0, \theta)$ is some measure of the discrepancy between the assumed model $p(D|\theta)$ and its closest approximation within the class $\{p(D|\theta_0), \theta_0 \in \Theta_0\}$, such that $d(\Theta_0, \theta) = 0$ whenever $\theta \in \Theta_0$, and (ii) d^* is a context dependent *utility constant* which measures the (necessarily positive) advantage of being able to work with the simpler model when it is true. Choices of both $d(\Theta_0, \theta)$ and d^* which may be appropriate for general use will now be described.

For reasons similar to those supporting its use in point estimation, an attractive choice for the function $d(\Theta_0, \theta)$ is an appropriate extension of the intrinsic discrepancy; when there are no nuisance parameters, this is given by

$$\delta(\Theta_0, \theta) = \inf_{\theta_0 \in \Theta_0} \min\{k(\theta_0 | \theta), k(\theta | \theta_0)\} \quad (38)$$

where $k(\theta_i | \theta_j) = \int_T p(t | \theta_j) \log\{p(t | \theta_j)/p(t | \theta_i)\} dt$, and $t = t(D) \in T$ is any sufficient statistic, which may well be the whole data set D . As before, if the data $D = \{x_1, \dots, x_n\}$ consist of a random sample from $p(x | \theta)$, then

$$k(\theta_i | \theta_j) = n \int_X p(x | \theta_j) \log \frac{p(x | \theta_j)}{p(x | \theta_i)} dx. \quad (39)$$

Naturally, the loss function $d(\Theta_0, \theta)$ reduces to the intrinsic discrepancy $\delta(\theta_0, \theta)$ of Example 6 when Θ_0 contains a single element θ_0 . Besides, as in the case of estimation, the definition is easily extended to problems with nuisance parameters, leading to

$$\delta(\Theta_0, (\theta, \lambda)) = \inf_{\theta_0 \in \Theta_0, \lambda_0 \in \Lambda} \delta\{(\theta_0, \lambda_0), (\theta, \lambda)\}. \quad (40)$$

The (null) hypothesis H_0 should be rejected if (and only if) the posterior expected advantage of rejecting is too large, *i.e.*, iff

$$d(\Theta_0 | D) = \int_{\Lambda} \int_{\Theta} \delta(\Theta_0, (\theta, \lambda)) p(\theta, \lambda | D) d\theta d\lambda > d^*, \quad (41)$$

for some $d^* > 0$. It is easily verified that the function $d(\Theta_0, D)$ is nonnegative. Moreover, if $\phi = \phi(\theta)$ is a one-to-one transformation of θ , then

$$d(\phi(\Theta_0), D) = d(\Theta_0, D),$$

so that—as one should surely require—the expected intrinsic loss of rejecting H_0 is invariant under reparametrization.

It may be shown that, as the sample size increases, the expected value of $d(\Theta_0, D)$ under sampling tends to one when H_0 is true, and tends to infinity otherwise; thus the *intrinsic statistic* $d(\Theta_0, D)$ may be regarded as a continuous, positive measure of how inappropriate (in loss of information units) it would be to simplify the model by accepting H_0 . In traditional language, $d(\Theta_0, D)$ is a *test statistic* for H_0 and the hypothesis should be rejected if the value of $d(\Theta_0, D)$ exceeds some *critical value* d^* . However, in sharp contrast to conventional hypothesis testing, this critical value d^* is a context specific, positive utility constant d^* , which

may precisely be described as the number of *information units* which the decision maker is prepared to lose in order to be able to work with the simpler model H_0 , and does not depend on the sampling properties of the probability model. The procedure may be used with standard, continuous regular priors even in *sharp* hypothesis testing, when Θ_0 is a zero-measure set (as would be the case if θ is continuous and Θ_0 contains a single point θ_0). Naturally, to implement the test, the utility constant d^* which defines the rejection region must be chosen.

It is easily verified from its definition that the intrinsic discrepancy between two probability models is the minimum expected value of their log-likelihood ration. Thus using, say $d^* = \log(100) \approx 4.6$ as the required threshold, is to reject a null hypothesis when the data may be expected to be about 100 times more likely under a model labelled by a value of θ which does *not* belong to Θ_0 than under a model labelled by the best estimate of θ .

All measurements are based on a comparison with a standard; comparison with the “canonical” problem of testing a value $\mu = \mu_0$ for the mean of a normal distribution with known variance (see below) makes it possible to *calibrate* this *information scale* with respect to a well known situation. Values of $d(\Theta_0, D)$ of about 1 should be regarded as an indication of no evidence against H_0 , since the expected value of $d(\Theta_0, D)$ under H_0 is exactly equal to one. Values of $f(\Theta_0, D)$ of about $2.5 \approx \log(12)$, and $5 \approx \log(148)$ should be respectively regarded as an indication of mild evidence against H_0 , and significant evidence against H_0 since (see the details below), in the canonical normal problem, these values correspond to the observed sample mean \bar{x} respectively lying 2 or 3 posterior standard deviations from the null value μ_0 , and the data would respectively be 12 and 148 times more likely to have been observed from a value $\mu \neq \mu_0$. Notice that, in sharp contrast to frequentist hypothesis testing, where it is hazily recommended to adjust the significance level for dimensionality and sample size, this provides an absolute scale (in information units) which remains valid for any sample size and any dimensionality.

Example 10. (Testing the value of a normal mean). Let the data $D = \{x_1, \dots, x_n\}$ be a random sample from a normal distribution $N(x | \mu, \sigma)$, where σ is assumed to be known, and consider the “canonical” problem of testing whether these data are or are not compatible with some specific sharp hypothesis $H_0 \equiv \{\mu = \mu_0\}$ on the value of the mean.

The conventional approach to this problem requires a non-regular prior which places a probability mass, say p_0 , on the value μ_0 to be tested, with the remaining $1 - p_0$ probability continuously distributed over \mathfrak{R} . If this prior is chosen to be $p(\mu | \mu \neq \mu_0) = N(\mu | \mu_0, \sigma_0)$, Bayes theorem may be used to obtain the corresponding posterior probability,

$$\Pr[\mu_0 | D, \lambda] = \frac{B_{01}(D, \lambda) p_0}{(1 - p_0) + p_0 B_{01}(D, \lambda)}, \quad (42)$$

$$B_{01}(D, \lambda) = \left(1 + \frac{n}{\lambda}\right)^{1/2} \exp\left[-\frac{1}{2} \frac{n}{n + \lambda} z^2\right], \quad (43)$$

where $z = (\bar{x} - \mu_0)/(\sigma/\sqrt{n})$ measures, in standard deviations, the distance between \bar{x} and μ_0 and $\lambda = \sigma^2/\sigma_0^2$ is the ratio of model to prior variance. The function $B_{01}(D, \lambda)$, a ratio of (integrated) likelihood functions, is called the *Bayes factor* in favour of H_0 . With a conventional zero-one loss function, H_0 should be rejected if $\Pr[\mu_0 | D, \lambda] < 1/2$. The choices $p_0 = 1/2$ and $\lambda = 1$ or $\lambda = 1/2$, describing particular forms of *sharp* prior knowledge, have been suggested in the literature for routine use. The conventional approach to sharp hypothesis testing deals with situations of *concentrated* prior probability; it *assumes* important prior knowledge about the value of μ and, hence, should *not* be used unless this is an appropriate assumption. Moreover, as pointed out in the 1950's by Bartlett, the resulting posterior probability is extremely sensitive to the specific prior specification. In most applications, H_0 is really a hazily defined small region rather than a point. For moderate sample sizes, the posterior probability $\Pr[\mu_0 | D, \lambda]$ is an *approximation* to the posterior probability $\Pr[\mu_0 - \epsilon < \mu < \mu_0 + \epsilon | D, \lambda]$ for some small interval around μ_0 which would have been obtained from a regular, continuous prior heavily concentrated around μ_0 ; however, this approximation *always* breaks down for sufficiently large sample sizes. One consequence (which is immediately apparent from the last two equations) is that for any *fixed* value of the pertinent statistic z , the posterior probability of the null, $\Pr[\mu_0 | D, \lambda]$, tends to one as $n \rightarrow \infty$. Far from being specific to this example, this unappealing behaviour of posterior probabilities based on sharp, non-regular priors (discovered by Lindley in the 1950's, and generally known as *Lindley's paradox*) is *always* present in the conventional Bayesian approach to *sharp* hypothesis testing.

The intrinsic approach may be used without assuming any sharp prior knowledge. The intrinsic discrepancy is $\delta(\mu_0, \mu) = n(\mu - \mu_0)^2/(2\sigma^2)$, a simple transformation of the standardized distance between μ and μ_0 . As later explained (Section 5), absence of initial information about the value of μ may formally be described in this problem by the (improper) uniform prior function $p(\mu) = 1$; Bayes' theorem may then be used to obtain the corresponding (proper) posterior distribution, $p(\mu | D) = N(\mu | \bar{x}, \sigma/\sqrt{n})$. The expected value of $d(\mu_0, \mu)$ with respect to this posterior is $d(\mu_0, D) = (1 + z^2)/2$, where $z = (\bar{x} - \mu_0)/(\sigma/\sqrt{n})$ is the standardized distance between \bar{x} and μ_0 . As foretold by the general theory, the expected value of $d(\mu_0, D)$ under repeated sampling is one if $\mu = \mu_0$, and increases linearly with n if $\mu \neq \mu_0$. Moreover, in this canonical example, to reject H_0 whenever $|z| > 2$ or $|z| > 3$, that is whenever μ_0 is 2 or 3 posterior standard deviations away from \bar{x} ,

respectively corresponds to rejecting H_0 whenever $d(\mu_0, D)$ is larger than 2.5, or larger than 5. But the information scale is independent of the problem, so that rejecting the null whenever its expected discrepancy from the true model is larger than $d^* = 5$ units of information is a *general* rule (and one which corresponds to the conventional ‘ 3σ ’ rule in the canonical normal case).

If σ is unknown, the intrinsic discrepancy becomes

$$d(\mu_0, \mu, \sigma) = \frac{n}{2} \log \left[1 + \left(\frac{\mu - \mu_0}{\sigma} \right)^2 \right]. \quad (44)$$

Moreover, as mentioned before, absence of initial information about both μ and σ may be described by the (improper) prior function $p(\mu, \sigma) = \sigma^{-1}$. The intrinsic test statistic $d(\mu_0, D)$ is found as the expected value of $d(\mu_0, \mu, \sigma)$ under the corresponding joint posterior distribution; this may be exactly expressed in terms of hypergeometric functions, and is approximated by

$$d(\mu_0, D) \approx \frac{1}{2} + \frac{n}{2} \log \left(1 + \frac{t^2}{n} \right), \quad (45)$$

where t is the traditional statistic $t = \sqrt{n-1}(\bar{x} - \mu_0)/s$, $ns^2 = \sum_j (x_j - \bar{x})^2$. For instance, for samples sizes 5, 30 and 1000, and using the utility constant $d^* = 5$, the hypothesis H_0 would be rejected whenever $|t|$ is respectively larger than 5.025, 3.240, and 3.007.

5. Reference Analysis

Under the Bayesian paradigm, the outcome of any inference problem (the posterior distribution of the quantity of interest) combines the information provided by the data with relevant available prior information. In many situations, however, either the available prior information on the quantity of interest is too vague to warrant the effort required to have it formalized in the form of a probability distribution, or it is too subjective to be useful in scientific communication or public decision making. It is therefore important to be able to identify the mathematical form of a “noninformative” prior, a prior that would have a minimal effect, relative to the data, on the posterior inference. More formally, suppose that the probability mechanism which has generated the available data D is assumed to be $p(D | \omega)$, for some $\omega \in \Omega$, and that the quantity of interest is some real-valued function $\theta = \theta(\omega)$ of the model parameter ω . Without loss of generality, it may be assumed that the probability model is of the form $p(D | \theta, \lambda)$, $\theta \in \Theta$, $\lambda \in \Lambda$, where λ is some appropriately chosen nuisance parameter vector. As described in Section 3, to obtain the required posterior distribution of the quantity of interest $p(\theta | D)$ it is necessary to specify a *joint* prior $p(\theta, \lambda)$. It is now required to identify the form of

that joint prior $\pi_\theta(\theta, \boldsymbol{\lambda})$, the θ -reference prior, which would have a *minimal effect* on the corresponding posterior distribution of θ ,

$$\pi(\theta | D) \propto \int_{\Lambda} p(D | \theta, \boldsymbol{\lambda}) \pi_\theta(\theta, \boldsymbol{\lambda}) d\boldsymbol{\lambda}, \quad (46)$$

a prior which, to use a conventional expression, “would let the data speak for themselves” about the likely value of θ . Properly defined, reference *posterior* distributions have an important role to play in scientific communication, for they provide the answer to a central question in the sciences: conditional on the assumed model $p(D | \theta, \boldsymbol{\lambda})$, and on any further assumptions of the value of θ on which there might be universal agreement, the reference posterior $\pi(\theta | D)$ should specify what *could* be said about θ if the only available information about θ were some well-documented data D .

Much work has been done to formulate “reference” priors which would make the idea described above mathematically precise. This section concentrates on an approach that is based on information theory to derive reference distributions which may be argued to provide the most advanced general procedure available. In the formulation described below, far from ignoring prior knowledge, the reference posterior exploits certain well-defined features of a *possible* prior, namely those describing a situation where relevant knowledge about the quantity of interest (beyond that universally accepted) may be held to be negligible compared to the information about that quantity which repeated experimentation (from a particular data generating mechanism) might possibly provide. Reference analysis is appropriate in contexts where the set of inferences which could be drawn in this *possible* situation is considered to be pertinent.

Any statistical analysis contains a fair number of subjective elements; these include (among others) the data selected, the model assumptions, and the choice of the quantities of interest. Reference analysis may be argued to provide an “objective” Bayesian solution to statistical inference problems in just the same sense that conventional statistical methods claim to be “objective”: in that the solutions only depend on model assumptions and observed data. The whole topic of objective Bayesian methods is, however, subject to polemic; interested readers will find in the bibliography some pointers to the relevant literature.

5.1. Reference Distributions

One parameter. Consider the experiment which consists of the observation of data D , generated by a random mechanism $p(D | \theta)$ which only depends on a real-valued parameter $\theta \in \Theta$, and let $\boldsymbol{t} = \boldsymbol{t}(D) \in T$ be *any* sufficient statistic (which may well be the complete data set D). In Shannon’s general information theory, the *amount of information* $I^\theta\{T, p(\theta)\}$ which may be expected to be provided by D , or (equivalently) by $\boldsymbol{t}(D)$, about the value of θ is defined by

$$I^\theta\{T, p(\theta)\} = \int_T \int_\Theta p(\mathbf{t}, \theta) \log \frac{p(\mathbf{t}, \theta)}{p(\mathbf{t})p(\theta)} d\theta d\mathbf{t}, \quad (47)$$

which is equal to

$$= \mathbf{E}_\mathbf{t} \left[\int_\Theta p(\theta | \mathbf{t}) \log \frac{p(\theta | \mathbf{t})}{p(\theta)} d\theta \right],$$

the expected logarithmic divergence of the prior from the posterior. This is naturally a *functional* of the prior $p(\theta)$: the larger the prior information, the smaller the information which the data may be expected to provide. The functional $I^\theta\{T, p(\theta)\}$ is concave, non-negative, and invariant under one-to-one transformations of θ . Consider now the amount of information $I^\theta\{T^k, p(\theta)\}$ about θ which may be expected from the experiment which consists of k conditionally independent replications $\{\mathbf{t}_1, \dots, \mathbf{t}_k\}$ of the original experiment. As $k \rightarrow \infty$, such an experiment would provide any *missing information* about θ which could possibly be obtained within this framework; thus, as $k \rightarrow \infty$, the functional $I^\theta\{T^k, p(\theta)\}$ will approach the missing information about θ associated with the prior $p(\theta)$. Intuitively, a θ -“noninformative” prior is one which *maximizes the missing information* about θ . Formally, if $\pi_k(\theta)$ denotes the prior density which maximizes $I^\theta\{T^k, p(\theta)\}$ in the class \mathcal{P} of strictly positive prior distributions which are compatible with accepted assumptions on the value of θ (which may well be the class of *all* strictly positive proper priors) then the θ -reference prior $\pi(\theta)$ is the limit as $k \rightarrow \infty$ (in a sense to be made precise) of the sequence of priors $\{\pi_k(\theta), k = 1, 2, \dots\}$.

Notice that this limiting procedure is *not* some kind of asymptotic approximation, but an essential element of the *definition* of a reference prior. In particular, this definition implies that reference distributions only depend on the *asymptotic* behaviour of the assumed probability model, a feature which greatly simplifies their actual derivation.

Example 11. (Maximum entropy). If θ may only take a *finite* number of values, so that the parameter space is $\Theta = \{\theta_1, \dots, \theta_m\}$ and $p(\theta) = \{p_1, \dots, p_m\}$, with $p_i = \Pr(\theta = \theta_i)$, then the missing information associated to $\{p_1, \dots, p_m\}$ may be shown to be

$$\lim_{k \rightarrow \infty} I^\theta\{T^k, p(\theta)\} = H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log(p_i), \quad (48)$$

that is, the *entropy* of the prior distribution $\{p_1, \dots, p_m\}$.

Thus, in the finite case, the reference prior is that with *maximum entropy* in the class \mathcal{P} of priors compatible with accepted assumptions. Consequently, the reference prior algorithm contains “maximum entropy” priors as the particular case which obtains when the parameter space is *finite*, the *only* case where the original concept of entropy (in statistical mechanics, as a measure of uncertainty)

is unambiguous and well-behaved. If, in particular, \mathcal{P} contains *all* priors over $\{\theta_1, \dots, \theta_m\}$, then the reference prior is $\pi(\theta) = \{1/m, \dots, 1/m\}$, the uniform prior (proposed by Laplace under the principle of insufficient reason argument).

Formally, the *reference prior function* $\pi(\theta)$ of a univariate parameter θ is defined to be the limit of the sequence of the proper priors $\pi_k(\theta)$ which maximize $I^\theta\{T^k, p(\theta)\}$ in the precise sense that, for any value of the sufficient statistic $\mathbf{t} = \mathbf{t}(D)$, the *reference posterior*, the pointwise limit $\pi(\theta | \mathbf{t})$ of the corresponding sequence of posteriors $\{\pi_k(\theta | \mathbf{t})\}$, may be obtained from $\pi(\theta)$ by formal use of Bayes theorem, so that $\pi(\theta | \mathbf{t}) \propto p(\mathbf{t} | \theta) \pi(\theta)$.

Reference prior *functions* are often simply called reference priors, even though they are usually *not* probability distributions. They should *not* be considered as expressions of belief, but technical devices to obtain (proper) posterior distributions which are a limiting form of the posteriors which could have been obtained from possible prior beliefs which were relatively uninformative with respect to the quantity of interest when compared with the information which data could provide.

If (i) the sufficient statistic $\mathbf{t} = \mathbf{t}(D)$ is a consistent estimator $\tilde{\theta}$ of a continuous parameter θ , and (ii) the class \mathcal{P} contains *all* strictly positive priors, then the reference prior may be shown to have a simple form in terms of any *asymptotic* approximation to the posterior distribution of θ . Notice that, by construction, an *asymptotic* approximation to the posterior does *not* depend on the prior. Specifically, if the posterior density $p(\theta | D)$ has an asymptotic approximation of the form $p(\theta | \tilde{\theta}, n)$, the reference prior is simply

$$\pi(\theta) \propto p(\theta | \tilde{\theta}, n)|_{\tilde{\theta}=\theta} \quad (49)$$

One-parameter reference priors are shown to be *invariant* under reparametrization; thus, if $\psi = \psi(\theta)$ is a piecewise one-to-one function of θ , then the ψ -reference prior is simply the appropriate probability transformation of the θ -reference prior.

Example 12. (Jeffreys' prior). If θ is univariate and continuous, and the posterior distribution of θ given $\{x_1, \dots, x_n\}$ is asymptotically normal with standard deviation $s(\tilde{\theta})/\sqrt{n}$, then, using (49), the reference prior function is $\pi(\theta) \propto s(\theta)^{-1}$. Under regularity conditions (often satisfied in practice, see Section 3.3), the posterior distribution of θ is asymptotically normal with variance $n^{-1} F^{-1}(\hat{\theta})$, where $F(\theta)$ is Fisher's information function and $\hat{\theta}$ is the MLE of θ . Hence, the reference prior function in these conditions is $\pi(\theta) \propto F(\theta)^{1/2}$, which is known as Jeffreys' prior. It follows that the reference prior algorithm contains Jeffreys' priors as the particular case which obtains when the probability model only depends on a single continuous univariate parameter, there are regularity conditions to guarantee asymptotic normality, and there is no additional information, so that the class of possible priors \mathcal{P} contains all strictly positive priors over Θ . These are precisely the

conditions under which there is general agreement on the use of Jeffreys' prior as a "noninformative" prior.

Example 2. (Inference on a binomial parameter, continued). Let data $D = \{x_1, \dots, x_n\}$ consist of a sequence of n independent Bernoulli trials, so that

$$p(x|\theta) = \theta^x(1-\theta)^{1-x}, \quad x \in \{0, 1\};$$

this is a regular, one-parameter continuous model, whose Fisher's function is $F(\theta) = \theta^{-1}(1-\theta)^{-1}$. Thus, the reference prior $\pi(\theta)$ is proportional to $\theta^{-1/2}(1-\theta)^{-1/2}$, so that the reference prior is the (proper) Beta distribution $\text{Be}(\theta | 1/2, 1/2)$. Since the reference algorithm is invariant under reparametrization, the reference prior of $\phi(\theta) = 2 \arcsin \sqrt{\theta}$ is $\pi(\phi) = \pi(\theta)/|\partial\phi/\partial\theta| = 1$; thus, the reference prior is uniform on the variance-stabilizing transformation $\phi(\theta) = 2 \arcsin \sqrt{\theta}$, a feature generally true under regularity conditions. In terms of the original parameter θ , the corresponding reference posterior is $\text{Be}(\theta | r + 1/2, n - r + 1/2)$, where $r = \sum x_j$ is the number of positive trials.

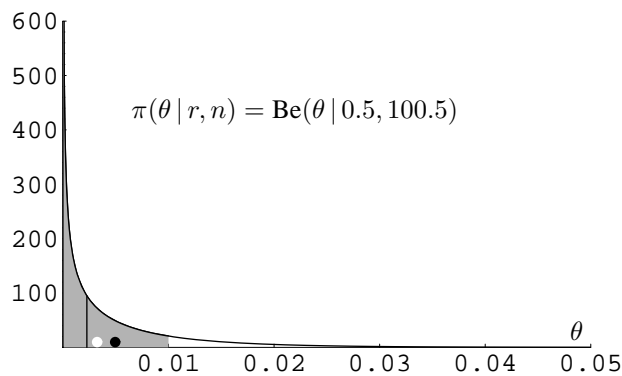


Figure 4. Posterior distribution of the proportion of infected people in the population, given the results of $n = 100$ tests, none of which were positive.

Suppose, for example, that $n = 100$ randomly selected people have been tested for an infection and that all tested negative, so that $r = 0$. The reference posterior distribution of the proportion θ of people infected is then the Beta distribution $\text{Be}(\theta | 0.5, 100.5)$, represented in Figure 4. It may well be known that the infection was rare, leading to the assumption that $\theta < \theta_0$, for some upper bound θ_0 ; the (restricted) reference prior would then be of the form $\pi(\theta) \propto \theta^{-1/2}(1-\theta)^{-1/2}$ if $\theta < \theta_0$, and zero otherwise. However, provided the likelihood is concentrated

in the region $\theta < \theta_0$, the corresponding posterior would virtually be identical to $\text{Be}(\theta | 0.5, 100.5)$. Thus, just on the basis of the observed experimental results, one may claim that the proportion of infected people is surely smaller than 5% (for the reference posterior probability of the event $\theta > 0.05$ is 0.001), that θ is smaller than 0.01 with probability 0.844 (area of the shaded region in Figure 4), that it is equally likely to be over or below 0.23% (for the median, represented by a vertical line, is 0.0023), and that the probability that a person randomly chosen from the population is infected is 0.005 (the posterior mean, represented in the figure by a black circle), since $\Pr(x = 1 | r, n) = E[\theta | r, n] = 0.005$. If a particular point estimate of θ is required (say a number to be quoted in the summary headline) the *intrinsic* estimator suggests itself; this is found to be $\theta^* = 0.0032$ (represented in the figure with a white circle). Notice that the traditional solution to this problem, based on the asymptotic behaviour of the MLE, here $\hat{\theta} = r/n = 0$, for any n , makes absolutely no sense in this scenario.

One nuisance parameter. The extension of the reference prior algorithm to the case of two parameters follows the usual mathematical procedure of reducing the problem to a sequential application of the established procedure for the single parameter case. Thus, if the probability model is $p(\mathbf{t} | \theta, \lambda)$, $\theta \in \Theta$, $\lambda \in \Lambda$ and a θ -reference prior $\pi_\theta(\theta, \lambda)$ is required, the reference algorithm proceeds in two steps:

(i) Conditional on θ , $p(\mathbf{t} | \theta, \lambda)$ only depends on the nuisance parameter λ and, hence, the one-parameter algorithm may be used to obtain the *conditional* reference prior $\pi(\lambda | \theta)$.

(ii) If $\pi(\lambda | \theta)$ is proper, this may be used to integrate out the nuisance parameter thus obtaining the one-parameter integrated model

$$p(\mathbf{t} | \theta) = \int_{\Lambda} p(\mathbf{t} | \theta, \lambda) \pi(\lambda | \theta) d\lambda,$$

to which the one-parameter algorithm may be applied again to obtain $\pi(\theta)$. The θ -reference prior is then $\pi_\theta(\theta, \lambda) = \pi(\lambda | \theta) \pi(\theta)$, and the required reference posterior is $\pi(\theta | \mathbf{t}) \propto p(\mathbf{t} | \theta) \pi(\theta)$.

If the conditional reference prior is *not* proper, then the procedure is performed within an increasing sequence $\{\Lambda_i\}$ of subsets converging to Λ over which $\pi(\lambda | \theta)$ is integrable. This makes it possible to obtain a corresponding sequence of θ -reference posteriors $\{\pi_i(\theta | \mathbf{t})\}$ for the quantity of interest θ , and the required reference posterior is the corresponding pointwise limit $\pi(\theta | \mathbf{t}) = \lim_i \pi_i(\theta | \mathbf{t})$. A θ -reference prior is then defined as a positive function $\pi_\theta(\theta, \lambda)$ which may be formally used in Bayes' theorem as a prior to obtain the reference posterior, *i.e.*, such that, for any $\mathbf{t} \in T$, $\pi(\theta | \mathbf{t}) \propto \int_{\Lambda} p(\mathbf{t} | \theta, \lambda) \pi_\theta(\theta, \lambda) d\lambda$. The approximating sequences should be *consistently* chosen within a given model. Thus, given a probability model $\{p(\mathbf{x} | \omega), \omega \in \Omega\}$ an appropriate approximating sequence $\{\Omega_i\}$ should be

chosen for the whole parameter space Ω . Thus, if the analysis is done in terms of, say, $\psi = \{\psi_1, \psi_2\} \in \Psi(\Omega)$, the approximating sequence should be chosen such that $\Psi_i = \psi(\Omega_i)$. A natural approximating sequence in location-scale problems is $\{\mu, \log \sigma\} \in [-i, i]^2$.

The θ -reference prior does *not* depend on the choice of the nuisance parameter λ ; thus, for any $\psi = \psi(\theta, \lambda)$ such that (θ, ψ) is a one-to-one function of (θ, λ) , the θ -reference prior in terms of (θ, ψ) is simply $\pi_\theta(\theta, \psi) = \pi_\theta(\theta, \lambda) / |\partial(\theta, \psi) / \partial(\theta, \lambda)|$, the appropriate probability transformation of the θ -reference prior in terms of (θ, λ) . Notice, however, that the reference prior *may* depend on the parameter of interest; thus, the θ -reference prior may differ from the ϕ -reference prior unless either ϕ is a piecewise one-to-one transformation of θ , or ϕ is asymptotically independent of θ . This is an expected consequence of the fact that the conditions under which the missing information about θ is maximized are not generally the same as the conditions which maximize the missing information about some function $\phi = \phi(\theta, \lambda)$.

The *non-existence* of a unique “noninformative prior” which would be appropriate for any inference problem within a given model was established in the 1970’s by Dawid, Stone and Zidek, when they showed that this is incompatible with *consistent marginalization*. Indeed, if given the model $p(D | \theta, \lambda)$, the reference posterior of the quantity of interest θ , $\pi(\theta | D) = \pi(\theta | \mathbf{t})$, only depends on the data through a statistic \mathbf{t} whose sampling distribution, $p(\mathbf{t} | \theta, \lambda) = p(\mathbf{t} | \theta)$, only depends on θ , one would expect the reference posterior to be of the form $\pi(\theta | \mathbf{t}) \propto \pi(\theta) p(\mathbf{t} | \theta)$ for some prior $\pi(\theta)$. However, examples were found where this cannot be the case if a *unique* joint “noninformative” prior were to be used for all possible quantities of interest.

Example 13. (Regular two dimensional continuous reference prior functions). If the joint posterior distribution of (θ, λ) is asymptotically normal, then the θ -reference prior may be derived in terms of the corresponding Fisher’s information matrix, $\mathbf{F}(\theta, \lambda)$. Indeed, if

$$\mathbf{F}(\theta, \lambda) = \begin{pmatrix} F_{\theta\theta}(\theta, \lambda) & F_{\theta\lambda}(\theta, \lambda) \\ F_{\theta\lambda}(\theta, \lambda) & F_{\lambda\lambda}(\theta, \lambda) \end{pmatrix}, \quad \text{and} \quad \mathbf{S}(\theta, \lambda) = \mathbf{F}^{-1}(\theta, \lambda), \quad (50)$$

then the θ -reference prior is $\pi_\theta(\theta, \lambda) = \pi(\lambda | \theta) \pi(\theta)$, where

$$\pi(\lambda | \theta) \propto F_{\lambda\lambda}^{-1/2}(\theta, \lambda), \quad \lambda \in \Lambda. \quad (51)$$

If $\pi(\lambda | \theta)$ is proper,

$$\pi(\theta) \propto \exp \left\{ \int_{\Lambda} \pi(\lambda | \theta) \log[S_{\theta\theta}^{-1/2}(\theta, \lambda)] d\lambda \right\}, \quad \theta \in \Theta. \quad (52)$$

If $\pi(\lambda | \theta)$ is not proper, integrations are performed on an approximating sequence $\{\Lambda_i\}$ to obtain a sequence $\{\pi_i(\lambda | \theta) \pi_i(\theta)\}$, (where $\pi_i(\lambda | \theta)$ is the proper renormalization of $\pi(\lambda | \theta)$ to Λ_i) and the θ -reference prior $\pi_\theta(\theta, \lambda)$ is defined as its appropriate limit. Moreover, if (i) both $F_{\lambda\lambda}^{1/2}(\theta, \lambda)$ and $S_{\theta\theta}^{-1/2}(\theta, \lambda)$ factorize, so that

$$S_{\theta\theta}^{-1/2}(\theta, \lambda) \propto f_\theta(\theta) g_\theta(\lambda), \quad F_{\lambda\lambda}^{1/2}(\theta, \lambda) \propto f_\lambda(\theta) g_\lambda(\lambda), \quad (53)$$

and (ii) the parameters θ and λ are *variation independent*, so that Λ does not depend on θ , then the θ -reference prior is simply $\pi_\theta(\theta, \lambda) = f_\theta(\theta) g_\lambda(\lambda)$, even if the conditional reference prior $\pi(\lambda | \theta) = \pi(\lambda) \propto g_\lambda(\lambda)$ (which will not depend on θ) is actually improper.

Example 3. (Inference on normal parameters, continued). The information matrix which corresponds to a normal model $N(x | \mu, \sigma)$ is

$$\mathbf{F}(\mu, \sigma) = \begin{pmatrix} \sigma^{-2} & 0 \\ 0 & 2\sigma^{-2} \end{pmatrix}, \quad \mathbf{S}(\mu, \sigma) = \mathbf{F}^{-1}(\mu, \sigma) = \begin{pmatrix} \sigma^2 & 0 \\ 0 & \frac{1}{2}\sigma^2 \end{pmatrix}; \quad (54)$$

hence $F_{\sigma\sigma}^{1/2}(\mu, \sigma) = \sqrt{2}\sigma^{-1} = f_\sigma(\mu) g_\sigma(\sigma)$, with $g_\sigma(\sigma) = \sigma^{-1}$; thus the conditional reference prior of σ is $\pi(\sigma | \mu) = \sigma^{-1}$. Similarly, $S_{\mu\mu}^{-1/2}(\mu, \sigma) = \sigma^{-1}$, which is of the form $f_\mu(\mu) g_\mu(\sigma)$, with $f_\mu(\mu) = 1$, and thus the marginal reference prior of μ is $\pi(\mu) = 1$. Therefore, the μ -reference prior is $\pi_\mu(\mu, \sigma) = \pi(\sigma | \mu) \pi(\mu) = \sigma^{-1}$, as already anticipated. Moreover, as one would expect from the fact that $\mathbf{F}(\mu, \sigma)$ is diagonal and also anticipated, it is similarly found that the σ -reference prior is $\pi_\sigma(\mu, \sigma) = \sigma^{-1}$, the same as $\pi_\mu(\mu, \sigma)$.

Suppose, however, that the quantity of interest is *not* the mean μ or the standard deviation σ , but the *standardized* mean $\phi = \mu/\sigma$. Fisher's information matrix in terms of the parameters ϕ and σ is $\mathbf{F}(\phi, \sigma) = J^t \mathbf{F}(\mu, \sigma) J$, where $J = (\partial(\mu, \sigma)/\partial(\phi, \sigma))$ is the Jacobian of the inverse transformation; this yields

$$\mathbf{F}(\phi, \sigma) = \begin{pmatrix} 1 & \phi\sigma^{-1} \\ \phi\sigma^{-1} & \sigma^{-2}(2 + \phi^2) \end{pmatrix}, \quad \mathbf{S}(\phi, \sigma) = \begin{pmatrix} 1 + \frac{1}{2}\phi^2 & -\frac{1}{2}\phi\sigma \\ -\frac{1}{2}\phi\sigma & \frac{1}{2}\sigma^2 \end{pmatrix}. \quad (55)$$

Thus, $S_{\phi\phi}^{-1/2}(\phi, \sigma) \propto (1 + \frac{1}{2}\phi^2)^{-1/2}$ and $F_{\sigma\sigma}^{1/2}(\phi, \sigma) \propto \sigma^{-1}(2 + \phi^2)^{1/2}$. Hence, using again the results in Example 13, $\pi_\phi(\phi, \sigma) = (1 + \frac{1}{2}\phi^2)^{-1/2}\sigma^{-1}$. In the original parametrization, this is $\pi_\phi(\mu, \sigma) = (1 + \frac{1}{2}(\mu/\sigma)^2)^{-1/2}\sigma^{-2}$, which is *very* different from $\pi_\mu(\mu, \sigma) = \pi_\sigma(\mu, \sigma) = \sigma^{-1}$. The corresponding reference posterior of ϕ is $\pi(\phi | x_1, \dots, x_n) \propto (1 + \frac{1}{2}\phi^2)^{-1/2} p(t | \phi)$ where $t = (\sum x_j)/(\sum x_j^2)^{1/2}$, a one-dimensional (marginally sufficient) statistic whose sampling distribution,

$p(t | \mu, \sigma) = p(t | \phi)$, only depends on ϕ . Thus, the reference prior algorithm is seen to be consistent under marginalization.

Many parameters. The reference algorithm is easily generalized to an arbitrary number of parameters. If the model is $p(\mathbf{t} | \omega_1, \dots, \omega_m)$, a joint reference prior

$$\pi(\theta_m | \theta_{m-1}, \dots, \theta_1) \times \dots \times \pi(\theta_2 | \theta_1) \times \pi(\theta_1) \quad (56)$$

may sequentially be obtained for each *ordered* parametrization $\{\theta_1(\omega), \dots, \theta_m(\omega)\}$ of interest, and these are invariant under reparametrization of any of the $\theta_i(\omega)$'s. The choice of the ordered parametrization $\{\theta_1, \dots, \theta_m\}$ precisely describes the particular prior required, namely that which *sequentially* maximizes the missing information about each of the θ_i 's, conditional on $\{\theta_1, \dots, \theta_{i-1}\}$, for $i = m, m-1, \dots, 1$.

Example 14. (Stein's paradox). Let D be a random sample from a m -variate normal distribution with mean $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_m\}$ and unitary variance matrix. The reference prior which corresponds to any permutation of the μ_i 's is uniform, and this prior leads indeed to appropriate reference posterior distributions for any of the μ_i 's, namely $\pi(\mu_i | D) = N(\mu_i | \bar{x}_i, 1/\sqrt{n})$. Suppose, however, that the quantity of interest is $\theta = \sum_i \mu_i^2$, the distance of $\boldsymbol{\mu}$ to the origin. As showed by Stein in the 1950's, the posterior distribution of θ based on that uniform prior (or in any "flat" *proper* approximation) has very undesirable properties; this is due to the fact that a uniform (or nearly uniform) prior, although "noninformative" with respect to each of the individual μ_i 's, is actually highly informative on the sum of their squares, introducing a severe positive bias (Stein's paradox). However, the reference prior which corresponds to a parametrization of the form $\{\theta, \lambda_1, \dots, \lambda_{m-1}\}$ produces, for any choice of the nuisance parameters $\lambda_i = \lambda_i(\boldsymbol{\mu})$, the reference posterior $\pi(\theta | D) = \pi(\theta | t) \propto \theta^{-1/2} \chi^2(nt | m, n\theta)$, where $t = \sum_i \bar{x}_i^2$, and this posterior is shown to have the appropriate consistency properties.

Far from being specific to Stein's example, the inappropriate behaviour in problems with many parameters of specific marginal posterior distributions derived from multivariate "flat" priors (proper or improper) is indeed very frequent. Hence, sloppy, uncontrolled use of "flat" priors (rather than the relevant reference priors), is strongly discouraged.

Limited information. Although often used in contexts where no universally agreed prior knowledge about the quantity of interest is available, the reference algorithm may be used to specify a prior which incorporates any acceptable prior knowledge; it suffices to maximize the missing information within the class \mathcal{P} of priors which is compatible with such accepted knowledge. Indeed, by progressive incorporation of further restrictions into \mathcal{P} , the reference prior algorithm becomes a method of (prior) *probability assessment*. As described below, the problem has a fairly simple analytical solution when those restrictions take the form of known expected

values. The incorporation of other type of restrictions usually involves numerical computations.

Example 15. (Univariate restricted reference priors). If the probability mechanism which is assumed to have generated the available data only depends on a univariate continuous parameter $\theta \in \Theta \subset \mathfrak{R}$, and the class \mathcal{P} of acceptable priors is a class of proper priors which satisfies some expected value restrictions, so that

$$\mathcal{P} = \left\{ p(\theta); \quad p(\theta) > 0, \int_{\Theta} p(\theta) d\theta = 1 \right\} \quad (57)$$

under the conditions

$$\int_{\Theta} p(\theta) d\theta = 1 \int_{\Theta} g_i(\theta) p(\theta) d\theta = \beta_i, \quad i = 1, \dots, m,$$

then the (restricted) reference prior is

$$\pi(\theta | \mathcal{P}) \propto \pi(\theta) \exp \left[\sum_{j=1}^m \gamma_j g_j(\theta) \right] \quad (58)$$

where $\pi(\theta)$ is the unrestricted reference prior and the γ_i 's are constants (the corresponding Lagrange multipliers), to be determined by the restrictions which define \mathcal{P} . Suppose, for instance, that data are considered to be a random sample from a location model centred at θ , and that it is further assumed that $E[\theta] = \mu_0$ and that $\text{Var}[\theta] = \sigma_0^2$. The unrestricted reference prior for any regular location problem may be shown to be uniform. Thus, the restricted reference prior must be of the form $\pi(\theta | \mathcal{P}) \propto \exp\{\gamma_1 \theta + \gamma_2 (\theta - \mu_0)^2\}$, with $\int_{\Theta} \theta \pi(\theta | \mathcal{P}) d\theta = \mu_0$ and $\int_{\Theta} (\theta - \mu_0)^2 \pi(\theta | \mathcal{P}) d\theta = \sigma_0^2$. Hence, $\pi(\theta | \mathcal{P})$ is a *normal* distribution with the specified mean and variance.

5.2. Frequentist Properties

Bayesian methods provide a *direct* solution to the problems typically posed in statistical inference; indeed, posterior distributions precisely state what can be said about unknown quantities of interest *given* available data and prior knowledge. In particular, unrestricted reference posterior distributions state what could be said if no prior knowledge about the quantities of interest were available.

A frequentist analysis of the behaviour of Bayesian procedures under repeated sampling may, however, be illuminating, for this provides some interesting connections between frequentist and Bayesian inference. It is found that the frequentist properties of Bayesian reference procedures are typically excellent, and may be used to provide a form of calibration for reference posterior probabilities.

Point Estimation. It is generally accepted that, as the sample size increases, a “good” estimator $\tilde{\theta}$ of θ ought to get the correct value of θ eventually, that is to be *consistent*. Under appropriate regularity conditions, any Bayes estimator ϕ^* of any function $\phi(\theta)$ converges in probability to $\phi(\theta)$, so that sequences of Bayes estimators are typically *consistent*. Indeed, it is known that if there is a consistent sequence of estimators, then Bayes estimators are consistent. The rate of convergence is often best for reference Bayes estimators.

It is also generally accepted that a “good” estimator should be *admissible*, that is, *not dominated* by any other estimator in the sense that its expected loss under sampling (conditional to θ) cannot be larger for all θ values than that corresponding to another estimator. Any *proper* Bayes estimator is admissible; moreover, as established by Wald in the 1950’s, a procedure *must* be Bayesian (proper or improper) to be admissible. Most published admissibility results refer to quadratic loss functions, but they often extend to more general loss functions. Reference Bayes estimators are typically admissible with respect to intrinsic loss functions.

Notice, however, that many other apparently intuitive frequentist ideas on estimation have been proved to be potentially misleading. For example, given a sequence of n Bernoulli observations with parameter θ resulting in r positive trials, the *best unbiased* estimate of θ^2 is found to be $r(r-1)/\{n(n-1)\}$, which yields $\tilde{\theta}^2 = 0$ when $r = 1$; but to estimate the probability of two positive trials as zero, when one positive trial has been observed, is not at all sensible. In marked contrast, any Bayes reference estimator provides a reasonable answer. For example, the intrinsic estimator of θ^2 is simply $(\theta^*)^2$, where θ^* is the intrinsic estimator of θ described in Section 4.1. In particular, if $r = 1$ and $n = 2$ the intrinsic estimator of θ^2 is (as one would naturally expect) $(\theta^*)^2 = 1/4$.

Interval Estimation. As the sample size increases, the frequentist coverage probability of a posterior q -credible region typically converges to q so that, for *large samples*, Bayesian credible intervals may (under regularity conditions) be interpreted as *approximate* frequentist confidence regions: under repeated sampling, a Bayesian q -credible region of θ based on a large sample will cover the true value of θ approximately $100q\%$ of times. Detailed results are readily available for univariate problems. For instance, consider the probability model $\{p(D|\omega), \omega \in \Omega\}$, let $\theta = \theta(\omega)$ be any univariate quantity of interest, and let $\mathbf{t} = \mathbf{t}(D) \in T$ be any sufficient statistic. If $\theta_q(\mathbf{t})$ denotes the $100q\%$ quantile of the posterior distribution of θ which corresponds to some unspecified prior, so that

$$\Pr[\theta \leq \theta_q(\mathbf{t}) | \mathbf{t}] = \int_{\theta \leq \theta_q(\mathbf{t})} p(\theta | \mathbf{t}) d\theta = q, \quad (59)$$

then the coverage probability of the q -credible interval $\{\theta; \theta \leq \theta_q(\mathbf{t})\}$,

$$\Pr[\theta_q(\mathbf{t}) \geq \theta | \boldsymbol{\omega}] = \int_{\theta_q(\mathbf{t}) \geq \theta} p(\mathbf{t} | \boldsymbol{\omega}) d\mathbf{t}, \quad (60)$$

is such that $\Pr[\theta_q(\mathbf{t}) \geq \theta | \boldsymbol{\omega}] = \Pr[\theta \leq \theta_q(\mathbf{t}) | \mathbf{t}] + O(n^{-1/2})$. This *asymptotic* approximation is true for *all* (sufficiently regular) positive priors. However, the approximation is better, actually $O(n^{-1})$, for a particular class of priors known as (first-order) *probability matching* priors. Reference priors are typically found to be probability matching priors, so that they provide this improved asymptotic agreement. As a matter of fact, the agreement (in regular problems) is typically quite good even for relatively small samples.

Example 16. (Product of normal means). Consider the case where independent random samples $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_m\}$ have respectively been taken from the normal densities $N(x | \omega_1, 1)$ and $N(y | \omega_2, 1)$, and suppose that the quantity of interest is the product of their means, $\phi = \omega_1 \omega_2$ (for instance, one may be interested in inferences about the area ϕ of a rectangular piece of land, given measurements $\{x_i\}$ and $\{y_j\}$ of its sides). Notice that this is a simplified version of a problem that it is often encountered in the sciences, where one is interested in the product of several magnitudes, all of which have been measured with error. Using the procedure described in Example 13, with the natural approximating sequence induced by $(\omega_1, \omega_2) \in [-i, i]^2$, the ϕ -reference prior is found to be

$$\pi_\phi(\omega_1, \omega_2) \propto (n\omega_1^2 + m\omega_2^2)^{-1/2}, \quad (61)$$

very different from the uniform prior $\pi_{\omega_1}(\omega_1, \omega_2) = \pi_{\omega_2}(\omega_1, \omega_2) = 1$ which should be used to make objective inferences about either ω_1 or ω_2 . The prior $\pi_\phi(\omega_1, \omega_2)$ may be shown to provide approximate agreement between Bayesian credible regions and frequentist confidence intervals for ϕ ; indeed, this prior (with $m = n$) was originally suggested by Stein in the 1980's to obtain such approximate agreement. The same example was later used by Efron to stress the fact that, even within a fixed probability model $\{p(D | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$, the prior required to make objective inferences about some function of the parameters $\phi = \phi(\boldsymbol{\omega})$ must generally depend on the function ϕ .

The numerical agreement between reference Bayesian credible regions and frequentist confidence intervals is actually perfect in special circumstances. Indeed, as Lindley pointed out in the 1950's, this is the case in those problems of inference which may be transformed to location-scale problems.

Example 3. (Inference on normal parameters, continued). Let $D = \{x_1, \dots, x_n\}$ be a random sample from a normal distribution $N(x | \mu, \sigma)$. As mentioned before, the reference posterior of the quantity of interest μ is the Student distribution $\text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$. Thus, normalizing μ , the *posterior* distribution of

$t(\mu) = \sqrt{n-1}(\bar{x} - \mu)/s$, as a function of μ given D , is the standard Student $\text{St}(t | 0, 1, n-1)$ with $n-1$ degrees of freedom. On the other hand, this function t is recognized to be precisely the conventional t statistic, whose *sampling distribution* is well known to *also* be standard Student with $n-1$ degrees of freedom. It follows that, *for all sample sizes*, posterior *reference* credible intervals for μ given the data will be *numerically identical* to frequentist confidence intervals based on the sampling distribution of t .

A similar result is obtained in inferences about the variance. Thus, the reference *posterior* distribution of $\lambda = \sigma^{-2}$ is the Gamma distribution $\text{Ga}(\lambda | (n-1)/2, ns^2/2)$ and, hence, the *posterior* distribution of $r = ns^2/\sigma^2$, as a function of σ^2 given D , is a (central) χ^2 with $n-1$ degrees of freedom. But the function r is recognized to be a conventional statistic for this problem, whose *sampling distribution* is well known to *also* be χ^2 with $n-1$ degrees of freedom. It follows that, *for all sample sizes*, posterior *reference* credible intervals for σ^2 (or any one-to-one function of σ^2) given the data will be *numerically identical* to frequentist confidence intervals based on the sampling distribution of r .

6. A Simplified Case Study

To further illustrate the main aspects of Bayesian methods, and to provide a detailed, worked out example, a simplified version of a problem in engineering is analyzed below.

To study the reliability of a large production batch, n randomly selected items were put to an expensive, destructive test, yielding $D = \{x_1, \dots, x_n\}$ as their observed lifetimes in hours of continuous use. Context considerations suggested that the lifetime x_i of each item could be assumed to be exponential with hazard rate θ , so that $p(x_i | \theta) = \text{Ex}[x_i | \theta] = \theta e^{-\theta x_i}$, $\theta > 0$, and that, given θ , the lifetimes of the n items are independent. Quality engineers were interested in information on the actual value of the hazard rate θ , and on prediction of the lifetime x of similar items. In particular, they were interested in the compatibility of the observed data with advertised values of the hazard rate, and on the proportion of items whose lifetime could be expected to be longer than some required industrial specification.

The statistical analysis of exponential data makes use of the exponential-gamma distribution $\text{Eg}(x | \alpha, \beta)$, obtained as a continuous mixture of exponentials with a gamma density,

$$\text{Eg}(x | \alpha, \beta) = \int_0^{\infty} \theta e^{-\theta x} \text{Ga}(\theta | \alpha, \beta) d\theta = \frac{\alpha \beta^{\alpha}}{(x + \beta)^{\alpha+1}} \quad (62)$$

where $x \geq 0$, $\alpha > 0$ and $\beta > 0$. This is a monotonically decreasing density with mode at zero; if $\alpha > 1$, it has a mean $\text{E}[x | \alpha, \beta] = \beta/(\alpha - 1)$. Moreover, tail probabilities have a simple expression; indeed,

$$\Pr[x > t | \alpha, \beta] = \left\{ \frac{\beta}{\beta + t} \right\}^\alpha. \quad (63)$$

Likelihood function. Under the accepted assumptions on the mechanism which generated the data, $p(D | \theta) = \prod_j \theta e^{-\theta x_j} = \theta^n e^{-\theta s}$, which only depends on $s = \sum_j x_j$, the sum of the observations. Thus, $t = (s, n)$ is a *sufficient* statistic for this model. The corresponding MLE estimator is $\hat{\theta} = n/s$ and Fisher's information function is $F(\theta) = \theta^{-2}$. Moreover, the sampling distribution of s is the Gamma distribution $p(s | \theta) = \text{Ga}(s | n, \theta)$.

The actual data consisted of $n = 25$ uncensored observed lifetimes which, in thousands of hours, yielded a sum $s = 41.574$, hence a mean $\bar{x} = 1.663$, and a MLE $\hat{\theta} = 0.601$. The standard deviation of the observed lifetimes was 1.286 and their range was $[0.136, 5.591]$, showing the large variation (from a few hundred to a few thousand hours) typically observed in exponential data.

Using the results of Section 3.3, and the form of Fisher's information function given above, the *asymptotic* posterior distribution of θ is

$$p(\theta | D) \approx \text{N}(\theta | \hat{\theta}, \hat{\theta}/\sqrt{n}) = \text{N}(\theta | 0.601, 0.120).$$

This provided a first, quick approximation to the possible values of θ which, for instance, could be expected to belong to the interval $0.601 \pm 1.96 * 0.120$, or $(0.366, 0.837)$, with probability close to 0.95.

6.1. Objective Bayesian Analysis

The firm was to be audited on behalf of a major client. A report had to be prepared about the available information on the value of the hazard rate θ , *exclusively* based on the *documented* data D , as if this were the *only* information available. Within a Bayesian framework, this "objective" analysis (objective in the sense of not using any information beyond that provided by the data under the assumed model) may be achieved by computing the corresponding *reference* posterior distribution.

Reference prior and reference posteriors. The exponential model meets all necessary regularity conditions. Thus, using the results in Example 12 and the form of Fisher's information function mentioned above, the *reference prior function* (which in this case is also Jeffreys' prior) is simply $\pi(\theta) \propto F(\theta)^{1/2} = \theta^{-1}$. Hence, using Bayes' theorem, the reference posterior is $\pi(\theta | D) \propto p(D|\theta) \theta^{-1} \propto \theta^{n-1} e^{-s\theta}$, the kernel of a gamma density, so that

$$\pi(\theta | D) = \text{Ga}(\theta | n, s), \quad \theta > 0, \quad (64)$$

which has mean $E[\theta | D] = n/s$ (which is also the MLE $\hat{\theta}$), mode $(n - 1)/s$, and standard deviation $\sqrt{n}/s = \hat{\theta}/\sqrt{n}$. Thus, the reference posterior of the hazard rate was found to be $\pi(\theta | D) = \text{Ga}(\theta | 25, 41.57)$ (represented in Figure 5) with mean 0.601, mode 0.577, and standard deviation 0.120. One-dimensional numerical integration further yields $\Pr[\theta < 0.593 | D] = 0.5$, $\Pr[\theta < 0.389 | D] = 0.025$ and $\Pr[\theta < 0.859 | D] = 0.975$; thus, the median is 0.593, and the interval (0.389, 0.859) is a 95% reference posterior credible region (shaded area in Figure 5). The intrinsic estimator (see below) was found to be 0.590 (dashed line in Figure 5).

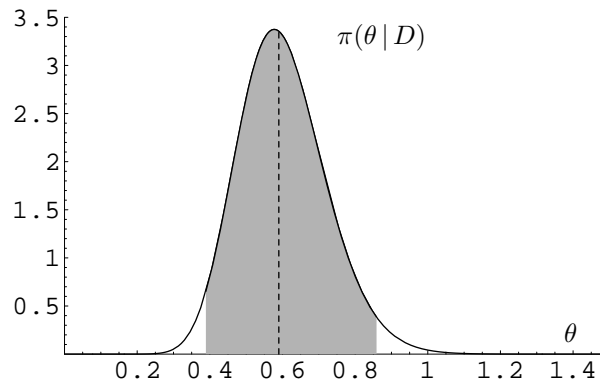


Figure 5. Reference posterior density of the hazard rate θ . The shaded region is a 95% credible interval. The dashed line indicates the position of the intrinsic estimator.

Under the accepted assumptions for the probability mechanism which has generated the data, the reference posterior distribution $\pi(\theta | D) = \text{Ga}(\theta | 25, 41.57)$ contained *all* that could be said about the value of the hazard rate θ on the exclusive basis of the observed data D . Figure 5 and the numbers quoted above respectively provided useful graphical and numerical summaries, but the fact that $\pi(\theta | D)$ is the *complete* answer (necessary for further work on prediction or decision making) was explained to the engineers by their consultant statistician.

Reference posterior predictive distribution. The reference predictive posterior density of a future lifetime x is

$$\pi(x | D) = \int_0^\infty \theta e^{-\theta x} \text{Ga}(\theta | n, s) d\theta = \text{Eg}(\theta | n, s) \quad (65)$$

with mean $s/(n - 1)$. Thus, the posterior predictive density of the lifetime of a random item produced in similar conditions was found to be

$$\pi(x | D) = \text{Eg}(x | 25, 41.57),$$

represented in Figure 6 against the background of a histogram of the observed data. The mean of this distribution is 1.732; hence, given data D , the expected lifetime of future similar items is 1.732 thousands of hours. The contract with their client specified a compensation for any item whose lifetime was smaller than 250 hours. Since

$$\Pr[x < b | D] = \int_0^b \text{Eg}(x | n, s) = 1 - \left\{ \frac{s}{s+b} \right\}^n, \quad (66)$$

the expected proportion of items with lifetime smaller than 250 hours is

$$\Pr[x < 0.250 | D] = 0.139,$$

the shaded area in Figure 6; thus, conditional on accepted assumptions, the engineers were advised to expect 14% of items to be non-conforming.

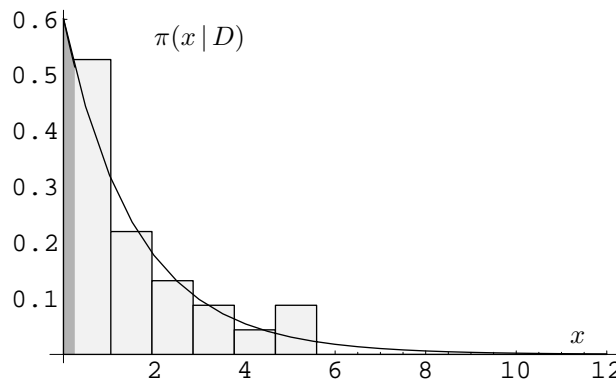


Figure 6. Reference predictive posterior density of lifetimes (in thousands of hours). The shaded region represents the probability of producing non conforming items, with lifetime smaller than 250 hours. The background is a histogram of the observed data.

Calibration. Consider $t = t(\theta) = (s/n)\theta$ as a function of θ , and its inverse transformation $\theta = \theta(t) = (n/s)t$. Since $t = t(\theta)$ is a one-to-one transformation of θ , if R_t is a q -posterior credible region for t , then $R_\theta = \theta(R_t)$ is a q -posterior credible region for θ . Moreover, changing variables, the reference *posterior* distribution of $t = t(\theta)$, as a function of θ conditional on s , is

$$\pi(t(\theta) | n, s) = \pi(\theta | n, s) / |\partial t(\theta) / \partial \theta| = \text{Ga}(t | n, n),$$

a gamma density which does not depend on s . On the other hand, the sampling distribution of the sufficient statistic s is $p(s | n, \theta) = \text{Ga}(\theta | n, \theta)$; therefore, the *sampling* distribution of $t = t(s) = (\theta/n)s$, as a function of s conditional to θ , is

$$p(t(s) | n, \theta) = p(s | n, \theta) / |\partial t(s) / \partial s| = \text{Ga}(t | n, n),$$

which does not contain θ and is precisely the *same* gamma density obtained before. It follows that, for *any* sample size n , all q -credible reference posterior regions of the hazard rate θ will *also* be frequentist confidence regions of level q . Any q -credible reference posterior region has, given the data, a (rational) degree of belief q of containing the true value of θ ; the result just obtained may be used to provide an exact calibration for this degree of belief. Indeed, for any $\theta > 0$ and any $q \in (0, 1)$, the limiting proportion of q -credible reference posterior regions which would cover the true value of θ under repeated sampling is precisely equal to q . It was therefore possible to explain to the engineers that, when reporting that the hazard rate θ of their production was expected to be within $(0.389, 0.859)$ with probability (rational degree of belief) 0.95, they could claim this to be a *calibrated* statement in the sense that hypothetical replications of the same *procedure* under controlled conditions, with samples simulated from *any* exponential distribution, would yield 95% of regions containing the value from which the sample was simulated.

Estimation. The commercial department could use any location measure of the reference posterior distribution of θ as an intuitive estimator $\tilde{\theta}$ of the hazard rate θ , but if a particular value has to be chosen with, say, some legal relevance, this would pose a decision problem for which an appropriate loss function $L(\tilde{\theta}, \theta)$ would have to be specified. Since no particular decision was envisaged, but the auditing firm nevertheless required that a particular estimator had to be quoted in the report, the attractive properties of the *intrinsic* estimator were invoked to justify its choice. The intrinsic discrepancy $d(\theta_i, \theta_j)$ between the *models* $\text{Ex}(x | \theta_i)$ and $\text{Ex}(x | \theta_j)$ is

$$d(\theta_i, \theta_j) = \min\{\delta(\theta_i | \theta_j), \delta(\theta_j | \theta_i)\}, \quad (67)$$

$$\delta(\theta_i | \theta_j) = (\theta_j / \theta_i) - 1 - \log(\theta_j / \theta_i).$$

As expected, $d(\theta_i, \theta_j)$ is a symmetric, non-negative concave function, which attains its minimum value zero if, and only if, $\theta_i = \theta_j$. The intrinsic estimator of the hazard rate is that $\theta^*(n, s)$ which minimizes the expected reference posterior loss,

$$d(\tilde{\theta} | n, s) = n \int_0^\infty d(\tilde{\theta}, \theta) \text{Ga}(\theta | n, s) d\theta. \quad (68)$$

To a very good approximation ($n > 1$), this is given by $\theta^*(n, s) \approx (2n - 1)/2s$, the arithmetic average of the reference posterior mean and the reference posterior

mode, quite close to the reference posterior median. With the available data, this approximation yielded $\theta^* \approx 0.5893$, while the exact value, found by numerical minimization was $\theta^* = 0.5899$. It was noticed that, since intrinsic estimation is an invariant procedure, the intrinsic estimate of any function $\phi(\theta)$ of the hazard rate would simply be $\phi(\theta^*)$.

Hypothesis Testing. A criterion of excellence in this industrial sector described first-rate production as one with a hazard rate smaller than 0.4, yielding an expected lifetime larger than 2500 hours. The commercial department was interested in whether or not the data obtained were *compatible* with the hypothesis that the actual hazard rate of the firm's production was that small. A direct answer was provided by the corresponding reference posterior probability $\Pr[\theta < 0.4 | D] = \int_0^{0.4} \text{Ga}(\theta | n, s) d\theta = 0.033$, suggesting that the hazard rate of present production might possibly be around 0.4, but it is actually unlikely to be that low.

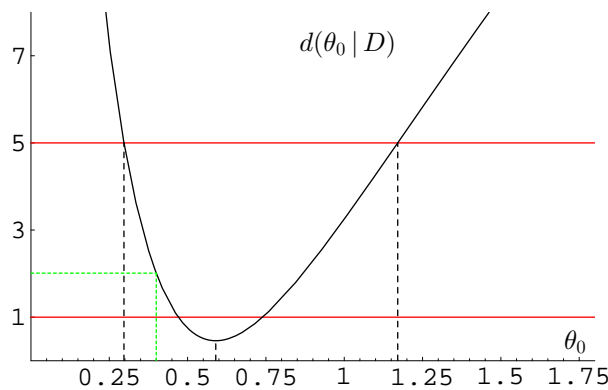


Figure 7. Expected reference posterior intrinsic loss for accepting θ_0 as a proxy for the true value of θ . The minimum is reached at the intrinsic estimator $\theta^* = 0.590$. Values of θ outside the interval $(0.297, 1.170)$ would be conventionally rejected.

Under pressure to provide a quantitative measure of the compatibility of the data with the *precise* value $\theta = \theta_0 = 0.4$, the statistician produced the expected intrinsic discrepancy $d(\theta_0 | n, s)$ from accepting θ_0 as a proxy for the true value of θ on the basis of data (n, s) by evaluating (68) at $\hat{\theta} = \theta_0$. It was recalled that the expected value of $d(\theta_0 | D)$ under repeated sampling is exactly equal to one when $\theta = \theta_0$, and that a large value of $d(\theta_0 | D)$ indicates strong evidence against θ_0 . Moreover, using a frequent language in engineering, the statistician explained that values of $d(\theta_0 | D) = d^*$ indicate, for $d^* = 2.5, 5.0$ or 8.5 , a level of evidence against $\theta = \theta_0$ comparable to the evidence against a zero mean that would be provided by

a normal observation x which was, respectively, 2, 3 or 4 standard deviations from zero. As indicated in Figure 7, values of θ_0 larger than 1.170 or smaller than 0.297 would be conventionally rejected by a “3 σ ” normal criterion. The actual value for θ_0 was found to be $d(0.4 | D) = 2.01$ (equivalent to 1.73 σ under normality). Thus, although there was some evidence suggesting that θ is likely to be larger than 0.4, the precise value $\theta = 0.4$ could not be definitely rejected on the exclusive basis of the information provided by the data D .

6.2. Sensitivity Analysis

Although conscious that this information could not be used in the report prepared for the client’s auditors, the firm’s management was interested in taping their engineers’ inside knowledge to gather further information on the actual lifetime of their products. This was done by exploring the consequences on the analysis of (i) introducing that information about the process which their engineers considered “beyond reasonable doubt” and (ii) introducing an “informed best guess” based on their experience with the product. The results, analyzed below and encapsulated in Figure 8, provide an analysis of the sensitivity of the final inferences on θ to changes in the prior information.

Limited prior information. When questioned by their consultant statistician, the production engineers claimed to know from past experience that the average lifetime $E[x]$ should be about 2250 hours, and that this average could not possibly be larger than 5000 or smaller than 650. Since $E[x | \theta] = \theta^{-1}$, those statements may directly be put in terms of conditions on the prior distribution of θ ; indeed, working in thousands of hours, they imply $E[\theta] = (2.25)^{-1} = 0.444$, and that $\theta \in \Theta_c = (0.20, 1.54)$. To describe mathematically this knowledge K_1 , the statistician used the corresponding *restricted* reference prior, that is the prior which maximizes the missing information about θ (*i.e.*, what it is unknown about its value) within the class of priors which satisfy those conditions. The reference prior restricted to $\theta \in \Theta_c$ and $E[\theta] = \mu$ is the solution of $\pi(\theta) \propto \theta^{-1} e^{-\lambda\theta}$, subject to the restrictions $\theta \in \Theta_c$ and $\int_{\Theta_c} \theta \pi(\theta | K_1) d\theta = \mu$. With the available data, this was numerically found to be $\pi(\theta | K_1) \propto \theta^{-1} e^{-2.088\theta}$, $\theta \in \Theta_c$. Bayes’ theorem was then used to obtain the corresponding posterior distribution $\pi(\theta | D, K_1) \propto p(D | \theta) \pi(\theta | K_1) \propto \theta^{24} e^{-43.69\theta}$, $\theta \in \Theta_c$, a gamma density $\text{Ga}(\theta | 25, 43.69)$ renormalized to $\theta \in \Theta_c$, which is represented by a thin line in Figure 8. Comparison with the unrestricted reference posterior, described by a solid line, suggests that, compared with the information provided by the data, the additional knowledge K_1 is relatively unimportant.

Detailed prior information. When further questioned by their consultant statistician, the production engineers guessed that the average lifetime is “surely” not larger than 3000 hours; when requested to be more precise they identified “surely” with

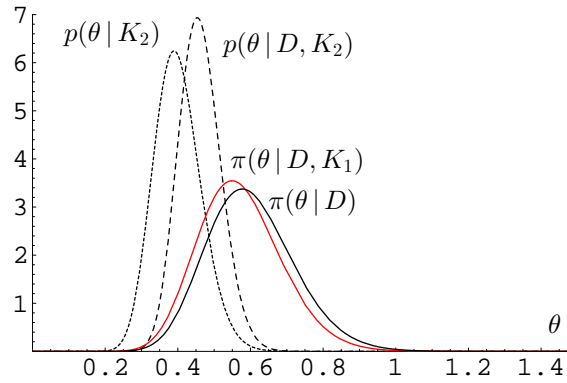


Figure 8. Probability densities of the hazard rate θ . Subjective prior (dotted line), subjective posterior (dashed line), partially informative reference posterior (thin line) and conventional reference posterior (solid line).

a 0.95 subjective degree of belief. Working in thousands of hours, this implies that $\Pr[\theta > 3^{-1}] = 0.95$. Together with their earlier claim on the expected lifetime, implying $E[\theta] = 0.444$, this was sufficient to completely specify a (subjective) prior distribution $p(\theta | K_2)$. To obtain a tractable form for such a prior, the statistician used a simple numerical routine to fit a restricted gamma distribution to those two statements, and found this to be $p(\theta | K_2) \propto \text{Ga}(\theta | \alpha, \beta)$, with $\alpha = 38.3$ and $\beta = 86.3$. Moreover, the statistician derived the corresponding *prior predictive* distribution $p(x | K_2) = \text{Eg}(x | \alpha, \beta)$ and found that the elicited prior $p(\theta)$ would imply, for instance, that $\Pr[x > 1 | K_2] = 0.64$, $\Pr[x > 3 | K_2] = 0.27$, and $\Pr[x > 10 | K_2] = 0.01$, so that the implied proportion of items with a lifetime over 1, 3, and 10 thousands of hours were, respectively, 64%, 27%, and 1%. The engineers declared that those numbers agreed with their experience and, hence, the statistician proceeded to accept $p(\theta) = \text{Ga}(\theta | 38.3, 86.3)$, represented with a dotted line in Figure 8, as a reasonable description of their prior beliefs. Using Bayes' theorem, the posterior density which corresponds to a $\text{Ga}(\theta | \alpha, \beta)$ prior is $p(\theta | D) = p(\theta | n, s) \propto \theta^n e^{-\theta s} \theta^{\alpha-1} e^{-\beta\theta} \propto \theta^{\alpha+n-1} e^{-(\beta+s)\theta}$, the kernel of a gamma density, so that

$$p(\theta | D) = \text{Ga}(\theta | \alpha + n, \beta + s), \quad \theta > 0. \quad (69)$$

Thus, the posterior distribution, combining the engineers' prior knowledge K_2 and data D was found to be $p(\theta | D, K_2) = \text{Ga}(\theta | 63.3, 127.8)$, represented with a dashed line in Figure 8. It is easily appreciated from Figure 8 that the 25

observations contained in the data analyzed do not represent a dramatic increase in information over that initially claimed by the production engineers, although the posterior distribution is indeed more concentrated than the prior, and it is displaced towards the values of θ suggested by the data. The firm's management would not be able to use this combined information in their auditing but, if they trusted their production engineers, they were advised to use $p(\theta | D, K_2)$ to further understand their production process, or to design policies intended to improve its performance.

7. Discussion and Further Issues

In writing a broad article it is always hard to decide what to leave out. This article focused on the basic concepts of the Bayesian paradigm; methodological topics which have unwillingly been omitted include design of experiments, sample surveys, linear models and sequential methods. The interested reader is referred to the bibliography for further information. This final section briefly reviews the main arguments for the Bayesian approach, and includes pointers to further issues which have not been discussed in more detail due to space limitations.

7.1. Coherence

By using probability distributions to characterize *all* uncertainties in the problem, the Bayesian paradigm reduces statistical inference to applied probability, thereby ensuring the coherence of the proposed solutions. There is no need to investigate, on a case by case basis, whether or not the solution to a particular problem is logically correct: a Bayesian result is only a *mathematical consequence of explicitly stated assumptions* and hence, unless a logical mistake has been committed in its derivation, it cannot be formally wrong. In marked contrast, conventional statistical methods are plagued with counterexamples. These include, among many others, negative estimators of positive quantities, q -confidence regions ($q < 1$) which consist of the whole parameter space, empty sets of "appropriate" solutions, and incompatible answers from alternative methodologies simultaneously supported by the theory.

The Bayesian approach does require, however, the specification of a (prior) probability distribution over the parameter space. The sentence "a prior distribution does not exist for this problem" is often stated to justify the use of non-Bayesian methods. However, the general representation theorem *proves the existence* of such a distribution whenever the observations are assumed to be exchangeable (and, if they are assumed to be a random sample then, *a fortiori*, they are assumed to be exchangeable). To ignore this fact, and to proceed as if a prior distribution did not exist, just because it is not easy to specify, is mathematically untenable.

7.2. Objectivity

It is generally accepted that any statistical analysis is subjective, in the sense that it is always conditional on accepted assumptions (on the structure of the data, on the probability model, and on the outcome space) and those assumptions, although possibly well founded, are definitely *subjective* choices. It is, therefore, mandatory to make all assumptions very explicit.

Users of conventional statistical methods rarely dispute the mathematical foundations of the Bayesian approach, but claim to be able to produce “objective” answers in contrast to the possibly subjective elements involved in the choice of the prior distribution.

Bayesian methods do indeed require the choice of a prior distribution, and critics of the Bayesian approach systematically point out that in many important situations, including scientific reporting and public decision making, the results must exclusively depend on documented data which might be subject to independent scrutiny. This is of course true, but those critics choose to ignore the fact that this particular case is covered within the Bayesian approach by the use of *reference* prior distributions which (i) are mathematically derived from the accepted probability model (and, hence, they are “objective” insofar as the choice of that model might be objective) and, (ii) by construction, they produce posterior probability distributions which, given the accepted probability model, *only* contain the information about their values which data may provide and, *optionally*, any further contextual information over which there might be universal agreement.

An issue related to objectivity is that of the operational meaning of reference posterior probabilities; it is found that the analysis of their behaviour under repeated sampling provides a suggestive form of calibration. Indeed,

$$\Pr[\theta \in R | D] = \int_R \pi(\theta | D) d\theta,$$

the reference posterior probability that $\theta \in R$, is *both* a measure of the conditional uncertainty (given the assumed model and the observed data D) about the event that the unknown value of θ belongs to $R \subset \Theta$, and the limiting proportion of the regions which would cover θ under repeated sampling conditional on data “sufficiently similar” to D . Under broad conditions (to guarantee regular asymptotic behaviour), all large data sets from the same model are “sufficiently similar” among themselves in this sense and hence, given those conditions, reference posterior credible regions are *approximate* unconditional frequentist confidence regions.

The conditions for this approximate *unconditional* equivalence to hold exclude, however, important special cases, like those involving “extreme” or “relevant” observations. In very special situations, when probability models may be transformed to location-scale models, there is an exact unconditional equivalence; in

those cases reference posterior credible intervals are, for any sample size, exact unconditional frequentist confidence intervals.

7.3. Applicability

In sharp contrast to most conventional statistical methods, which may only be exactly applied to a handful of relatively simple stylized situations, Bayesian methods are (in theory) totally general. Indeed, for a given probability model and prior distribution over its parameters, the derivation of posterior distributions is a well-defined mathematical exercise. In particular, Bayesian methods do not require any particular regularity conditions on the probability model, do not depend on the existence of sufficient statistics of finite dimension, do not rely on asymptotic relations, and do not require the derivation of any sampling distribution, nor (a fortiori) the existence of a “pivotal” statistic whose sampling distribution is independent of the parameters.

However, when used in complex models with many parameters, Bayesian methods often require the computation of multidimensional definite integrals and, for a long time in the past, this requirement effectively placed practical limits on the complexity of the problems which could be handled. This has dramatically changed in recent years with the general availability of large computing power, and the parallel development of simulation-based numerical integration strategies like *importance sampling* or *Markov chain Monte Carlo* (MCMC). These methods provide a structure within which many complex models may be analyzed using generic software. MCMC is numerical integration using Markov chains. Monte Carlo integration proceeds by drawing samples from the required distributions, and computing sample averages to approximate expectations. MCMC methods draw the required samples by running appropriately defined Markov chains for a long time; specific methods to construct those chains include the Gibbs sampler and the Metropolis algorithm, originated in the 1950's in the literature of statistical physics. The development of improved algorithms and appropriate diagnostic tools to establish their convergence, remains a very active research area.

Actual scientific research often requires the use of models that are far too complex for conventional statistical methods. This article concludes with a glimpse at some of them.

Hierarchical structures. Consider a situation where a possibly variable number n_i of observations, $\{\mathbf{x}_{ij}, j = 1, \dots, n_i\}, i = 1, \dots, m$, are made on each of m internally homogeneous subsets of some population. For instance, a firm might have chosen m production lines for inspection, and n_i items might have been randomly selected among those made by production line i , so that \mathbf{x}_{ij} is the result of the measurements made on item j of production line i . As another example, animals of some species are captured to study their metabolism, and a blood sample taken before releasing

them again; the procedure is repeated in the same habitat for some time, so that some of the animals are recaptured several times, and \mathbf{x}_{ij} is the result of the analysis of the j -th blood sample taken from animal i . In those situations, it is often appropriate to assume that the n_i observations on subpopulation i are exchangeable, so that they may be treated as a random sample from some model $p(\mathbf{x} | \boldsymbol{\theta}_i)$ indexed by a parameter $\boldsymbol{\theta}_i$ which depends on the subpopulation observed, and that the parameters which label the subpopulations may also be assumed to be exchangeable, so that $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m\}$ may be treated as a random sample from some distribution $p(\boldsymbol{\theta} | \boldsymbol{\omega})$. Thus, the complete *hierarchical* model which is assumed to have generated the observed data $D = \{\mathbf{x}_{11}, \dots, \mathbf{x}_{mn_m}\}$ is of the form

$$p(D | \boldsymbol{\omega}) = \int_{\Theta^m} \left[\prod_{j=1}^{n_i} p(\mathbf{x}_{ij} | \boldsymbol{\theta}_i) \right] \left[\prod_{i=1}^m p(\boldsymbol{\theta}_i | \boldsymbol{\omega}) \right] \left[\prod_{i=1}^m d\boldsymbol{\theta}_i \right]. \quad (70)$$

Hence, under the Bayesian paradigm, a family of conventional probability models, say $p(\mathbf{x} | \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \Theta$, and an appropriate “structural” prior $p(\boldsymbol{\theta} | \boldsymbol{\omega})$, may be naturally combined to produce a versatile, complex model $\{p(D | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$ whose analysis is often well beyond the scope of conventional statistics. The Bayesian solution only requires the specification a prior distribution $p(\boldsymbol{\omega})$, the use Bayes’ theorem to obtain the corresponding posterior $p(\boldsymbol{\omega} | D) \propto p(D | \boldsymbol{\omega}) p(\boldsymbol{\omega})$, and the performance of the appropriate probability transformations to derive the posterior distributions of the quantities of interest (which may well be functions of $\boldsymbol{\omega}$, functions of the $\boldsymbol{\theta}_i$ ’s, or functions of future observations). As in any other Bayesian analysis, the prior distribution $p(\boldsymbol{\omega})$ has to describe available knowledge about $\boldsymbol{\omega}$; if none is available, or if an objective analysis is required, an appropriate reference prior function $\pi(\boldsymbol{\omega})$ may be used.

Contextual information. In many problems of statistical inference, objective and universally agreed contextual information is available on the parameter values. This information is usually very difficult to handle within the framework of conventional statistics, but it is easily incorporated into a Bayesian analysis by simply restricting the prior distribution to the class $\{\mathcal{P}\}$ of priors which are compatible with such information. As an example, consider the frequent problem in archaeology of trying to establish the occupation period $[\alpha, \beta]$ of a site by some past culture on the basis of the radiocarbon dating of organic samples taken from the excavation. Radiocarbon dating is not precise, so that each dating x_i is typically taken to be a normal observation from a distribution $N(x | \mu(\theta_i), \sigma_i)$, where θ_i is the actual, unknown calendar date of the sample, $\mu(\theta)$ is an internationally agreed calibration curve, and σ_i is a known standard error quoted by the laboratory. The actual calendar dates $\{\theta_1, \dots, \theta_m\}$ of the samples are typically assumed to be uniformly distributed within the occupation period $[\alpha, \beta]$; however, stratigraphic evidence indicates some

partial orderings for, if sample i was found on top of sample j in undisturbed layers, then $\theta_i > \theta_j$. Thus, if \mathcal{C} denotes the class of values of $\{\theta_1, \dots, \theta_m\}$ which satisfy those known restrictions, data may be assumed to have been generated by the hierarchical model

$$p(x_1, \dots, x_m | \alpha, \beta) = \int_{\mathcal{C}} \left[\prod_{i=1}^m N(x_i | \mu(\theta_i), \sigma_i^2) \right] (\beta - \alpha)^{-m} d\theta_1 \dots d\theta_m. \quad (71)$$

Often, contextual information further indicates an absolute lower bound α_0 and an absolute upper bound β_0 for the period investigated, so that $\alpha_0 < \alpha < \beta < \beta_0$. If no further documented information is available, the corresponding restricted reference prior for the quantities of interest, $\{\alpha, \beta\}$ should be used; this is found to be $\pi(\alpha, \beta) \propto (\beta - \alpha)^{-1}$ whenever $\alpha_0 < \alpha < \beta < \beta_0$ and zero otherwise. The corresponding reference posterior

$$\pi(\alpha, \beta | x_1, \dots, x_m) \propto p(x_1, \dots, x_m | \alpha, \beta) \pi(\alpha, \beta)$$

summarizes all available information on the occupation period.

Covariate information. Over the last 30 years, both linear and non-linear regression models have been analyzed from a Bayesian point of view at increasing levels of sophistication. These studies range from the elementary objective Bayesian analysis of simple linear regression structures (which parallel their frequentist counterparts) to the sophisticated analysis of time series involved in dynamic forecasting which often make use of complex hierarchical structures. The field is far too large to be reviewed in this article, but the bibliography contains some relevant pointers.

Model Criticism. It has been stressed that *any* statistical analysis is conditional on the accepted assumptions of the probability model which is presumed to have generated the data. Recent years have shown a huge effort into the development of Bayesian procedures for *model criticism* and *model choice*. Most of these are sophisticated elaborations of the procedures described in Section 4.2 under the heading of hypothesis testing. Again, this is too large a topic to be reviewed here, but some key references are included in the bibliography.

Annotated Bibliography

- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. ringer. [A thorough account of Bayesian methods emphasizing its decision-theoretical aspects]
- Bernardo, J. M. (1979a). Expected information as expected utility. *Ann. Statist.* **7**, 686–690. [Establishes statistical inference as a decision problem with an information-based utility function]
- Bernardo, J. M. (1979b). Reference posterior distributions for Bayesian inference. *J. R. Statist. Soc. B* **41**, 113-147 (with discussion). Reprinted in *Bayesian Inference I* (G. C. Tiao and N. G. Polson, eds). Oxford: Edward Elgar, 229-263. [The original paper on reference analysis]
- Bernardo, J. M. (1997). Noninformative priors do not exist. *J. Statist. Planning and Inference* **65**, 159–189 (con discusión). [A non-technical analysis of the polemic on objective Bayesian statistics]
- Bernardo, J. M. (2005). Reference analysis. *Handbook of Statistics* **25** (D. Dey & C. R. Rao, eds). Amsterdam: North Holland (in press).
<www.uv.es/~bernardo/RefAna.pdf>
[A long introductory paper to reference analysis, which also includes a discussion of estimation and hypothesis testing using an information-theoretical based loss, the intrinsic discrepancy]
- Bernardo, J. M. and Juárez, M. (2003). Intrinsic Estimation. *Bayesian Statistics 7* (J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith and M. West, eds.). Oxford: University Press, 465-476. [A new decision-oriented approach to invariant point estimation]
- Bernardo, J. M. and Ramón, J. M. (1998). An introduction to Bayesian reference analysis: inference on the ratio of multinomial parameters. *The Statistician* **47**, 1–35. [An elementary introduction to objective Bayesian analysis]
- Bernardo, J. M. and Rueda, R. (2002). Bayesian hypothesis testing: A reference approach. *International Statistical Review* **70**, 351-372. [A new decision-oriented approach to sharp hypothesis testing]
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*, Chichester: Wiley. [A thorough account of key concepts and theoretical results in Bayesian statistics at a graduate level, with extensive bibliography]
- Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P., Heckerman, D., Smith, A. F. M. and West, M. (eds). (2003). *Bayesian Statistics 7*. Oxford: Oxford University Press. [The Proceedings of the 7th Valencia International Meeting on Bayesian Statistics; the Valencia meetings, held every four years, provide definite up-to-date overviews on current research within the Bayesian paradigm.]

- Berry, D. A. (1996). *Statistics, a Bayesian Perspective*. Belmont, CA: Wadsworth. [A very good introduction to Bayesian statistics from a subjectivist viewpoint]
- Box, G. E. P. and Tiao, G. C. (1973). *Bayesian Inference in Statistical Analysis*. Reading, MA: Addison-Wesley. [An excellent objective Bayesian account of standard statistical problems]
- Dawid, A. P., Stone, M. and Zidek, J. V. (1973). Marginalization paradoxes in Bayesian and structural inference. *J. R. Statist. Soc. B* **35**, 189-233 (with discussion). [Proves that a unique “non-informative prior for all parameters of interest within a given model is not possible]
- DeGroot, M. H. (1970). *Optimal Statistical Decisions*, New York: McGraw-Hill. [A thorough account of Bayesian decision theory and Bayesian inference with a rigorous treatment of foundations]
- Efron, B. (1986). Why isn't everyone a Bayesian? *Amer. Statist.* **40**, 1-11 (con discusión). [A good example of the polemic between Bayesian and non-Bayesian approaches to statistics]
- de Finetti, B. (1970). *Teoria delle Probabilità*, Turin: Einaudi. English translation as *Theory of Probability* in 1975, Chichester: Wiley. [An outstanding book on probability and statistics from a subjective viewpoint]
- Geisser, S. (1993). *Predictive Inference: an Introduction*. London: Chapman and Hall. [A comparative account of frequentist and objective Bayesian methods of prediction]
- Gelfand, A. E. and Smith, A. F. M. (1990). Sampling based approaches to calculating marginal densities. *J. Amer. Statist. Assoc.* **85**, 398-409. [An excellent primer on simulation-based techniques to numerical integration in the context of Bayesian statistics]
- Gelman, A., Carlin, J. B., Stern, H. and Rubin, D. B. (1995). *Bayesian Data Analysis*. London: Chapman and Hall. [A comprehensive treatment of Bayesian data analysis emphasizing computational tools]
- Gilks, W. R., Richardson, S. and Spiegelhalter, D. J. (1996). *Markov Chain Monte Carlo in Practice*. London: Chapman and Hall. [An excellent introduction to MCMC methods and their applications]
- Jaynes, E. T. (1976). Confidence intervals vs. Bayesian intervals. *Foundations of Probability Theory, Statistical Inference and Statistical Theories of Science* **2** (W. L. Harper and C. A. Hooker, eds). Dordrecht: Reidel, 175-257 (with discussion). [A provocative collection of counter-examples to conventional statistical methods]
- Kass, R. E. and Raftery, A. E. (1995). Bayes factors. *J. Amer. Statist. Assoc.* **90**, 773-795. [A very good review of Bayes factor methods for hypothesis testing]

- Lindley, D. V. (1972). *Bayesian Statistics, a Review*. Philadelphia, PA: SIAM. [A sharp comprehensive review of the whole subject up to the 1970's, emphasizing its internal consistency]
- Lindley, D. V. (1985). *Making Decisions*. (2nd ed.) Chichester: Wiley. [The best elementary introduction to Bayesian decision analysis]
- Lindley, D. V. (1990). The 1988 Wald memorial lecture: The present position in Bayesian Statistics. *Statist. Sci.* **5**, 44-89 (con discusión). [An informative account of the Bayesian paradigm and its relationship with other attitudes to inference]
- Lindley, D. V. (2000). The philosophy of statistics. *The Statistician* **49**, 293–337 (con discusión). [A recent description of the Bayesian paradigm from a subjectivist viewpoint]
- O'Hagan, A. (1994). *Bayesian Inference* London: Edward Arnold. [A good account of Bayesian inference integrated into Kendall's Library of Statistics]
- Press, S. J. (1972). *Applied Multivariate Analysis: using Bayesian and Frequentist Methods of Inference*. Melbourne, FL: Krieger. [A comprehensive and comparative account of frequentist and objective Bayesian methods of inference in multivariate problems]
- West, M. and Harrison, P. J. (1989). *Bayesian Forecasting and Dynamic Models*. ringer. [An excellent thorough account of Bayesian time series analysis]
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics*. New York: Wiley. Reprinted in 1987, Melbourne, FL: Krieger. [A detailed objective Bayesian analysis of linear models]