

# An Introduction to Objective Bayesian Statistics

**José M. Bernardo**

*Universitat de València, Spain*

<jose.m.bernardo@uv.es>

<http://www.uv.es/bernardo>

Université de Neuchâtel, Switzerland

March 15th–March 17th, 2006

# Summary

## 1. *Concept of Probability*

*Introduction.* Notation. Statistical models.

*Intrinsic discrepancy.* Intrinsic convergence of distributions.

*Foundations.* Probability as a rational degree of belief.

## 2. *Basics of Bayesian Analysis*

*Parametric inference.* The learning process.

*Reference analysis.* No relevant initial information.

*Inference summaries.* Point and region estimation.

*Prediction.* Regression.

*Hierarchical models.* Exchangeability.

## 3. *Decision Making*

*Structure of a decision problem.* Intrinsic loss functions.

*Point and region estimation.* Intrinsic estimators and credible regions.

*Hypothesis testing.* Bayesian reference criterion (BRC).

# 1. Concept of Probability

## 1.1. Introduction

- Tentatively accept a *formal* statistical model
  - Typically suggested by informal descriptive evaluation
  - Conclusions conditional on the assumption that model is correct
- Bayesian approach firmly based on *axiomatic foundations*
  - Mathematical need to describe by probabilities all uncertainties
  - Parameters *must* have a (*prior*) distribution describing available information about their values
  - Not* a description of their variability (*fixed unknown* quantities), but a description of the *uncertainty* about their true values.
- Important particular case: no relevant (or subjective) initial information: scientific and industrial reporting, public decision making, ...
  - Prior *exclusively* based on model assumptions and available, well-documented data: *Objective Bayesian Statistics*

- *Notation*

- Under conditions  $C$ ,  $p(\mathbf{x} | C)$ ,  $\pi(\boldsymbol{\theta} | C)$  are, respectively, *probability densities* (or mass) functions of *observables*  $\mathbf{x}$  and *parameters*  $\boldsymbol{\theta}$   
 $p(\mathbf{x} | C) \geq 0$ ,  $\int_{\mathcal{X}} p(\mathbf{x} | C) d\mathbf{x} = 1$ ,  $E[\mathbf{x} | C] = \int_{\mathcal{X}} \mathbf{x} p(\mathbf{x} | C) d\mathbf{x}$ ,  
 $\pi(\boldsymbol{\theta} | C) \geq 0$ ,  $\int_{\Theta} \pi(\boldsymbol{\theta} | C) d\boldsymbol{\theta} = 1$ ,  $E[\boldsymbol{\theta} | C] = \int_{\Theta} \boldsymbol{\theta} \pi(\boldsymbol{\theta} | C) d\boldsymbol{\theta}$ .
- Special densities (or mass) functions use specific notation, as  $N(x | \mu, \sigma)$ ,  $Bi(x | n, \theta)$ , or  $Pn(x | \lambda)$ . Other examples:

---

Beta       $\{\text{Be}(x | \alpha, \beta), \quad 0 < x < 1, \quad \alpha > 0, \beta > 0\}$

$$\text{Be}(x | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$


---

Gamma       $\{\text{Ga}(x | \alpha, \beta), \quad x > 0, \quad \alpha > 0, \beta > 0\}$

$$\text{Ga}(x | \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$$


---

Student       $\{\text{St}(x | \mu, \sigma, \alpha), \quad x \in \mathfrak{R}, \quad \mu \in \mathfrak{R}, \sigma > 0, \alpha > 0\}$

$$\text{St}(x | \mu, \sigma, \alpha) = \frac{\Gamma\{(\alpha+1)/2\}}{\Gamma(\alpha/2)} \frac{1}{\sigma\sqrt{\alpha\pi}} \left[ 1 + \frac{1}{\alpha} \left( \frac{x-\mu}{\sigma} \right)^2 \right]^{-(\alpha+1)/2}$$


---

- *Statistical Models*

- *Statistical model* generating  $\mathbf{x} \in \mathcal{X}$ ,  $\{p(\mathbf{x} | \boldsymbol{\theta}), \mathbf{x} \in \mathcal{X}, \boldsymbol{\theta} \in \Theta\}$   
*Parameter vector*  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_k\} \in \Theta$ . *Parameter space*  $\Theta \subset \mathbb{R}^k$ .  
*Data set*  $\mathbf{x} \in \mathcal{X}$ . *Sampling (Outcome) space*  $\mathcal{X}$ , of arbitrary structure.
- *Likelihood function* of  $\mathbf{x}$ ,  $l(\boldsymbol{\theta} | \mathbf{x})$ .  
 $l(\boldsymbol{\theta} | \mathbf{x}) = p(\mathbf{x} | \boldsymbol{\theta})$ , as a function of  $\boldsymbol{\theta} \in \Theta$ .
- *Maximum likelihood estimator (mle)* of  $\boldsymbol{\theta}$   
 $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{x}) = \arg \sup_{\boldsymbol{\theta} \in \Theta} l(\boldsymbol{\theta} | \mathbf{x})$
- *Data*  $\mathbf{x} = \{x_1, \dots, x_n\}$  *random sample* (iid) from model if  
 $p(\mathbf{x} | \boldsymbol{\theta}) = \prod_{j=1}^n p(x_j | \boldsymbol{\theta})$ ,  $x_j \in \mathcal{X}$ ,  $\mathcal{X} = \mathcal{X}^n$
- *Behaviour under repeated sampling* (general, not iid data)  
 Considering  $\{\mathbf{x}_1, \mathbf{x}_2, \dots\}$ , a (possibly infinite) sequence of possible replications of the *complete* data set  $\mathbf{x}$ .  
 Denote by  $\mathbf{x}^{(m)} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  a finite set of  $m$  such replications.
- *Asymptotic results* obtained as  $m \rightarrow \infty$

## 1.2. Intrinsic Divergence

- *Logarithmic divergences*

- The logarithmic divergence (Kullback-Leibler)  $k\{\hat{p} | p\}$  of a density  $\hat{p}(\mathbf{x})$ ,  $\mathbf{x} \in \mathcal{X}$  from its true density  $p(\mathbf{x})$ , is

$$\kappa\{\hat{p} | p\} = \int_{\mathcal{X}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{\hat{p}(\mathbf{x})} d\mathbf{x}, \text{ (provided this exists)}$$

The functional  $\kappa\{\hat{p} | p\}$  is non-negative, (zero iff,  $\hat{p}(\mathbf{x}) = p(\mathbf{x})$  a.e.) and *invariant* under one-to-one transformations of  $\mathbf{x}$ .

- But  $\kappa\{p_1 | p_2\}$  is *not symmetric* and diverges if, strictly,  $\mathcal{X}_2 \subset \mathcal{X}_1$ .

- *Intrinsic discrepancy between distributions*

- $\delta\{p_1, p_2\} = \min \left\{ \int_{\mathcal{X}_1} p_1(\mathbf{x}) \log \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} d\mathbf{x}, \int_{\mathcal{X}_2} p_2(\mathbf{x}) \log \frac{p_2(\mathbf{x})}{p_1(\mathbf{x})} d\mathbf{x} \right\}$

The *intrinsic discrepancy*  $\delta\{p_1, p_2\}$  is non-negative (zero iff,  $p_1 = p_2$  a.e.), and *invariant* under one-to-one transformations of  $\mathbf{x}$ ,

- Defined if  $\mathcal{X}_2 \subset \mathcal{X}_1$  or  $\mathcal{X}_1 \subset \mathcal{X}_2$ , operative interpretation as the minimum amount of information (in *nits*) required to discriminate.

- *Interpretation and calibration of the intrinsic discrepancy*

- Let  $\{p_1(\mathbf{x} | \boldsymbol{\theta}_1), \boldsymbol{\theta}_1 \in \Theta_1\}$  or  $\{p_2(\mathbf{x} | \boldsymbol{\theta}_2), \boldsymbol{\theta}_2 \in \Theta_2\}$  be two alternative statistical models for  $\mathbf{x} \in X$ , one of which is assumed to be true. The intrinsic divergence  $\delta\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\} = \delta\{p_1, p_2\}$  is then *minimum expected log-likelihood ratio in favour of the true model*.

Indeed, if  $p_1(\mathbf{x} | \boldsymbol{\theta}_1)$  true model, the expected log-likelihood ratio in its favour is  $E_1[\log\{p_1(\mathbf{x} | \boldsymbol{\theta}_1)/p_2(\mathbf{x} | \boldsymbol{\theta}_1)\}] = \kappa\{p_2 | p_1\}$ . If the true model is  $p_2(\mathbf{x} | \boldsymbol{\theta}_2)$ , the expected log-likelihood ratio in favour of the true model is  $\kappa\{p_2 | p_1\}$ . But  $\delta\{p_2 | p_1\} = \min[\kappa\{p_2 | p_1\}, \kappa\{p_1 | p_2\}]$ .

- *Calibration*.  $\delta = \log[100] \approx 4.6$  nits, likelihood ratios for the true model larger than 100 making *discrimination very easy*.

$\delta = \log(1 + \varepsilon) \approx \varepsilon$  nits, likelihood ratios for the true model may about  $1 + \varepsilon$  making *discrimination very hard*.

Intrinsic Discrepancy $\delta$	0.01	0.69	2.3	4.6	6.9
Average Likelihood Ratio for <b>true</b> model $\exp[\delta]$	1.01	2	10	100	1000

- *Example.* Conventional Poisson approximation  $P_n(r | n\theta)$  of Binomial probabilities  $\text{Bi}(r | n, \theta)$

Intrinsic discrepancy between Binomial and Poisson distributions

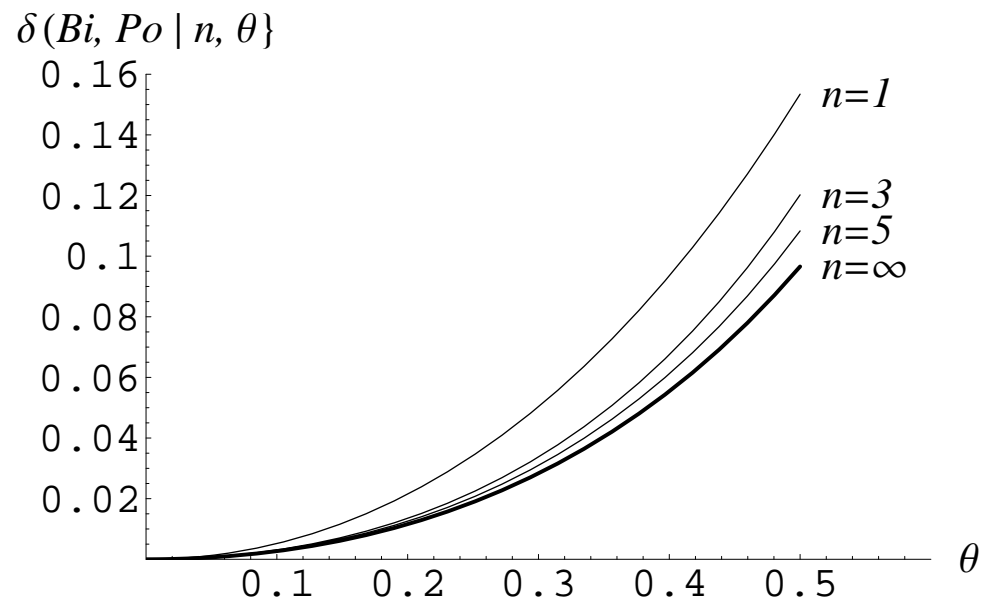
$$\begin{aligned} \delta\{\text{Bi}(r | n, \theta), \text{Po}(r | n\theta)\} &= \min[k\{\text{Bi} | \text{Po}\}, k\{\text{Po} | \text{Bi}\}] = k\{\text{Bi} | \text{Po}\} \\ &= \sum_{r=0}^n \text{Bi}(r | n, \theta) \log[\text{Bi}(r | n, \theta) / \text{Po}(r | n\theta)] = \delta\{n, \theta\} \end{aligned}$$

$$\delta\{3, 0.05\} = 0.00074$$

$$\delta\{5000, 0.05\} = 0.00065$$

$$\delta\{\infty, \theta\} = \frac{1}{2}[-\theta - \log(1 - \theta)]$$

Good Poisson approximations are *impossible* if  $\theta$  is not small, however large  $n$  might be.



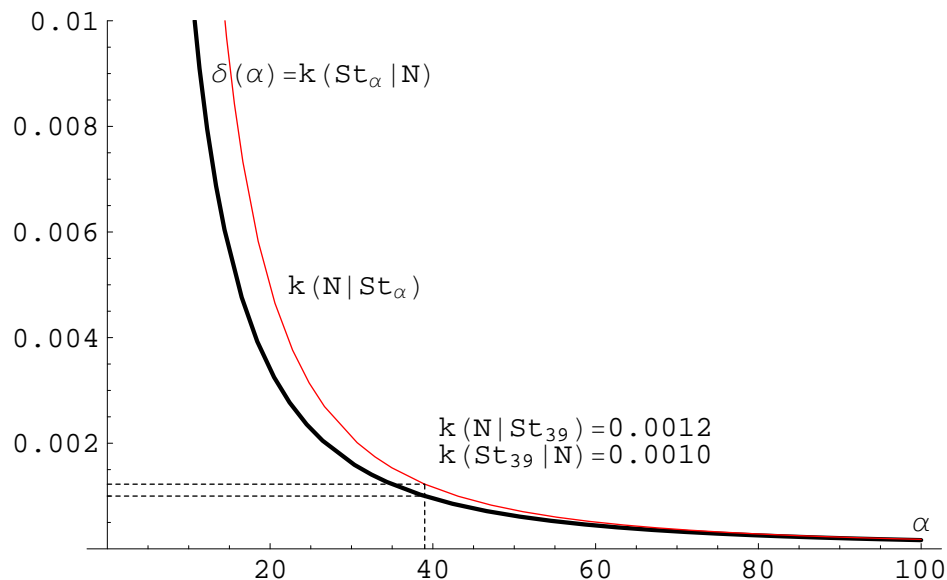


- *Intrinsic Convergence of Distributions*

□ *Intrinsic convergence.* A sequence of probability densities (or mass) functions  $\{p_i(\mathbf{x})\}_{i=1}^{\infty}$  converges *intrinsically* to  $p(\mathbf{x})$  if (and only if) the intrinsic divergence between  $p_i(x)$  and  $p(x)$  converges to zero. *i.e.*, iff  $\lim_{i \rightarrow \infty} \delta(p_i, p) = 0$ .

□ *Example.* Normal approximation to a Student distribution.

$$\begin{aligned} \delta(\alpha) &= \delta\{\text{St}(x \mid \mu, \sigma, \alpha), \text{N}(x \mid \mu, \sigma)\} = \min[k\{\text{St}_\alpha \mid \text{N}\}, k\{\text{N} \mid \text{St}_\alpha\}] \\ &= k\{\text{St}_\alpha \mid \text{N}\} = \int_{\mathcal{R}} \text{N}(x \mid 0, 1) \log \frac{\text{N}(x \mid 0, 1)}{\text{St}(x \mid 0, 1, \alpha)} dx \approx \frac{7}{\alpha(22 + 4\alpha)} \end{aligned}$$



$k\{\text{N} \mid \text{St}_\alpha\}$  diverges for  $\alpha \leq 2$

$k\{\text{St}_\alpha \mid \text{N}\}$  is finite for all  $\alpha > 0$ .

$\delta(18) \approx 0.04$      $\delta(25) \approx 0.02$

Expected log-density ratios  
at least 0.001 when  $\alpha < 40$ .

## 1.3. Foundations

- *Foundations of Statistics*

- Axiomatic foundations on rational description of uncertainty imply that the uncertainty about all unknown quantities should be measured with *probability* distributions  $\{\pi(\boldsymbol{\theta} | C), \boldsymbol{\theta} \in \Theta\}$  describing the plausibility of their given available conditions  $C$ .
- Axioms have a strong intuitive appeal; examples include
  - *Transitivity of plausibility.*  
If  $E_1 \succ E_2 | C$ , and  $E_2 \succ E_3 | C$ , then  $E_1 \succ E_3 | C$
  - *The sure-thing principle.*  
If  $E_1 \succ E_2 | A, C$  and  $E_1 \succ E_2 | \bar{A}, C$ , then  $E_1 \succ E_2 | C$ .
- Axioms are not a *description* of actual human activity, but a *normative* set of principles for those aspiring to rational behaviour.
- “Absolute” probabilities do not exist. Typical applications produce  $\Pr(E | \boldsymbol{x}, A, K)$ , a measure of rational belief in the occurrence of the *event*  $E$ , given data  $\boldsymbol{x}$ , assumptions  $A$  and available knowledge  $K$ .

- *Probability as a Measure of Conditional Uncertainty*

- Axiomatic foundations imply that  $\Pr(E | C)$ , the *probability* of an event  $E$  given  $C$  is *always* a conditional measure of the (presumably rational) uncertainty, on a  $[0, 1]$  scale, about the occurrence of  $E$  in conditions  $C$ .

- *Probabilistic diagnosis.*  $V$  is the event that a person carries a virus and  $+$  a positive test result. *All* related probabilities, *e.g.*,

$$\Pr(+ | V) = 0.98, \Pr(+ | \bar{V}) = 0.01, \Pr(V | K) = 0.002,$$

$$\Pr(+ | K) = \Pr(+ | V)\Pr(V | K) + \Pr(+ | \bar{V})\Pr(\bar{V} | K) = 0.012$$

$$\Pr(V | +, A, K) = \frac{\Pr(+ | V)\Pr(V | K)}{\Pr(+ | K)} = 0.164 \text{ (Bayes' Theorem)}$$

are conditional uncertainty measures (and proportion estimates).

- *Estimation of a proportion.* Survey conducted to estimate the proportion  $\theta$  of positive individuals in a population.

Random sample of size  $n$  with  $r$  positive.

$\Pr(a < \theta < b | r, n, A, K)$ , a conditional measure of the uncertainty about the event that  $\theta$  belongs to  $[a, b]$  *given* assumptions  $A$ , initial knowledge  $K$  and data  $\{r, n\}$ .

- *Measurement of a physical constant.* Measuring the unknown value of physical constant  $\mu$ , with data  $\boldsymbol{x} = \{x_1, \dots, x_n\}$ , considered to be measurements of  $\mu$  subject to error. Desired to find  $\Pr(a < \mu < b \mid x_1, \dots, x_n, A, K)$ , the *probability* that the unknown value of  $\mu$  (fixed in nature, but unknown to the scientists) belongs to  $[a, b]$  given the information provided by the data  $\boldsymbol{x}$ , assumptions  $A$  made, and available knowledge  $K$ .
- The statistical model may include *nuisance* parameters, unknown quantities, which have to be eliminated in the statement of the final results. For instance, the precision of the measurements described by unknown standard deviation  $\sigma$  in a  $\mathbf{N}(x \mid \mu, \sigma)$  normal model
- Relevant scientific information may impose *restrictions* on the admissible values of the quantities of interest. These must be taken into account. For instance, in measuring the value of the gravitational field  $g$  in a laboratory, it is known that it must lie between  $9.7803 \text{ m/sec}^2$  (average value at the Equator) and  $9.8322 \text{ m/sec}^2$  (average value at the poles).

- *Future discrete observations.* Experiment counting the number  $r$  of times that an event  $E$  takes place in each of  $n$  replications. Desired to forecast the number of times  $r$  that  $E$  will take place in a future, similar situation,  $\Pr(r \mid r_1, \dots, r_n, A, K)$ . For instance, no accidents in each of  $n = 10$  consecutive months may yield  $\Pr(r = 0 \mid \mathbf{x}, A, K) = 0.953$ .
- *Future continuous observations.* Data  $\mathbf{x} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ . Desired to forecast the value of a future observation  $\mathbf{y}$ ,  $p(\mathbf{y} \mid \mathbf{x}, A, K)$ . For instance, from breaking strengths  $\mathbf{x} = \{y_1, \dots, y_n\}$  of  $n$  randomly chosen safety belt webbings, the engineer may find  $\Pr(y > y^* \mid \mathbf{x}, A, K) = 0.9987$ .
- *Regression.* Data set consists of pairs  $\mathbf{x} = \{(\mathbf{y}_1, \mathbf{v}_1), \dots, (\mathbf{y}_n, \mathbf{v}_n)\}$  of quantity  $\mathbf{y}_j$  observed in conditions  $\mathbf{v}_j$ . Desired to forecast the value of  $\mathbf{y}$  in conditions  $\mathbf{v}$ ,  $p(\mathbf{y} \mid \mathbf{v}, \mathbf{x}, A, K)$ . For instance,  $y$  contamination levels,  $v$  wind speed from source; environment authorities interested in  $\Pr(y > y^* \mid v, \mathbf{x}, A, K)$

## 2. Basics of Bayesian Analysis

### 2.1. Parametric Inference

- *Bayes Theorem*

- Let  $\mathcal{M} = \{p(\mathbf{x} | \boldsymbol{\theta}), \mathbf{x} \in \mathcal{X}, \boldsymbol{\theta} \in \Theta\}$  be an statistical model, let  $\pi(\boldsymbol{\theta} | K)$  be a probability density for  $\boldsymbol{\theta}$  given prior knowledge  $K$  and let  $\mathbf{x}$  be some available data.

$$\pi(\boldsymbol{\theta} | \mathbf{x}, \mathcal{M}, K) = \frac{p(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | K)}{\int_{\Theta} p(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | K) d\boldsymbol{\theta}},$$

encapsulates all information about  $\boldsymbol{\theta}$  given data and prior knowledge.

- Simplifying notation, Bayes' theorem may be expressed as

$$\pi(\boldsymbol{\theta} | \mathbf{x}) \propto p(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) :$$

*The posterior is proportional to the likelihood times the prior.* The missing proportionality constant  $[\int_{\Theta} p(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}]^{-1}$  may be deduced from the fact that  $\pi(\boldsymbol{\theta} | \mathbf{x})$  must integrate to one. To identify a posterior distribution it suffices to identify a *kernel*  $k(\boldsymbol{\theta}, \mathbf{x})$  such that  $\pi(\boldsymbol{\theta} | \mathbf{x}) = c(\mathbf{x}) k(\boldsymbol{\theta}, \mathbf{x})$ . This is a very common technique.

- *Bayesian Inference with a Finite Parameter Space*

- Model  $\{p(\mathbf{x} | \theta_i), \mathbf{x} \in \mathcal{X}, \theta_i \in \Theta\}$ , with  $\Theta = \{\theta_1, \dots, \theta_m\}$ , so that  $\theta$  may only take a *finite* number  $m$  of different values. Using the finite form of Bayes' theorem,

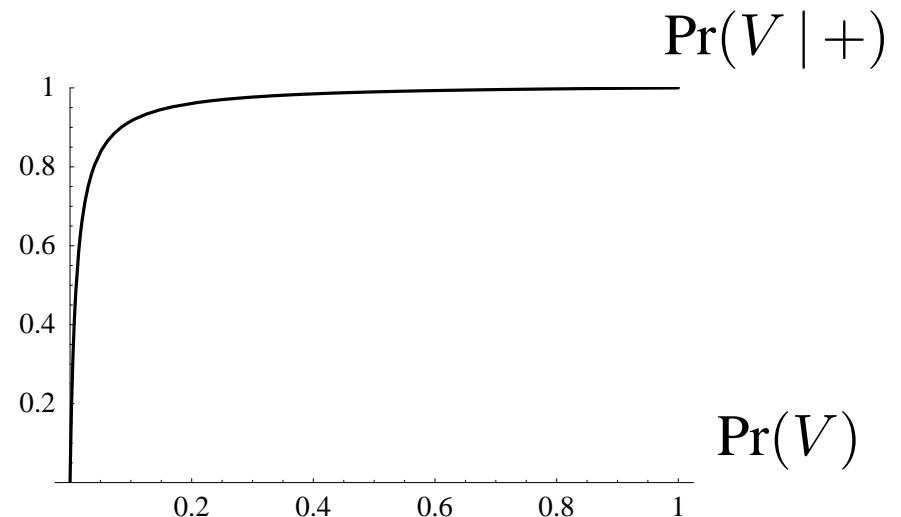
$$\Pr(\theta_i | \mathbf{x}) = \frac{p(\mathbf{x} | \theta_i) \Pr(\theta_i)}{\sum_{j=1}^m p(\mathbf{x} | \theta_j) \Pr(\theta_j)}, \quad i = 1, \dots, m.$$

- *Example: Probabilistic diagnosis.* A test to detect a virus, is known from laboratory research to give a positive result in 98% of the infected people and in 1% of the non-infected. The posterior probability that a person who tested positive is infected is

$$\Pr(V | +) = \frac{0.98 p}{0.98 p + 0.01 (1 - p)}$$

as a function of  $p = \Pr(V)$ .

- Notice sensitivity of posterior  $\Pr(V | +)$  to changes in the prior  $p = \Pr(V)$ .



- *Example: Inference about a binomial parameter*

□ Let data  $\mathbf{x}$  be  $n$  Bernoulli observations with parameter  $\theta$  which contain  $r$  positives, so that  $p(\mathbf{x} | \theta, n) = \theta^r (1 - \theta)^{n-r}$ .

□ If  $\pi(\theta) = \text{Be}(\theta | \alpha, \beta)$ , then

$$\pi(\theta | \mathbf{x}) \propto \theta^{r+\alpha-1} (1 - \theta)^{n-r+\beta-1}$$

kernel of  $\text{Be}(\theta | r + \alpha, n - r + \beta)$ .

□ Prior information ( $K$ )

$P(0.4 < \theta < 0.6) = 0.95$ ,  
and symmetric, yields  $\alpha = \beta = 47$ ;

□ **No prior information**  $\alpha = \beta = 1/2$

□  $n = 1500, r = 720$

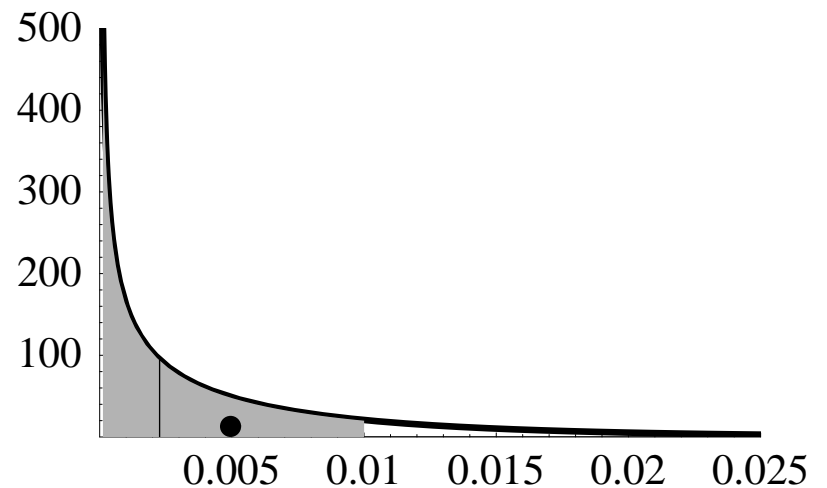
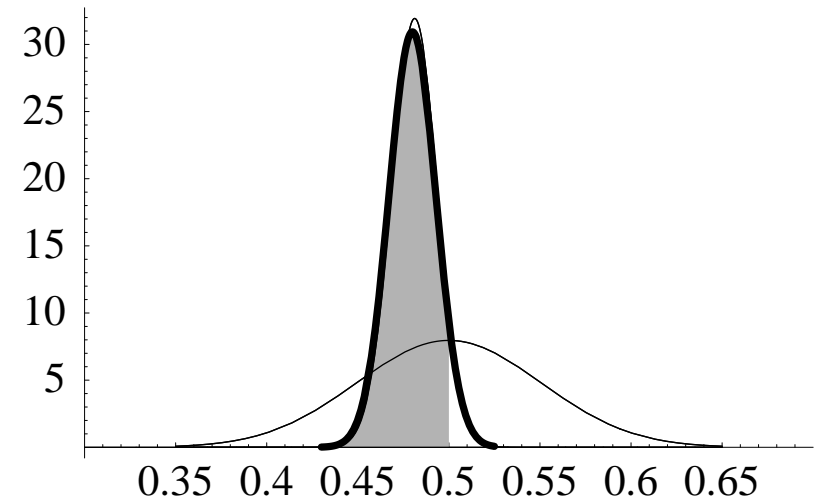
$$P(\theta < 0.5 | \mathbf{x}, K) = 0.933$$

$$P(\theta < 0.5 | \mathbf{x}) = 0.934$$

□  $n = 100, r = 0$

$$P(\theta < 0.01 | \mathbf{x}) = 0.844$$

Notice:  $\hat{\theta} = 0$ , but  $\text{Me}[\theta | \mathbf{x}] = 0.0023$





- *Sufficiency*

- Given a model  $p(\mathbf{x} | \boldsymbol{\theta})$ , a function of the data  $\mathbf{t} = \mathbf{t}(\mathbf{x})$ , is a *sufficient* statistic if it encapsulates all information about  $\boldsymbol{\theta}$  available in  $\mathbf{x}$ .
- Formally,  $\mathbf{t} = \mathbf{t}(\mathbf{x})$  is *sufficient* if (and only if), for any prior  $\pi(\boldsymbol{\theta})$   $\pi(\boldsymbol{\theta} | \mathbf{x}) = \pi(\boldsymbol{\theta} | \mathbf{t})$ . Hence,  $\pi(\boldsymbol{\theta} | \mathbf{x}) = \pi(\boldsymbol{\theta} | \mathbf{t}) \propto p(\mathbf{t} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})$ .
- This is equivalent to the frequentist definition; thus  $\mathbf{t} = \mathbf{t}(\mathbf{x})$  is sufficient iff  $p(\mathbf{x} | \boldsymbol{\theta}) = f(\boldsymbol{\theta}, \mathbf{t})g(\mathbf{x})$ .
- A sufficient statistic always exists, for  $\mathbf{t}(\mathbf{x}) = \mathbf{x}$  is obviously sufficient  
 A much simpler sufficient statistic, with fixed dimensionality independent of the sample size, often exists.  
 This is case whenever the statistical model belongs to the *generalized exponential family*, which includes many of the more frequently used statistical models.
- In contrast to frequentist statistics, Bayesian methods are independent on the possible existence of a sufficient statistic of fixed dimensionality.  
 For instance, if data come from an **Student** distribution, there is *no sufficient statistic* of fixed dimensionality: *all data are needed*.

- *Example: Inference from Cauchy observations*

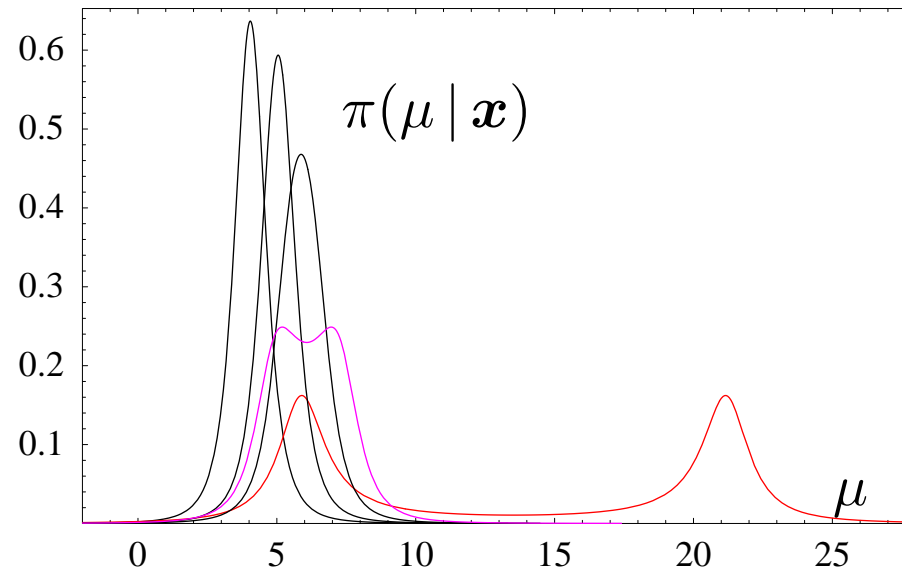
- Data  $\mathbf{x} = \{x_1, \dots, x_n\}$  random from  $\text{Ca}(x | \mu, 1) = \text{St}(x | \mu, 1, 1)$ .
- Objective reference prior for the location parameter  $\mu$  is  $\pi(\mu) = 1$ .
- By Bayes' theorem,

$$\pi(\mu | \mathbf{x}) \propto \prod_{j=1}^n \text{Ca}(x_j | \mu, 1) \pi(\mu) \propto \prod_{j=1}^n \frac{1}{1 + (x_j - \mu)^2}.$$

Proportionality constant easily obtained by numerical integration.

- Five samples of size  $n = 2$  simulated from  $\text{Ca}(x | 5, 1)$ .

$x_1$	$x_2$
4.034	4.054
21.220	5.831
5.272	6.475
4.776	5.317
7.409	4.743



- *Improper prior functions*

- Objective Bayesian methods often use functions which play the role of prior distributions but are *not* probability distributions.
- An *improper prior function* is a non-negative function  $\pi(\boldsymbol{\theta})$  such that  $\int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}$  is not finite.

The Cauchy example uses the improper prior function  $\pi(\mu) = 1, \mu \in \mathfrak{R}$ .

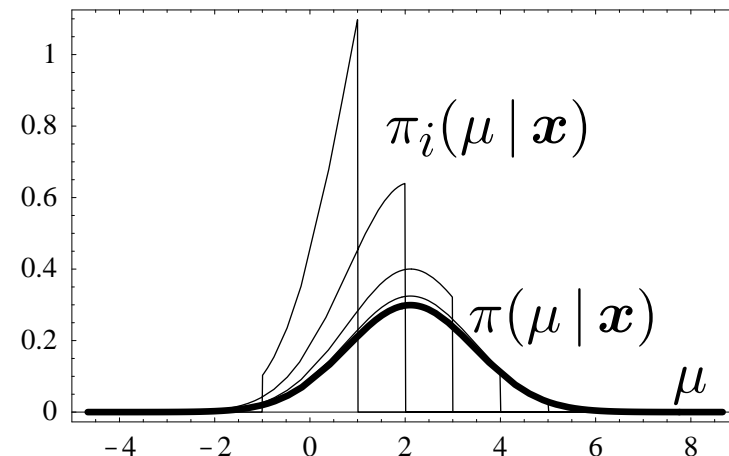
- $\pi(\boldsymbol{\theta})$  is an improper prior function,  $\{\Theta_i\}_{i=1}^{\infty}$  an increasing sequence approximating  $\Theta$ , such that  $\int_{\Theta_i} \pi(\boldsymbol{\theta}) < \infty$ , and  $\{\pi_i(\boldsymbol{\theta})\}_{i=1}^{\infty}$  the proper priors obtained by *renormalizing*  $\pi(\boldsymbol{\theta})$  within the  $\Theta_i$ 's.
- For any data  $\boldsymbol{x}$  with likelihood  $p(\boldsymbol{x} | \boldsymbol{\theta})$ , the sequence of posteriors  $\pi_i(\boldsymbol{\theta} | \boldsymbol{x})$  converges intrinsically to  $\pi(\boldsymbol{\theta} | \boldsymbol{x}) \propto p(\boldsymbol{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})$ .

- Normal data,  $\sigma$  known,  $\pi(\mu) = 1$ .

$$\begin{aligned} \pi(\mu | \boldsymbol{x}) &\propto p(\boldsymbol{x} | \mu, \sigma) \pi(\mu) \\ &\propto \exp\left[-\frac{n}{2\sigma^2}(\bar{x} - \mu)^2\right] \end{aligned}$$

$$\pi(\mu | \boldsymbol{x}) = \mathbf{N}(\mu | \bar{x}, \sigma/\sqrt{n})$$

Example:  $n = 9, \bar{x} = 2.11, \sigma = 4$



- *Sequential updating*

- Prior and posterior are terms *relative* to a set of data.
- If data  $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are sequentially presented, the final result will be the same whether data are globally or sequentially processed.

$$\pi(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_{i+1}) \propto p(\mathbf{x}_{i+1} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_i).$$

The “posterior” at a given stage becomes the “prior” at the next.

- **Typically** (but not always), the new **posterior**,  $\pi(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_{i+1})$ , is **more concentrated** around the true value than  $\pi(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_i)$ .

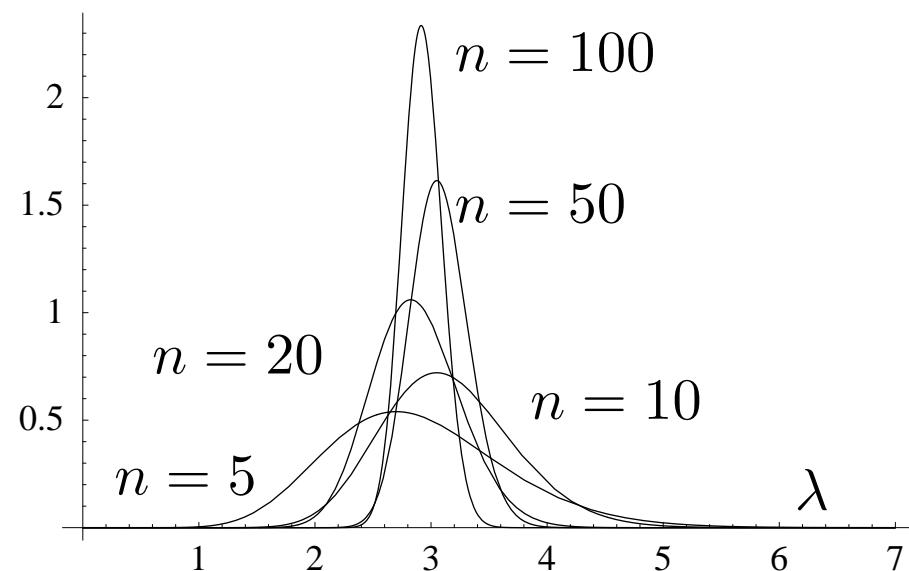
- Posteriors  $\pi(\lambda | x_1, \dots, x_i)$   
from increasingly large  
simulated data from Poisson

$P_n(x | \lambda)$ , with  $\lambda = 3$

$$\pi(\lambda | x_1, \dots, x_i)$$

$$= \text{Ga}(\lambda | r_i + 1/2, i)$$

$$r_i = \sum_{j=1}^i x_j$$



- *Nuisance parameters*

- In general the *vector of interest* is not the whole parameter vector  $\theta$ , but some function  $\phi = \phi(\theta)$  of possibly lower dimension.
- By Bayes' theorem  $\pi(\theta | x) \propto p(x | \theta) \pi(\theta)$ . Let  $\omega = \omega(\theta) \in \Omega$  be another function of  $\theta$  such that  $\psi = \{\phi, \omega\}$  is a bijection of  $\theta$ , and let  $J(\psi) = (\partial\theta / \partial\psi)$  be the Jacobian of the inverse function  $\psi = \psi(\theta)$ .

From probability theory,  $\pi(\psi | x) = |J(\psi)| [\pi(\theta | x)]_{\theta=\theta(\psi)}$

and  $\pi(\phi | x) = \int_{\Omega} \pi(\phi, \omega | x) d\omega$ .

- Any valid conclusion on  $\phi$  will be contained in  $\pi(\phi | x)$ .
- Particular case: *marginal posteriors*

Often model directly expressed in terms of vector of interest  $\phi$ , and vector of nuisance parameters  $\omega$ ,  $p(x | \theta) = p(x | \phi, \omega)$ .

Specify the prior  $\pi(\theta) = \pi(\phi) \pi(\omega | \phi)$

Get the joint posterior  $\pi(\phi, \omega | x) \propto p(x | \phi, \omega) \pi(\omega | \phi) \pi(\phi)$

Integrate out  $\omega$ ,  $\pi(\phi | x) \propto \pi(\phi) \int_{\Omega} p(x | \phi, \omega) \pi(\omega | \phi) d\omega$

- *Example: Inferences about a Normal mean*

□ Data  $\mathbf{x} = \{x_1, \dots, x_n\}$  random from  $\mathbf{N}(x | \mu, \sigma)$ . Likelihood function  $p(\mathbf{x} | \mu, \sigma) \propto \sigma^{-n} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)]$ , with  $n\bar{x} = \sum_i x_i$ , and  $ns^2 = \sum_i (x_i - \bar{x})^2$ .

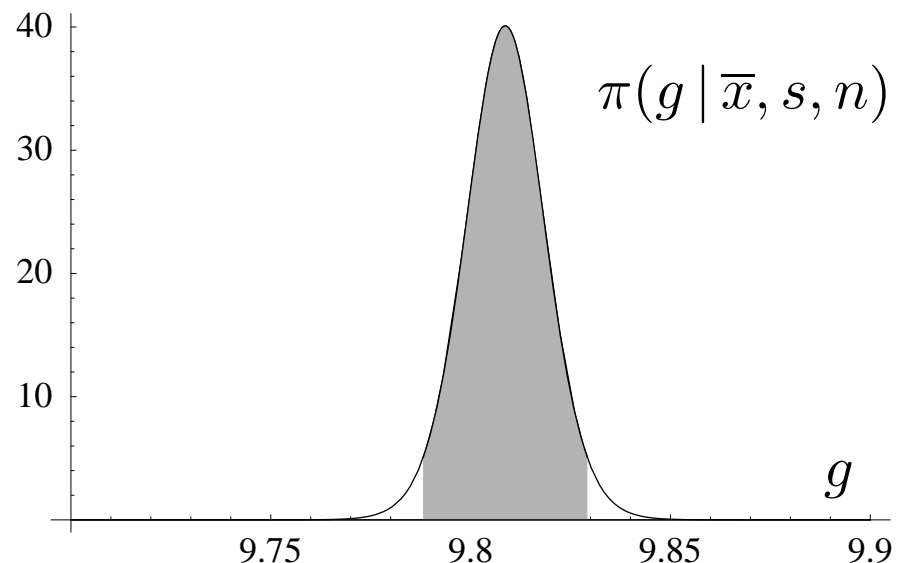
□ Objective prior is uniform in both  $\mu$  and  $\log(\sigma)$ , i.e.,  $\pi(\mu, \sigma) = \sigma^{-1}$ . Joint posterior  $\pi(\mu, \sigma | \mathbf{x}) \propto \sigma^{-(n+1)} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)]$ .

□ Marginal posterior  $\pi(\mu | \mathbf{x}) \propto \int_0^\infty \pi(\mu, \sigma | \mathbf{x}) d\sigma \propto [s^2 + (\bar{x} - \mu)^2]^{-n/2}$ , kernel of the Student density  $\text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$

□ Classroom experiment to measure gravity  $g$  yields  $\bar{x} = 9.8087$ ,  $s = 0.0428$  with  $n = 20$  measures.

$$\begin{aligned} \pi(g | \bar{x}, s, n) \\ = \text{St}(g | 9.8087, 0.0098, 19) \end{aligned}$$

$$\begin{aligned} \Pr(9.788 < g < 9.829 | \mathbf{x}) \\ = 0.95 \quad (\text{shaded area}) \end{aligned}$$



- *Restricted parameter space*

□ Range of values of  $\theta$  restricted by contextual considerations.

If  $\theta$  known to belong to  $\Theta_c \subset \Theta$ ,  $\pi(\theta) > 0$  iff  $\theta \in \Theta_c$

By Bayes' theorem,

$$\pi(\theta | \mathbf{x}, \theta \in \Theta_c) = \begin{cases} \frac{\pi(\theta | \mathbf{x})}{\int_{\Theta_c} \pi(\theta | \mathbf{x}) d\theta}, & \text{if } \theta \in \Theta_c \\ 0 & \text{otherwise} \end{cases}$$

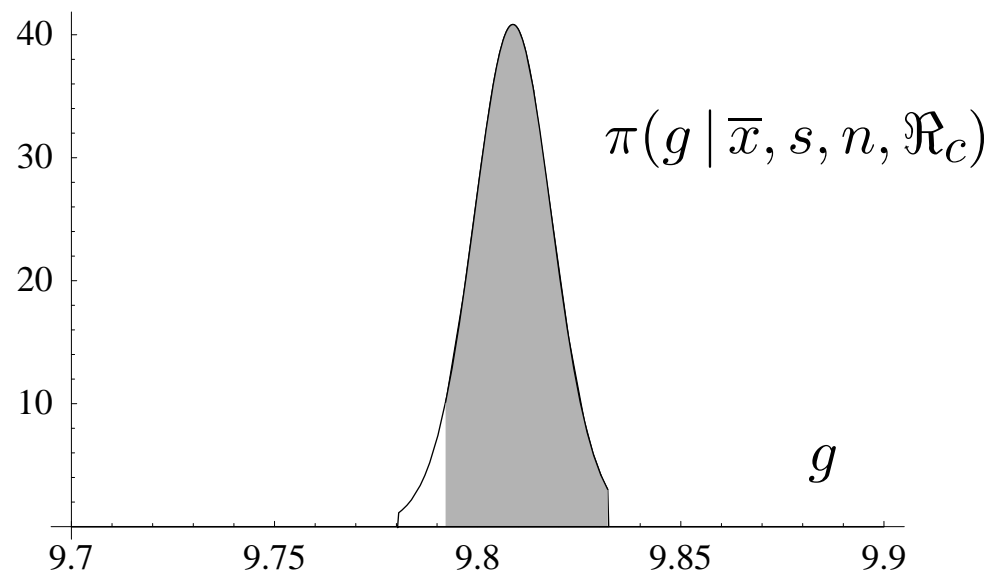
□ To incorporate a restriction, it suffices to *renormalize* the unrestricted posterior distribution to the set  $\Theta_c \subset \Theta$  of admissible parameter values.

□ Classroom experiment to measure gravity  $g$  with restriction to lie between

$g_0 = 9.7803$  (equator)

$g_1 = 9.8322$  (poles).

$\Pr(9.7921 < g < 9.8322 | \mathbf{x})$   
 $= 0.95$  (shaded area)



- *Asymptotic behaviour, discrete case*

- If the parameter space  $\Theta = \{\theta_1, \theta_2, \dots\}$  is *countable* and  
The true parameter value  $\theta_t$  is *distinguishable* from the others, *i.e.*,  
 $\delta\{p(\mathbf{x} | \theta_t), p(\mathbf{x} | \theta_i)\} > 0, i \neq t,$

$$\lim_{n \rightarrow \infty} \pi(\theta_t | \mathbf{x}_1, \dots, \mathbf{x}_n) = 1$$

$$\lim_{n \rightarrow \infty} \pi(\theta_i | \mathbf{x}_1, \dots, \mathbf{x}_n) = 0, \quad i \neq t$$

- To prove this, take logarithms in Bayes' theorem,  
define  $z_i = \log[p(\mathbf{x} | \theta_i)/p(\mathbf{x} | \theta_t)],$   
and use the strong law of large numbers on the  $n$   
i.i.d. random variables  $z_1, \dots, z_n.$
- For instance, in probabilistic diagnosis the posterior probability of the true disease converges to one as new relevant information accumulates, *provided* the model distinguishes the probabilistic behaviour of data under the true disease from its behaviour under the other alternatives.



- *Asymptotic behaviour, continuous case*

- If the parameter  $\theta$  is *one-dimensional and continuous*, so that  $\Theta \subset \mathfrak{R}$ , and the model  $\{p(\mathbf{x} \mid \theta), \mathbf{x} \in \mathcal{X}\}$  is *regular*: basically,
  - $\mathcal{X}$  does not depend on  $\theta$ ,
  - $p(\mathbf{x} \mid \theta)$  is twice differentiable with respect to  $\theta$
- Then, as  $n \rightarrow \infty$ ,  $\pi(\theta \mid \mathbf{x}_1, \dots, \mathbf{x}_n)$  converges intrinsically to a *normal* distribution with mean at the mle estimator  $\hat{\theta}$ , and with variance  $v(\mathbf{x}_1, \dots, \mathbf{x}_n, \hat{\theta})$ , where
 
$$v^{-1}(\mathbf{x}_1, \dots, \mathbf{x}_n, \hat{\theta}) = - \sum_{j=1}^n \frac{\partial^2}{\partial \theta^2} \log[p(\mathbf{x}_j \mid \theta)]$$
- To prove this, express is Bayes' theorem as
 
$$\pi(\theta \mid \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \exp[\log \pi(\theta) + \sum_{j=1}^n \log p(\mathbf{x}_j \mid \theta)],$$
 and expand  $\sum_{j=1}^n \log p(\mathbf{x}_j \mid \theta)$  about its maximum, the mle  $\hat{\theta}$
- The result is easily extended to the multivariate case  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_k\}$ , to obtain a limiting  $k$ -variate normal centered at  $\hat{\boldsymbol{\theta}}$ , and with a dispersion matrix  $\mathbf{V}(\mathbf{x}_1, \dots, \mathbf{x}_n, \hat{\boldsymbol{\theta}})$  which generalizes  $v(\mathbf{x}_1, \dots, \mathbf{x}_n, \hat{\theta})$ .

- *Asymptotic behaviour, continuous case. Simpler form*

- Using the strong law of large numbers on the sums above a simpler, less precise approximation is obtained:

- If the parameter  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_k\}$  is continuous, so that  $\Theta \subset \mathbb{R}^k$  and the model  $\{p(\mathbf{x} | \boldsymbol{\theta}), \mathbf{x} \in \mathcal{X}\}$  is *regular*, so that  $\mathcal{X}$  does not depend on  $\boldsymbol{\theta}$  and  $p(\mathbf{x} | \boldsymbol{\theta})$  is twice differentiable with respect to each of the  $\theta_i$ 's, then, as  $n \rightarrow \infty$ ,  $\pi(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_n)$  converges intrinsically to a *multivariate normal* distribution with mean the mle  $\hat{\boldsymbol{\theta}}$  and precision matrix (inverse of the dispersion or variance-covariance matrix)  $n \mathbf{F}(\hat{\boldsymbol{\theta}})$ , where  $\mathbf{F}(\boldsymbol{\theta})$  is Fisher's matrix, of general element

$$F_{ij}(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x} | \boldsymbol{\theta}} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(\mathbf{x} | \boldsymbol{\theta}) \right]$$

- The properties of the multivariate normal yield from this result the asymptotic forms for the *marginal* and the *conditional* posterior distributions of any subgroup of the  $\theta_j$ 's.

- In one dimension,  $\pi(\theta | \mathbf{x}_1, \dots, \mathbf{x}_n) \approx \mathbf{N}(\theta | \hat{\theta}, (nF(\theta))^{-1/2})$ , where  $F(\theta) = -\mathbb{E}_{\mathbf{x} | \theta} [\partial^2 \log p(\mathbf{x} | \theta) / \partial \theta^2]$

- *Example: Asymptotic approximation with Poisson data*

- Data  $\mathbf{x} = \{x_1, \dots, x_n\}$  random from  $\text{Pn}(x | \lambda) \propto e^{-\lambda} \lambda^x / x!$   
hence,  $p(\mathbf{x} | \lambda) \propto e^{-n\lambda} \lambda^r$ ,  $r = \sum_j x_j$ , and  $\hat{\lambda} = r/n$ .

Fisher's function is  $F(\lambda) = -\mathbf{E}_{\mathbf{x} | \lambda} \left[ \frac{\partial^2}{\partial \lambda^2} \log \text{Pn}(x | \lambda) \right] = \frac{1}{\lambda}$

- The objective prior function is  $\pi(\lambda) = F(\lambda)^{1/2} = \lambda^{-1/2}$

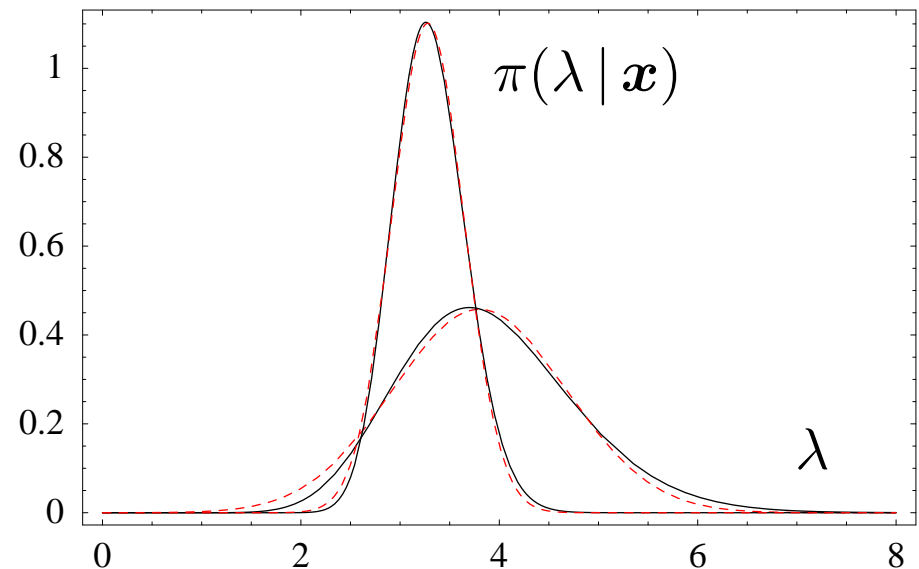
Hence  $\pi(\lambda | \mathbf{x}) \propto e^{-n\lambda} \lambda^{r-1/2}$

the kernel of  $\text{Ga}(\lambda | r + \frac{1}{2}, n)$

- The Normal approximation is

$$\begin{aligned} \pi(\lambda | \mathbf{x}) &\approx \mathbf{N}\{\lambda | \hat{\lambda}, (n F(\hat{\lambda}))^{-1/2}\} \\ &= \mathbf{N}\{\lambda | r/n, \sqrt{r}/n\} \end{aligned}$$

- Samples  $n = 5$  and  $n = 25$   
simulated from Poisson  $\lambda = 3$   
yielded  $r = 19$  and  $r = 82$



## 2.2. Reference Analysis

- *No Relevant Initial Information*

- Identify the mathematical form of a “noninformative” prior. One with *minimal effect, relative to the data, on the posterior distribution of the quantity of interest.*

- Intuitive basis:

Use *information theory* to measure the amount of information about the quantity of interest to be expected from data. This depends on prior knowledge: the more it is known, the less the amount of information the data may be expected to provide.

Define the *missing information* about the quantity of interest as that which infinite independent replications of the experiment could possibly provide.

Define the *reference prior* as that which *maximizes the missing information about the quantity of interest.*

- *Expected information from the data*

□ Given model  $\{p(\mathbf{x} | \theta), \mathbf{x} \in \mathcal{X}, \theta \in \Theta\}$ , the *amount of information*  $I^\theta\{\mathcal{X}, \pi(\theta)\}$  which may be expected to be provided by  $\mathbf{x}$ , about the value of  $\theta$  is defined by

$$I^\theta\{\mathcal{X}, \pi(\theta)\} = \delta\{p(\mathbf{x}, \theta), p(\mathbf{x})\pi(\theta)\},$$

the intrinsic discrepancy between the joint distribution  $p(\mathbf{x}, \theta)$  and the product of their marginals  $p(\mathbf{x})\pi(\theta)$ , which is the *intrinsic association* between the random quantities  $\mathbf{x}$  and  $\theta$ .

□ Consider  $I^\theta\{\mathcal{X}^k, \pi(\theta)\}$  the information about  $\theta$  which may be expected from  $k$  conditionally independent replications of the original setup.

As  $k \rightarrow \infty$ , this would provide any *missing information* about  $\theta$ . Hence, as  $k \rightarrow \infty$ , the functional  $I^\theta\{\mathcal{X}^k, \pi(\theta)\}$  will approach the missing information about  $\theta$  associated with the prior  $\pi(\theta)$ .

□ Let  $\pi_k(\theta)$  be the prior which maximizes  $I^\theta\{\mathcal{X}^k, \pi(\theta)\}$  in the class  $\mathcal{P}$  of strictly positive prior distributions compatible with accepted assumptions on the value of  $\theta$  (which be the class of *all* strictly positive priors).

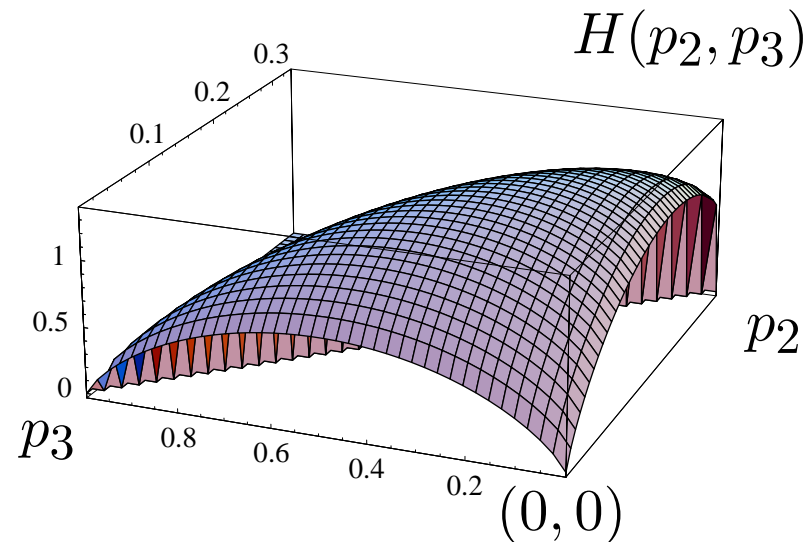
The *reference prior*  $\pi^*(\theta)$  is the limit as  $k \rightarrow \infty$  (in a sense to be made precise) of the sequence of priors  $\{\pi_k(\theta), k = 1, 2, \dots\}$ .

- *Reference priors in the finite case*

- If  $\theta$  may only take a *finite* number  $m$  of different values  $\{\theta_1, \dots, \theta_m\}$  and  $\pi(\theta) = \{p_1, \dots, p_m\}$ , with  $p_i = \Pr(\theta = \theta_i)$ , then  $\lim_{k \rightarrow \infty} I^\theta \{ \mathcal{X}^k, \pi(\theta) \} = H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log(p_i)$ , that is, the *entropy* of the prior distribution  $\{p_1, \dots, p_m\}$ .
- In the finite case, the reference prior is that with *maximum entropy* within the class  $\mathcal{P}$  of priors compatible with accepted assumptions. (cf. Statistical Physics)
- If, in particular,  $\mathcal{P}$  contains *all* priors over  $\{\theta_1, \dots, \theta_m\}$ , the reference prior is the *uniform* prior,  $\pi(\theta) = \{1/m, \dots, 1/m\}$ . (cf. Bayes-Laplace postulate of insufficient reason)

- Prior  $\{p_1, p_2, p_3, p_4\}$  in genetics problem where  $p_1 = 2p_2$ .

Reference prior is  $\{0.324, 0.162, 0.257, 0.257\}$



- *Reference priors in one-dimensional continuous case*

- Let  $\pi_k(\theta)$  be the prior which maximizes  $I^\theta\{\mathcal{X}^k, \pi(\theta)\}$  in the class  $\mathcal{P}$  of acceptable priors.

For any data  $\mathbf{x} \in \mathcal{X}$ , let  $\pi_k(\theta | \mathbf{x}) \propto p(\mathbf{x} | \theta) \pi_k(\theta)$  be the corresponding posterior.

- The *reference posterior density*  $\pi^*(\theta | \mathbf{x})$  is defined to be the intrinsic limit of the sequence  $\{\pi_k(\theta | \mathbf{x}), k = 1, 2, \dots\}$

A *reference prior function*  $\pi^*(\theta)$  is any positive function such that, for all  $\mathbf{x} \in \mathcal{X}$ ,  $\pi^*(\theta | \mathbf{x}) \propto p(\mathbf{x} | \theta) \pi^*(\theta)$ .

This is defined up to an (irrelevant) arbitrary constant.

- Let  $\mathbf{x}^{(k)} \in \mathcal{X}^k$  be the result of  $k$  independent replications of  $\mathbf{x} \in \mathcal{X}$ . The exact expression for  $\pi_k(\theta)$  (which may be obtained with calculus of variations) is

$$\pi_k(\theta) = \exp \left[ \mathbf{E}_{\mathbf{x}^{(k)} | \theta} \{ \log \pi_k(\theta | \mathbf{x}^{(k)}) \} \right]$$

- This formula may be used, by repeated simulation from  $p(\mathbf{x} | \theta)$  for different  $\theta$  values, to obtain a *numerical approximation* to the reference prior.

- *Reference priors under regularity conditions*

□ Let  $\tilde{\theta}_k = \tilde{\theta}(x^{(k)})$  be a consistent, asymptotically sufficient estimator of  $\theta$ . In regular problems this is often the case with the mle estimator  $\hat{\theta}$ .

The exact expression for  $\pi_k(\theta)$  then becomes, for large  $k$ ,

$$\square \pi_k(\theta) \approx \exp[\mathbf{E}_{\tilde{\theta}_k | \theta} \{\log \pi_k(\theta | \tilde{\theta}_k)\}]$$

As  $k \rightarrow \infty$  this converges to  $\pi_k(\theta | \tilde{\theta}_k)|_{\tilde{\theta}_k=\theta}$

□ Let  $\tilde{\theta}_k = \tilde{\theta}(x^{(k)})$  be a consistent, asymptotically sufficient estimator of  $\theta$ . Let  $\pi(\theta | \tilde{\theta}_k)$  be any asymptotic approximation to  $\pi(\theta | x^{(k)})$ , the posterior distribution of  $\theta$ .

Hence,  $\pi^*(\theta) = \pi(\theta | \tilde{\theta}_k)|_{\tilde{\theta}_k=\theta}$

□ Under regularity conditions, the posterior distribution of  $\theta$  is asymptotically Normal, with mean  $\hat{\theta}$  and precision  $n F(\hat{\theta})$ , where  $F(\theta) = -\mathbf{E}_{\mathbf{x} | \theta} [\partial^2 \log p(\mathbf{x} | \theta) / \partial \theta^2]$  is Fisher's information function.

Hence,  $\pi^*(\theta) = F(\theta)^{1/2}$  (Jeffreys' rule).



- *One nuisance parameter*

□ *Two parameters*: reduce the problem to a *sequential* application of the one parameter case. Probability model is  $\{p(\mathbf{x} | \theta, \lambda, \theta \in \Theta, \lambda \in \Lambda)\}$  and a  $\theta$ -reference prior  $\pi_{\theta}^*(\theta, \lambda)$  is required. Two steps:

(i) Conditional on  $\theta$ ,  $p(\mathbf{x} | \theta, \lambda)$  only depends on  $\lambda$ , and it is possible to obtain the *conditional* reference prior  $\pi^*(\lambda | \theta)$ .

(ii) If  $\pi^*(\lambda | \theta)$  is proper, integrate out  $\lambda$  to get the one-parameter model  $p(\mathbf{x} | \theta) = \int_{\Lambda} p(\mathbf{x} | \theta, \lambda) \pi^*(\lambda | \theta) d\lambda$ , and use the one-parameter solution to obtain  $\pi^*(\theta)$ .

The  $\theta$ -reference prior is then  $\pi_{\theta}^*(\theta, \lambda) = \pi^*(\lambda | \theta) \pi^*(\theta)$ .

The required reference posterior is  $\pi^*(\theta | \mathbf{x}) \propto p(\mathbf{x} | \theta) \pi^*(\theta)$ .

□ If  $\pi^*(\lambda | \theta)$  is an *improper* prior function, proceed within an increasing sequence  $\{\Lambda_i\}$  over which  $\pi^*(\lambda | \theta)$  is integrable and, for given data  $\mathbf{x}$ , obtain the corresponding sequence of reference posteriors  $\{\pi_i^*(\theta | \mathbf{x})\}$ .

The required reference posterior  $\pi^*(\theta | \mathbf{x})$  is their intrinsic limit.

A  $\theta$ -reference prior is any positive function such that, for any data  $\mathbf{x}$ ,  $\pi^*(\theta | \mathbf{x}) \propto \int_{\Lambda} p(\mathbf{x} | \theta, \lambda) \pi_{\theta}^*(\theta, \lambda) d\lambda$ .

- *The regular two-parameter continuous case*

- Model  $p(\mathbf{x} \mid \theta, \lambda)$ . If the joint posterior of  $(\theta, \lambda)$  is asymptotically normal, the  $\theta$ -reference prior may be derived in terms of the corresponding Fisher's information matrix,  $\mathbf{F}(\theta, \lambda)$ .

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

The  $\theta$ -reference prior is  $\pi_{\theta}^*(\theta, \lambda) = \pi^*(\lambda \mid \theta) \pi^*(\theta)$ , where

$\pi^*(\lambda \mid \theta) \propto F_{\lambda\lambda}^{1/2}(\theta, \lambda)$ ,  $\lambda \in \Lambda$ , and, if  $\pi^*(\lambda \mid \theta)$  is proper,

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

- If  $\pi^*(\lambda \mid \theta)$  is not proper, integrations are performed within an approximating sequence  $\{\Lambda_i\}$  to obtain a sequence  $\{\pi_i^*(\lambda \mid \theta) \pi_i^*(\theta)\}$ , and the  $\theta$ -reference prior  $\pi_{\theta}^*(\theta, \lambda)$  is defined as its intrinsic limit.
- Even if  $\pi^*(\lambda \mid \theta)$  is improper, if  $\theta$  and  $\lambda$  are variation independent,  $S_{\theta\theta}^{-1/2}(\theta, \lambda) \propto f_{\theta}(\theta) g_{\theta}(\lambda)$ , and  $F_{\lambda\lambda}^{1/2}(\theta, \lambda) \propto f_{\lambda}(\theta) g_{\lambda}(\lambda)$ ,  
Then  $\pi_{\theta}^*(\theta, \lambda) = f_{\theta}(\theta) g_{\lambda}(\lambda)$ .

- *Examples: Inference on normal parameters*

□ The information matrix for the normal model  $N(x | \mu, \sigma)$  is

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

Since  $\mu$  and  $\sigma$  are variation independent, and both  $F_{\sigma\sigma}$  and  $S_{\mu\mu}$  factorize,

$$\pi^*(\sigma | \mu) \propto F_{\sigma\sigma}^{1/2} \propto \sigma^{-1}, \quad \pi^*(\mu) \propto S_{\mu\mu}^{-1/2} \propto 1.$$

The  $\mu$ -reference prior, as anticipated, is

$$\pi_{\mu}^*(\mu, \sigma) = \pi^*(\sigma | \mu) \pi^*(\mu) = \sigma^{-1},$$

*i.e.*, uniform on both  $\mu$  and  $\log \sigma$

□ Since  $\mathbf{F}(\mu, \sigma)$  is diagonal the  $\sigma$ -reference prior is

$$\pi_{\sigma}^*(\mu, \sigma) = \pi^*(\mu | \sigma) \pi^*(\sigma) = \sigma^{-1}, \text{ the same as } \pi_{\mu}^*(\mu, \sigma) = \pi_{\sigma}^*(\mu, \sigma).$$

□ In fact, it may be shown that, for location-scale models,

$$p(x | \mu, \sigma) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right),$$

the reference prior for the location and scale parameters are always

$$\pi_{\mu}^*(\mu, \sigma) = \pi_{\sigma}^*(\mu, \sigma) = \sigma^{-1}.$$

- Within any given model  $p(\mathbf{x} | \boldsymbol{\theta})$  the  $\phi$ -reference prior  $\pi_{\phi}^*(\boldsymbol{\theta})$  maximizes the missing information about  $\phi = \phi(\boldsymbol{\theta})$  and, in multiparameter problems, that prior *may change with the quantity of interest*  $\phi$ .
- For instance, within a normal  $\mathbf{N}(x | \mu, \sigma)$  model, let the *standardized mean*  $\phi = \mu/\sigma$  be the quantity of interest.

Fisher's information matrix in terms of the parameters  $\phi$  and  $\sigma$  is

$\mathbf{F}(\phi, \sigma) = \mathbf{J}^t \mathbf{F}(\mu, \sigma) \mathbf{J}$ , where  $\mathbf{J} = (\partial(\mu, \sigma)/\partial(\phi, \sigma))$  is the Jacobian of the inverse transformation; this yields

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

with  $F_{\sigma\sigma}^{1/2} \propto \sigma^{-1}$ , and  $S_{\phi\phi}^{-1/2} \propto (1 + \phi^2/2)^{-1/2}$ .

- The  $\phi$ -reference prior is,  $\pi_{\phi}^*(\phi, \sigma) = (1 + \phi^2/2)^{-1/2} \sigma^{-1}$ .

In the original parametrization,  $\pi_{\phi}^*(\mu, \sigma) = (1 + (\mu/\sigma)^2/2)^{-1/2} \sigma^{-2}$ , which is different from  $\pi_{\mu}^*(\mu, \sigma) = \pi_{\sigma}^*(\mu, \sigma)$ .

This prior is shown to lead to a reference posterior for  $\phi$  with *consistent marginalization properties*.

- *Many parameters*

- The reference algorithm generalizes to any number of parameters. If the model is  $p(\mathbf{x} | \boldsymbol{\theta}) = p(\mathbf{x} | \theta_1, \dots, \theta_m)$ , a joint reference prior  $\pi^*(\phi_m | \phi_{m-1}, \dots, \phi_1) \times \dots \times \pi^*(\phi_2 | \phi_1) \times \pi^*(\phi_1)$  may sequentially be obtained for each *ordered parametrization*,  $\{\phi_1(\boldsymbol{\theta}), \dots, \phi_m(\boldsymbol{\theta})\}$ . Reference priors are *invariant* under reparametrization of the  $\phi_i(\boldsymbol{\theta})$ 's.
- The choice of the ordered parametrization  $\{\phi_1, \dots, \phi_m\}$  describes the particular prior required, namely that which *sequentially* maximizes the missing information about each of the  $\phi_i$ 's, conditional on  $\{\phi_1, \dots, \phi_{i-1}\}$ , for  $i = m, m - 1, \dots, 1$ .
- Example: *Stein's paradox*. Data random from a  $m$ -variate normal  $N_m(\mathbf{x} | \boldsymbol{\mu}, \mathbf{I})$ . The reference prior function for any permutation of the  $\mu_i$ 's is uniform, and leads to appropriate posterior distributions for any of the  $\mu_i$ 's, but cannot be used if the quantity of interest is  $\theta = \sum_i \mu_i^2$ , the distance of  $\boldsymbol{\mu}$  to the origin.

The reference prior for  $\{\theta, \lambda_1, \dots, \lambda_{m-1}\}$  produces, for any choice of the  $\lambda_i$ 's, an appropriate the reference posterior for  $\theta$ .

## 2.3. Inference Summaries

- *Summarizing the posterior distribution*
  - *The* Bayesian final *outcome* of a problem of inference about any unknown quantity  $\theta$  *is* precisely the *posterior density*  $\pi(\theta | \mathbf{x}, C)$ .
  - Bayesian inference may be described as the problem of stating a probability distribution for the quantity of interest encapsulating all available information about its value.
  - In one or two dimensions, a *graph of the posterior probability density* of the quantity of interest conveys an intuitive summary of the main conclusions. This is greatly appreciated by users, and is an important asset of Bayesian methods.
  - However, graphical methods not easily extend to more than two dimensions and elementary *quantitative* conclusions are often required.

The simplest forms to *summarize* the information contained in the posterior distribution are closely related to the conventional concepts of point estimation and interval estimation.

- *Point Estimation: Posterior mean and posterior mode*

- It is often required to provide point estimates of relevant quantities. Bayesian point estimation is best described as a *decision problem* where one has to *choose* a particular value  $\tilde{\theta}$  as an approximate proxy for the actual, unknown value of  $\theta$ .
- Intuitively, any location measure of the posterior density  $\pi(\theta | \mathbf{x})$  may be used as a point estimator. When they exist, either  $E[\theta | \mathbf{x}] = \int_{\Theta} \theta \pi(\theta | \mathbf{x}) d\theta$  (*posterior mean*), or  $\text{Mo}[\theta | \mathbf{x}] = \arg \sup_{\theta \in \Theta} \pi(\theta | \mathbf{x})$  (*posterior mode*) are often regarded as natural choices.
- *Lack of invariance*. Neither the posterior mean nor the posterior mode are invariant under reparametrization. The point estimator  $\tilde{\psi}$  of a bijection  $\psi = \psi(\theta)$  of  $\theta$  will generally not be equal to  $\psi(\tilde{\theta})$ .

In pure “inferential” applications, where one is requested to provide a point estimate of the vector of interest without an specific application in mind, it is difficult to justify a non-invariant solution:

The best estimate of, say,  $\phi = \log(\theta)$  should be  $\phi^* = \log(\theta^*)$ .

- *Point Estimation: Posterior median*

- A summary of a multivariate density  $\pi(\boldsymbol{\theta} | \mathbf{x})$ , where  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_k\}$ , should contain summaries of:
  - (i) each of the marginal densities  $\pi(\theta_i | \mathbf{x})$ ,
  - (ii) the densities  $\pi(\phi | \mathbf{x})$  of other functions of interest  $\phi = \phi(\boldsymbol{\theta})$ .

- In *one-dimensional continuous* problems the *posterior median*, is easily defined and computed as

$$\text{Me}[\theta | \mathbf{x}] = q; \quad \Pr[\theta \leq q | \mathbf{x}] = \int_{\{\theta \leq q\}} \pi(\theta | \mathbf{x}) d\theta = 1/2$$

The one-dimensional posterior median has many attractive properties:

- (i) it is *invariant* under bijections,  $\text{Me}[\phi(\theta) | \mathbf{x}] = \phi(\text{Me}[\theta | \mathbf{x}])$ .
  - (ii) it *exists* and it is *unique* under very wide conditions
  - (iii) it is rather *robust* under moderate perturbations of the data.
- The posterior median is often considered to be the best ‘automatic’ Bayesian point estimator in one-dimensional continuous problems.
  - The posterior median is not easily used to a multivariate setting. The natural extension of its definition produces *surfaces* (not points).

General invariant multivariate definitions of point estimators is possible using Bayesian *decision theory*



- *General Credible Regions*

- To describe  $\pi(\boldsymbol{\theta} | \boldsymbol{x})$  it is often convenient to quote regions  $\Theta_p \subset \Theta$  of given probability content  $p$  under  $\pi(\boldsymbol{\theta} | \boldsymbol{x})$ . This is the intuitive basis of graphical representations like boxplots.
- A subset  $\Theta_p$  of the parameter space  $\Theta$  such that
 
$$\int_{\Theta_p} \pi(\boldsymbol{\theta} | \boldsymbol{x}) d\boldsymbol{\theta} = p, \quad \text{so that } \Pr(\boldsymbol{\theta} \in \Theta_p | \boldsymbol{x}) = p,$$
 is a *posterior  $p$ -credible region* for  $\boldsymbol{\theta}$ .
- A credible region is invariant under reparametrization:  
If  $\Theta_p$  is  $p$ -credible for  $\boldsymbol{\theta}$ ,  $\phi(\Theta_p)$  is a  $p$ -credible for  $\phi = \phi(\boldsymbol{\theta})$ .
- For any given  $p$  there are generally infinitely many credible regions.  
Credible regions may be 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.
- HPD regions are *not invariant*: the image  $\phi(\Theta_p)$  of an HPD region  $\Theta_p$  will be a credible region for  $\phi$ , but will not generally be HPD.  
There is no reason to restrict attention to HPD credible regions.

- *Credible Intervals*

- In *one-dimensional continuous* problems, posterior quantiles are often used to derive credible intervals.
- If  $\theta_q = Q_q[\theta | \mathbf{x}]$  is the  $q$ -quantile of the posterior distribution of  $\theta$ , the interval  $\Theta_p = \{\theta; \theta \leq \theta_p\}$  is a  $p$ -credible region, and it is invariant under reparametrization.
- *Equal-tailed*  $p$ -credible intervals of the form  $\Theta_p = \{\theta; \theta_{(1-p)/2} \leq \theta \leq \theta_{(1+p)/2}\}$  are typically unique, and they invariant under reparametrization.
- Example: Model  $N(x | \mu, \sigma)$ . *Credible intervals for the normal mean.* The reference posterior for  $\mu$  is  $\pi(\mu | \mathbf{x}) = \text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$ . Hence the reference *posterior* distribution of  $\tau = \sqrt{n-1}(\mu - \bar{x})/s$ , *a function of  $\mu$* , is  $\pi(\tau | \bar{x}, s, n) = \text{St}(\tau | 0, 1, n-1)$ .

Thus, the equal-tailed  $p$ -credible intervals for  $\mu$  are

$$\{\mu; \mu \in \bar{x} \pm q_{n-1}^{(1-p)/2} s/\sqrt{n-1}\},$$

where  $q_{n-1}^{(1-p)/2}$  is the  $(1-p)/2$  quantile of a standard Student density with  $n-1$  degrees of freedom.

- *Calibration*

- In the normal example above, the expression  $t = \sqrt{n-1}(\mu - \bar{x})/s$  may *also* be analyzed, for fixed  $\mu$ , as a *function of the data*.

The fact that the *sampling* distribution of the statistic  $t = t(\bar{x}, s | \mu, n)$  is *also* an standard Student  $p(t | \mu, n) = \text{St}(t | 0, 1, n-1)$  with the same degrees of freedom implies that, in this example, objective Bayesian credible intervals are *also* be *exact* frequentist confidence intervals.

- *Exact numerical agreement* between Bayesian credible intervals and frequentist confidence intervals is the *exception, not the norm*.
- For *large samples*, convergence to normality implies *approximate numerical agreement*. This provides a frequentist *calibration* to objective Bayesian methods.
- Exact numerical *agreement* is obviously *impossible when the data are discrete*: Precise (non randomized) frequentist confidence intervals do not exist in that case for most confidence levels.

The computation of Bayesian credible regions for continuous parameters is however *precisely the same* whether the data are *discrete or continuous*.

## 2.4. Prediction

- *Posterior predictive distributions*

- Data  $\mathbf{x} = \{x_1, \dots, x_n\}$ ,  $x_i \in \mathcal{X}$ , set of “homogeneous” observations. Desired to predict the value of a future observation  $x \in \mathcal{X}$  generated by the same mechanism.
- From the foundations arguments the solution *must* be a probability distribution  $p(x | \mathbf{x}, K)$  describing the uncertainty on the value that  $x$  will take, given data  $\mathbf{x}$  and any other available knowledge  $K$ . This is called the (posterior) *predictive density* of  $x$ .
- To derive  $p(x | \mathbf{x}, K)$  it is necessary to specify the *precise sense* in which the  $x_i$ 's are judged to be *homogeneous*.
- It is often directly assumed that the data  $\mathbf{x} = \{x_1, \dots, x_n\}$  consist of a *random sample* from some specified model,  $\{p(x | \boldsymbol{\theta}), x \in \mathcal{X}, \boldsymbol{\theta} \in \Theta\}$ , so that  $p(\mathbf{x} | \boldsymbol{\theta}) = p(x_1, \dots, x_n | \boldsymbol{\theta}) = \prod_{j=1}^n p(x_j | \boldsymbol{\theta})$ .

If this is the case, the solution to the prediction problem is immediate once a prior distribution  $\pi(\boldsymbol{\theta})$  has been specified.

- *Posterior predictive distributions from random samples*

□ Let  $\mathbf{x} = \{x_1, \dots, x_n\}$ ,  $x_i \in \mathcal{X}$  a random sample of size  $n$  from the statistical model  $\{p(x | \boldsymbol{\theta}), x \in \mathcal{X}, \boldsymbol{\theta} \in \Theta\}$

Let  $\pi(\boldsymbol{\theta})$  a prior distribution describing available knowledge (in any) about the value of the parameter vector  $\boldsymbol{\theta}$ .

The *posterior predictive distribution* is

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

This encapsulates all available information about the outcome of any future observation  $x \in \mathcal{X}$  from the same model.

□ To prove this, make use the total probability theorem, to have

$$p(x | \mathbf{x}) = \int_{\Theta} p(x | \boldsymbol{\theta}, \mathbf{x}) \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta}$$

and notice the new observation  $x$  has been assumed to be conditionally independent of the observed data  $\mathbf{x}$ , so that  $p(x | \boldsymbol{\theta}, \mathbf{x}) = p(x | \boldsymbol{\theta})$ .

□ The observable values  $x \in \mathcal{X}$  may be either *discrete* or *continuous* random quantities. In the discrete case, the predictive distribution will be described by its probability *mass* function; in the continuous case, by its probability *density* function. Both are denoted  $p(x | \mathbf{x})$ .

- *Prediction in a Poisson process*

- Data  $\mathbf{x} = \{r_1, \dots, r_n\}$  random from  $\text{Pn}(r | \lambda)$ . The reference posterior density of  $\lambda$  is  $\pi^*(\lambda | \mathbf{x}) = \text{Ga}(\lambda | t + 1/2, n)$ , where  $t = \sum_j r_j$ .

The (reference) posterior predictive distribution is

$$\begin{aligned} p(r | \mathbf{x}) &= \Pr[r | t, n] = \int_0^\infty \text{Pn}(r | \lambda) \text{Ga}(\lambda | t + \frac{1}{2}, n) d\lambda \\ &= \frac{n^{t+1/2}}{\Gamma(t + 1/2)} \frac{1}{r!} \frac{\Gamma(r + t + 1/2)}{(1 + n)^{r+t+1/2}}, \end{aligned}$$

an example of a Poisson-Gamma probability mass function.

- For example, no flash floods have been recorded on a particular location in 10 consecutive years. Local authorities are interested in forecasting possible future flash floods. Using a Poisson model, and assuming that meteorological conditions remain similar, the probabilities that  $r$  flash floods will occur next year in that location are given by the Poisson-Gamma mass function above, with  $t = 0$  and  $n = 10$ . This yields,  $\Pr[0 | t, n] = 0.953$ ,  $\Pr[1 | t, n] = 0.043$ , and  $\Pr[2 | t, n] = 0.003$ .

Many other situations may be described with the same model.

- *Prediction of Normal measurements*

□ Data  $\mathbf{x} = \{x_1, \dots, x_n\}$  random from  $N(x | \mu, \sigma)$ . Reference prior  $\pi^*(\mu, \sigma) = \sigma^{-1}$  or, in terms of the precision  $\lambda = \sigma^{-2}$ ,  $\pi^*(\mu, \lambda) = \lambda^{-1}$ .

The *joint* reference posterior,  $\pi^*(\mu, \lambda | \mathbf{x}) \propto p(\mathbf{x} | \mu, \lambda) \pi^*(\mu, \lambda)$ , is  $\pi^*(\mu, \lambda | \mathbf{x}) = N(\mu | \bar{x}, (n\lambda)^{-1/2}) \text{Ga}(\lambda | (n-1)/2, ns^2/2)$ .

□ The predictive distribution is

$$\begin{aligned} \pi^*(x | \mathbf{x}) &= \int_0^\infty \int_{-\infty}^\infty \mathbf{N}(x | \mu, \lambda^{-1/2}) \pi^*(\mu, \lambda | \mathbf{x}) \, d\mu \, d\lambda \\ &\propto \{(1+n)s^2 + (\mu - \bar{x})^2\}^{-n/2}, \end{aligned}$$

a kernel of the *Student* density  $\pi^*(x | \mathbf{x}) = \text{St}(x | \bar{x}, s \sqrt{\frac{n+1}{n-1}}, n-1)$ .

□ *Example.* Production of safety belts. Observed breaking strengths of 10 randomly chosen webbings have mean  $\bar{x} = 28.011$  kN and standard deviation  $s = 0.443$  kN. Specification requires  $x > 26$  kN.

Reference posterior predictive  $p(x | \mathbf{x}) = \text{St}(x | 28.011, 0.490, 9)$ .

$$\Pr(x > 26 | \mathbf{x}) = \int_{26}^\infty \text{St}(x | 28.011, 0.490, 9) \, dx = 0.9987.$$

- *Regression*

- Often *additional information* from relevant covariates. Data structure, set of pairs  $\mathbf{x} = \{(\mathbf{y}_1, \mathbf{v}_1), \dots, (\mathbf{y}_n, \mathbf{v}_n)\}$ ;  $\mathbf{y}_i, \mathbf{v}_i$ , both vectors. Given a new observation, with  $\mathbf{v}$  known, predict the corresponding value of  $\mathbf{y}$ . Formally, compute  $p\{\mathbf{y} | \mathbf{v}, (\mathbf{y}_1, \mathbf{v}_1), \dots, (\mathbf{y}_n, \mathbf{v}_n)\}$ .
- Need a model  $\{p(\mathbf{y} | \mathbf{v}, \boldsymbol{\theta}), \mathbf{y} \in \mathbf{Y}, \boldsymbol{\theta} \in \Theta\}$  which makes precise the probabilistic relationship between  $\mathbf{y}$  and  $\mathbf{v}$ . The simplest option assumes a *linear dependency* of the form  $p(\mathbf{y} | \mathbf{v}, \boldsymbol{\theta}) = \mathbf{N}(\mathbf{y} | \mathbf{V}\boldsymbol{\beta}, \Sigma)$ , but far more complex structures are common in applications.
- *Univariate linear regression on  $k$  covariates*.  $Y \subset \mathfrak{R}, \mathbf{v} = \{v_1, \dots, v_k\}$ .  $p(y | \mathbf{v}, \boldsymbol{\beta}, \sigma) = \mathbf{N}(y | \mathbf{v}\boldsymbol{\beta}, \sigma^2)$ ,  $\boldsymbol{\beta} = \{\beta_1, \dots, \beta_k\}^t$ . Data  $\mathbf{x} = \{\mathbf{y}, \mathbf{V}\}$ ,  $\mathbf{y} = \{y_1, \dots, y_n\}^t$ , and  $\mathbf{V}$  is the  $n \times k$  matrix with the  $\mathbf{v}_i$ 's as rows.  $p(\mathbf{y} | \mathbf{V}, \boldsymbol{\beta}, \sigma) = \mathbf{N}_n(\mathbf{y} | \mathbf{V}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$ ; reference prior  $\pi^*(\boldsymbol{\beta}, \sigma) = \sigma^{-1}$ .

Predictive posterior is the Student density

$$p(y | \mathbf{v}, \mathbf{y}, \mathbf{V}) = \text{St}(y | \mathbf{v}\hat{\boldsymbol{\beta}}, s \sqrt{f(\mathbf{v}, \mathbf{V}) \frac{n}{n-k}}, n - k)$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{V}^t \mathbf{V})^{-1} \mathbf{V}^t \mathbf{y}, \quad ns^2 = (\mathbf{y} - \mathbf{v}\hat{\boldsymbol{\beta}})^t (\mathbf{y} - \mathbf{v}\hat{\boldsymbol{\beta}})$$

$$f(\mathbf{v}, \mathbf{V}) = 1 + \mathbf{v}(\mathbf{V}^t \mathbf{V})^{-1} \mathbf{v}^t$$



- *Example: Simple linear regression*

- One covariate and a constant term;  $p(y | v, \beta, \sigma) = \mathbf{N}(y | \beta_1 + \beta_2 v, \sigma)$   
Sufficient statistic is  $\mathbf{t} = \{\bar{v}, \bar{y}, s_{vy}, s_{vv}\}$ , with  $n\bar{v} = \sum v_j$ ,  $n\bar{y} = \sum y_j$ ,  
 $s_{yv} = \sum v_j y_j / n - \bar{v} \bar{y}$ ,  $s_{vv} = \sum v_j^2 / n - \bar{v}^2$ .

$$p(y | v, \mathbf{t}) = \text{St}(y | \hat{\beta}_1 + \hat{\beta}_2 v, s \sqrt{f(v, \mathbf{t}) \frac{n}{n-2}}, n-2)$$

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{v}, \quad \hat{\beta}_2 = \frac{s_{vy}}{s_{vv}},$$

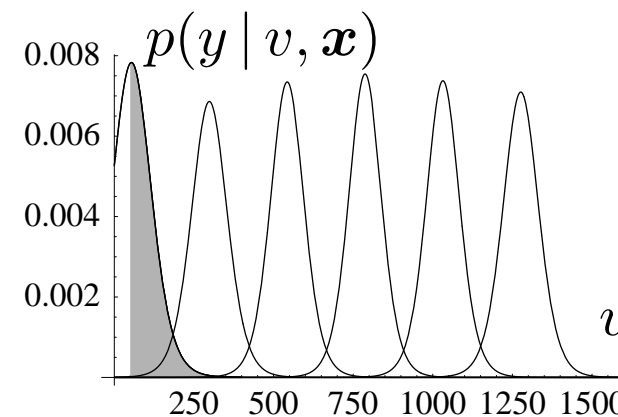
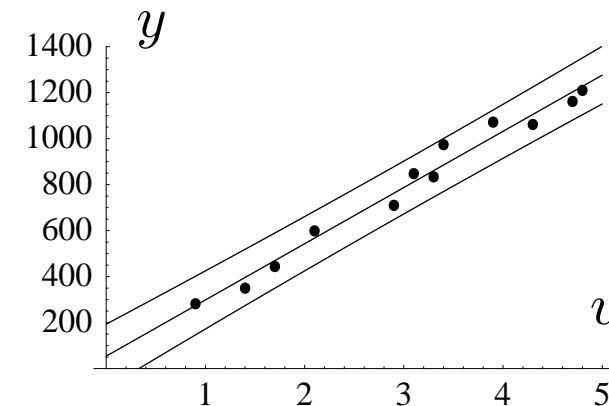
$$ns^2 = \sum_{j=1}^n (y_j - \hat{\beta}_1 - \hat{\beta}_2 x_j)^2$$

$$f(v, \mathbf{t}) = 1 + \frac{1}{n} \frac{(v - \bar{v})^2 + s_{vv}}{s_{vv}}$$

- Pollution density ( $\mu\text{gr}/\text{m}^3$ ), and wind speed from source ( $\text{m}/\text{s}$ ).

$y_j$	1212	836	850	446	1164	601
$v_j$	4.8	3.3	3.1	1.7	4.7	2.1
$y_j$	1074	284	352	1064	712	976
$v_j$	3.9	0.9	1.4	4.3	2.9	3.4

$$\Pr[y > 50 | v = 0, \mathbf{x}] = 0.66$$



## 2.4. Hierarchical Models

- *Exchangeability*

- Random quantities are often “homogeneous” in the precise sense that only their *values* matter, not the *order* in which they appear. Formally, this is captured by the notion of *exchangeability*. The set of random vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is exchangeable if their joint distribution is invariant under permutations. An infinite sequence  $\{\mathbf{x}_j\}$  of random vectors is exchangeable if all its finite subsequences are exchangeable.
- *Any random sample from any model is exchangeable.* The *representation theorem* establishes that if observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are exchangeable, they are a *random sample* from some model  $\{p(\mathbf{x} | \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$ , labeled by a *parameter vector*  $\boldsymbol{\theta}$ , defined as the limit (as  $n \rightarrow \infty$ ) of some function of the  $\mathbf{x}_i$ 's. Information about  $\boldsymbol{\theta}$  in prevailing conditions  $C$  is *necessarily* described by *some* probability distribution  $\pi(\boldsymbol{\theta} | C)$ .
- Formally, the joint density of any finite set of exchangeable observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  has an *integral representation* of the form
 
$$p(\mathbf{x}_1, \dots, \mathbf{x}_n | C) = \int_{\Theta} \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | C) d\boldsymbol{\theta}.$$

- *Structured Models*

- Complex data structures may often be usefully described by partial exchangeability assumptions.
- *Example: Public opinion.* Sample  $k$  different regions in the country. Sample  $n_i$  citizens in region  $i$  and record whether or not ( $y_{ij} = 1$  or  $y_{ij} = 0$ ) citizen  $j$  would vote  $A$ . Assuming exchangeable citizens within each region implies

$$p(y_{i1}, \dots, y_{in_i}) = \prod_{j=1}^{n_i} p(y_{ij} | \theta_i) = \theta_i^{r_i} (1 - \theta_i)^{n_i - r_i},$$

where  $\theta_i$  is the (unknown) proportion of citizens in region  $i$  voting  $A$  and  $r_i = \sum_j y_{ij}$  the number of citizens voting  $A$  in region  $i$ .

Assuming regions exchangeable within the country similarly leads to

$$p(\theta_1, \dots, \theta_k) = \prod_{i=1}^k \pi(\theta_i | \phi)$$

for some probability distribution  $\pi(\theta | \phi)$  describing the political variation within the regions. Often choose  $\pi(\theta | \phi) = \text{Be}(\theta | \alpha, \beta)$ .

- The resulting *two-stages hierarchical Binomial-Beta model*  $\mathbf{x} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ ,  $\mathbf{y}_i = \{y_{i1}, \dots, y_{in_i}\}$ , random from  $\text{Bi}(y | \theta_i)$ ,  $\{\theta_1, \dots, \theta_k\}$ , random from  $\text{Be}(\theta | \alpha, \beta)$  provides a far richer model than (unrealistic) simple binomial sampling.

- *Example: Biological response.* Sample  $k$  different animals of the same species in specific environment. Control  $n_i$  times animal  $i$  and record his responses  $\{\mathbf{y}_{i1}, \dots, \mathbf{y}_{in_i}\}$  to prevailing conditions. Assuming exchangeable observations within each animal implies

$$p(\mathbf{y}_{i1}, \dots, \mathbf{y}_{in_i}) = \prod_{j=1}^{n_i} p(\mathbf{y}_{ij} | \boldsymbol{\theta}_i).$$

Often choose  $p(\mathbf{y}_{ij} | \boldsymbol{\theta}_i) = \mathbf{N}_r(\mathbf{y} | \boldsymbol{\mu}_i, \Sigma_1)$ , where  $r$  is the number of biological responses measured.

Assuming exchangeable animals within the environment leads to

$$p(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \prod_{i=1}^k \pi(\boldsymbol{\mu}_i | \phi)$$

for some probability distribution  $\pi(\boldsymbol{\mu} | \phi)$  describing the biological variation within the species. Often choose  $\pi(\boldsymbol{\mu} | \phi) = \mathbf{N}_r(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \Sigma_2)$ .

- The *two-stages hierarchical multivariate Normal-Normal model*  
 $\mathbf{x} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ ,  $\mathbf{y}_i = \{\mathbf{y}_{i1}, \dots, \mathbf{y}_{in_i}\}$ , random from  $\mathbf{N}_r(\mathbf{y} | \boldsymbol{\mu}_i, \Sigma_1)$ ,  
 $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}$ , random from  $\mathbf{N}_r(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \Sigma_2)$   
 provides a far richer model than (unrealistic) simple multivariate normal sampling.
- Finer subdivisions, *e.g.*, subspecies within each species, similarly lead to hierarchical models with more stages.

- *Bayesian analysis of hierarchical models*

- A *two-stages hierarchical model* has the general form

$$\mathbf{x} = \{\mathbf{y}_1, \dots, \mathbf{y}_k\}, \mathbf{y}_i = \{z_{i1}, \dots, z_{in_i}\}$$

$\mathbf{y}_i$  random sample of size  $n_i$  from  $p(\mathbf{z} | \boldsymbol{\theta}_i)$ ,  $\boldsymbol{\theta}_i \in \Theta$ ,  
 $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k\}$ , random of size  $k$  from  $\pi(\boldsymbol{\theta} | \phi)$ ,  $\phi \in \Phi$ .

- Specify a *prior distribution* (or a reference prior function)  $\pi(\phi)$  for the *hyperparameter vector*  $\phi$ .

- Use *standard probability theory* to compute all desired *posterior distributions*:

$\pi(\phi | \mathbf{x})$  for inferences about the hyperparameters,

$\pi(\boldsymbol{\theta}_i | \mathbf{x})$  for inferences about the parameters,

$\pi(\psi | \mathbf{x})$  for inferences about the any function  $\psi = \psi(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)$   
of the parameters,

$\pi(\mathbf{y} | \mathbf{x})$  for predictions on future observations,

$\pi(t | \mathbf{x})$  for predictions on any function  $t = t(\mathbf{y}_1, \dots, \mathbf{y}_m)$   
of  $m$  future observations

- *Markov Chain Monte Carlo* based *software* available for the necessary computations.

## 3. Decision Making

### 3.1 Structure of a Decision Problem

- *Alternatives, consequences, relevant events*
  - A decision problem if two or more possible courses of action;  $\mathcal{A}$  is the class of possible *actions*.
  - For each  $a \in \mathcal{A}$ ,  $\Theta_a$  is the set of *relevant events*, those may affect the result of choosing  $a$ .
  - Each pair  $\{a, \theta\}$ ,  $\theta \in \Theta_a$ , produces a consequence  $c(a, \theta) \in \mathcal{C}_a$ . In this context,  $\theta$  is often referred to as the *parameter of interest*.
  - 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 the appropriate Cartesian product may be used.
  - In many problems the class of relevant events  $\Theta_a$  is the same for all  $a \in \mathcal{A}$ . Even if this is not the case, a comprehensive *parameter space*  $\Theta$  may be defined as the union of all the  $\Theta_a$ .

- *Foundations of decision theory*

- Different sets of principles capture a minimum collection of logical rules required for “rational” decision-making.

These are axioms with strong intuitive appeal.

Their basic structure consists of:

- The *Transitivity* of preferences:

If  $a_1 \succ a_2$  given  $C$ , and  $a_2 \succ a_3$  given  $C$ ,  
then  $a_1 \succ a_3$  given  $C$ .

- The *Sure-thing principle*:

If  $a_1 \succ a_2$  given  $C$  and  $E$ , and  $a_1 \succ a_2$  given  $C$  and not  $E$   
then  $a_1 \succ a_2$  given  $C$ .

- The existence of *Standard events*:

There are events of known plausibility.

These may be used as a unit of measurement, and  
have the properties of a probability measure

- These axioms are not a description of human decision-making,  
but a *normative* set of principles defining *coherent* decision-making.

- *Decision making*

- Many different axiom sets.

All lead basically to the same set of conclusions, namely:

- The consequences of wrong actions should be evaluated in terms of a real-valued *loss* function  $\ell(a, \boldsymbol{\theta})$  which specifies, on a numerical scale, their undesirability.
- The uncertainty about the parameter of interest  $\boldsymbol{\theta}$  should be measured with a *probability distribution*  $\pi(\boldsymbol{\theta} | C)$

$$\pi(\boldsymbol{\theta} | C) \geq 0, \quad \boldsymbol{\theta} \in \Theta, \quad \int_{\Theta} \pi(\boldsymbol{\theta} | C) d\boldsymbol{\theta} = 1,$$

describing all available knowledge about its value, given the conditions  $C$  under which the decision must be taken.

- The relative undesirability of available actions  $a \in \mathcal{A}$  is measured by their expected loss: *the optimal action minimizes the expected loss.*

$$\ell[a | C] = \int_{\Theta} \ell(a, \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | C) d\boldsymbol{\theta}, \quad a \in \mathcal{A}.$$

(alternatively, one may *maximize expected utility*)



- *Intrinsic loss functions: Intrinsic discrepancy*

- The loss function is typically *context dependent*.
- In mathematical statistics, *intrinsic* loss functions are used to measure the distance between between statistical models.

They measure the *divergence between the models*  $\{p_1(\mathbf{x} | \boldsymbol{\theta}_1), \mathbf{x} \in \mathcal{X}\}$  and  $\{p_2(\mathbf{x} | \boldsymbol{\theta}_2), \mathbf{x} \in \mathcal{X}\}$  as some *non-negative* function of the form  $\ell\{p_1, p_2\}$  which is zero if (and only if) the two distributions are equal almost everywhere.

- The *intrinsic discrepancy* between two statistical models is simply the intrinsic discrepancy between their sampling distributions, *i.e.*,

$$\begin{aligned} \delta\{p_1, p_2\} &= \delta\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\} \\ &= \min \left\{ \int_{\mathcal{X}_1} p_1(\mathbf{x} | \boldsymbol{\theta}_1) \log \frac{p_1(\mathbf{x} | \boldsymbol{\theta}_1)}{p_2(\mathbf{x} | \boldsymbol{\theta}_2)} d\mathbf{x}, \int_{\mathcal{X}_2} p_2(\mathbf{x} | \boldsymbol{\theta}_2) \log \frac{p_2(\mathbf{x} | \boldsymbol{\theta}_2)}{p_1(\mathbf{x} | \boldsymbol{\theta}_1)} d\mathbf{x} \right\} \end{aligned}$$

- The intrinsic discrepancy is an *information-based, symmetric, invariant intrinsic loss*.

## 3.2 Point and Region Estimation

- *Point estimation as a decision problem*

- Given statistical model  $\{p(\mathbf{x} | \omega), \mathbf{x} \in \mathcal{X}, \omega \in \Omega\}$ , quantity of interest  $\theta = \theta(\omega) \in \Theta$ . A *point estimator*  $\tilde{\theta} = \tilde{\theta}(\mathbf{x})$  of  $\theta$  is some function of the data to be regarded as a proxy for the unknown value of  $\theta$ .
- To choose a point estimate for  $\theta$  is a *decision problem*, where the action space is  $\mathcal{A} = \Theta$ .
- Given a *loss function*  $\ell(\tilde{\theta}, \theta)$ , the posterior expected loss is

$$\ell[\tilde{\theta} | \mathbf{x}] = \int_{\Theta} \ell(\tilde{\theta}, \theta) \pi(\theta | \mathbf{x}) d\theta,$$

The corresponding *Bayes estimator* is the function of the data,

$$\theta^* = \theta^*(\mathbf{x}) = \arg \inf_{\tilde{\theta} \in \Theta} \ell[\tilde{\theta} | \mathbf{x}],$$

which minimizes that expectation.

- *Conventional estimators*

- The *posterior mean* and the *posterior mode* are the Bayes estimators which respectively correspond to a *quadratic* and a *zero-one* loss functions.
  - If  $\ell(\tilde{\theta}, \theta) = (\tilde{\theta} - \theta)^t (\tilde{\theta} - \theta)$ , then, assuming that the mean exists, the Bayes estimator is the *posterior mean*  $E[\theta | \mathbf{x}]$ .
- • If the loss function is a zero-one function, so that  $\ell(\tilde{\theta}, \theta) = 0$  if  $\tilde{\theta}$  belongs to a ball of radius  $\varepsilon$  centered in  $\theta$  and  $\ell(\tilde{\theta}, \theta) = 1$  otherwise then, assuming that a unique mode exists, the Bayes estimator converges to the *posterior mode*  $\text{Mo}[\theta | \mathbf{x}]$  as the ball radius  $\varepsilon$  tends to zero.
- If  $\theta$  is *univariate and continuous*, and the loss function is *linear*,

$$\ell(\tilde{\theta}, \theta) = \begin{cases} c_1(\tilde{\theta} - \theta) & \text{if } \tilde{\theta} \geq \theta \\ c_2(\theta - \tilde{\theta}) & \text{if } \tilde{\theta} < \theta \end{cases}$$

then the Bayes estimator is the *posterior quantile* of order  $c_2/(c_1 + c_2)$ , so that  $\text{Pr}[\theta < \theta^*] = c_2/(c_1 + c_2)$ .

In particular, if  $c_1 = c_2$ , the Bayes estimator is the *posterior median*.

- Any  $\theta$  value may be optimal: *it all depends on the loss function*.

- *Intrinsic point estimation*

- Given the statistical model  $\{p(\mathbf{x} | \boldsymbol{\theta}), \mathbf{x} \in \mathcal{X}, \boldsymbol{\theta} \in \Theta\}$  the intrinsic discrepancy  $\delta(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  between two parameter values  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\theta}_2$  is the intrinsic discrepancy  $\delta\{p(\mathbf{x} | \boldsymbol{\theta}_1), p(\mathbf{x} | \boldsymbol{\theta}_2)\}$  between the corresponding probability models.

This is symmetric, non-negative (and zero iff  $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$ ), invariant under reparametrization and invariant under bijections of  $\mathbf{x}$ .

- The intrinsic estimator is the *reference* Bayes estimator which corresponds to the loss defined by the *intrinsic discrepancy*:

- The expected loss with respect to the reference posterior distribution

$$d(\tilde{\boldsymbol{\theta}} | \mathbf{x}) = \int_{\Theta} \delta\{\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}\} \pi^*(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta}$$

is an objective measure, in information units, of the *expected* discrepancy between the model  $p(\mathbf{x} | \tilde{\boldsymbol{\theta}})$  and the true (unknown) model  $p(\mathbf{x} | \boldsymbol{\theta})$ .

- The *intrinsic estimator*  $\boldsymbol{\theta}^* = \boldsymbol{\theta}^*(\mathbf{x})$  is the value which minimizes such expected discrepancy,

$$\boldsymbol{\theta}^* = \arg \inf_{\tilde{\boldsymbol{\theta}} \in \Theta} d(\tilde{\boldsymbol{\theta}} | \mathbf{x}).$$

- *Example: Intrinsic estimation of the Binomial parameter*

- Data  $\mathbf{x} = \{x_1, \dots, x_n\}$ , random from  $p(x | \theta) = \theta^x (1 - \theta)^{1-x}$ ,  $r = \sum x_j$ . Intrinsic discrepancy  $\delta(\tilde{\theta}, \theta) = n \min\{k(\tilde{\theta} | \theta), k(\theta | \tilde{\theta})\}$ ,  
 $k(\theta_1 | \theta_2) = \theta_2 \log \frac{\theta_2}{\theta_1} + (1 - \theta_2) \log \frac{1 - \theta_2}{1 - \theta_1}$ ,  $\pi^*(\theta) = \text{Be}(\theta | \frac{1}{2}, \frac{1}{2})$   
 $\pi^*(\theta | r, n) = \text{Be}(\theta | r + \frac{1}{2}, n - r + \frac{1}{2})$ .

- Expected reference discrepancy  
 $d(\tilde{\theta}, r, n) = \int_0^1 \delta(\tilde{\theta}, \theta) \pi^*(\theta | r, n) d\theta$

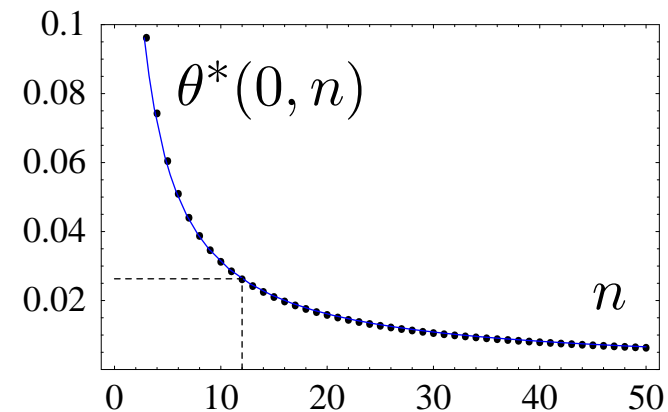
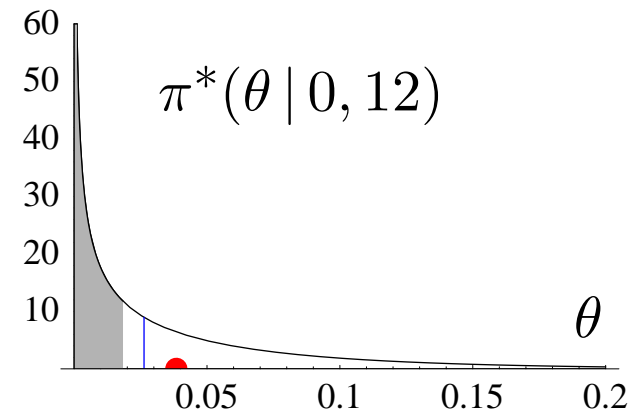
- Intrinsic estimator  
 $\theta^*(r, n) = \arg \min_{0 < \tilde{\theta} < 1} d(\tilde{\theta}, r, n)$

From invariance, for any bijection  
 $\phi = \phi(\theta)$ ,  $\phi^* = \phi(\theta^*)$ .

- Analytic approximation

$$\theta^*(r, n) \approx \frac{r + 1/3}{n + 2/3}, \quad n > 2$$

- $n = 12, r = 0, \theta^*(0, 12) = 0.026$   
 $\text{Me}[\theta | \mathbf{x}] = 0.018, \mathbf{E}[\theta | \mathbf{x}] = 0.038$



- *Intrinsic region (interval) estimation*

□ The *intrinsic  $q$ -credible region*  $R^*(q) \subset \Theta$  is that  $q$ -credible reference region which corresponds to minimum expected intrinsic loss:

(i)  $\int_{R^*(q)} \pi^*(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} = q$

(ii)  $\forall \boldsymbol{\theta}_i \in R^*(q), \forall \boldsymbol{\theta}_j \notin R^*(q), \quad d(\boldsymbol{\theta}_i | \mathbf{x}) < d(\boldsymbol{\theta}_j | \mathbf{x})$

□ Binomial examples:  $d(\boldsymbol{\theta}_i | \mathbf{x}) = d(\theta_i | r, n)$

$r = 0, n = 12,$

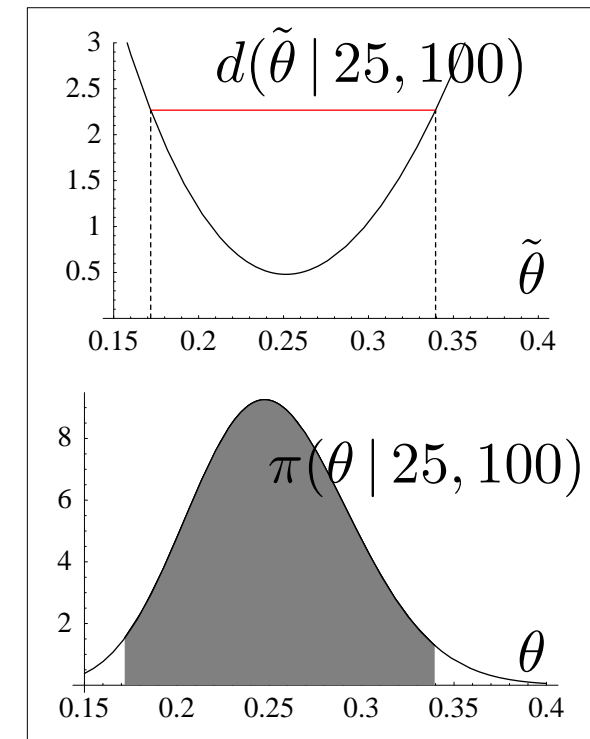
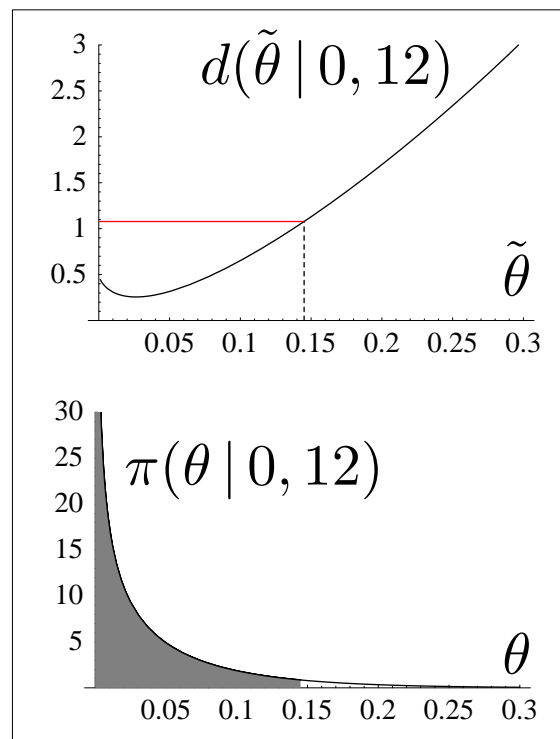
$\theta^* = 0.0263;$

$R_{0.95}^* = [0, 0.145];$

$r = 25, n = 100,$

$\theta^* = 0.2514;$

$R_{0.95}^* = [0.172, 0.340];$



## 3.3 Hypothesis Testing

- *Precise hypothesis testing as a decision problem*
  - The posterior  $\pi(\boldsymbol{\theta} | D)$  conveys intuitive information on the values of  $\boldsymbol{\theta}$  which are *compatible* with the observed data  $\boldsymbol{x}$ : those with a *relatively high probability density*.
  - Often a particular value  $\boldsymbol{\theta}_0$  is suggested for special consideration:
    - Because  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$  would greatly simplify the model
    - Because there are context specific arguments suggesting that  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$More generally, one may analyze the *restriction* of parameter space  $\Theta$  to a subset  $\Theta_0$  which may contain more than one value.
  - Formally, testing the hypothesis  $H_0 \equiv \{\boldsymbol{\theta} = \boldsymbol{\theta}_0\}$  is a *decision problem* with just two possible actions:
    - $a_0$ : to *accept*  $H_0$  and work with  $p(\boldsymbol{x} | \boldsymbol{\theta}_0)$ .
    - $a_1$ : to *reject*  $H_0$  and keep the general model  $p(\boldsymbol{x} | \boldsymbol{\theta})$ .
  - To proceed, a *loss* function  $\ell(a_i, \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} \in \Theta$ , describing the possible consequences of both actions, must be specified.

- *Structure of the loss function*

- Given data  $\mathbf{x}$ , optimal action is to reject  $H_0$  (action  $a_1$ ) *iff* the expected posterior loss of accepting,  $\int_{\Theta} \ell(a_0, \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta}$ , is *larger* than the expected posterior loss of rejecting,  $\int_{\Theta} \ell(a_1, \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta}$ , *i.e.*, iff
 
$$\int_{\Theta} [\ell(a_0, \boldsymbol{\theta}) - \ell(a_1, \boldsymbol{\theta})] \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} = \int_{\Theta} \Delta\ell(\boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} > 0.$$

Therefore, only the loss difference  $\Delta\ell(\boldsymbol{\theta}) = \ell(a_0, \boldsymbol{\theta}) - \ell(a_1, \boldsymbol{\theta})$ , which measures the *advantage* of rejecting  $H_0$  as a function of  $\boldsymbol{\theta}$ , has to be specified: The hypothesis should be rejected whenever the *expected* advantage of rejecting is positive.

- The advantage  $\Delta\ell(\boldsymbol{\theta})$  of rejecting  $H_0$  as a function of  $\boldsymbol{\theta}$  should be of the form  $\Delta\ell(\boldsymbol{\theta}) = l(\boldsymbol{\theta}_0, \boldsymbol{\theta}) - l^*$ , for some  $l^* > 0$ , where
  - $l(\boldsymbol{\theta}_0, \boldsymbol{\theta})$  measures the *discrepancy* between  $p(\mathbf{x} | \boldsymbol{\theta}_0)$  and  $p(\mathbf{x} | \boldsymbol{\theta})$ ,
  - $l^*$  is a positive *utility constant* which measures the advantage working with the simpler model when it is true.
- The Bayes criterion will then be: *Reject*  $H_0$  if (and only if)
 
$$\int_{\Theta} l(\boldsymbol{\theta}_0, \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} > l^*,$$
 that is if (and only if) the *expected discrepancy* between  $p(\mathbf{x} | \boldsymbol{\theta}_0)$  and  $p(\mathbf{x} | \boldsymbol{\theta})$  is *too large*.



- *Bayesian Reference Criterion*

□ An good choice for the function  $l(\boldsymbol{\theta}_0, \boldsymbol{\theta})$  is the *intrinsic discrepancy*,  
 $\delta(\boldsymbol{\theta}_0, \boldsymbol{\theta}) = \min \{k(\boldsymbol{\theta}_0 | \boldsymbol{\theta}), k(\boldsymbol{\theta} | \boldsymbol{\theta}_0)\},$

where  $k(\boldsymbol{\theta}_0 | \boldsymbol{\theta}) = \int_{\mathcal{X}} p(\boldsymbol{x} | \boldsymbol{\theta}) \log\{p(\boldsymbol{x} | \boldsymbol{\theta})/p(\boldsymbol{x} | \boldsymbol{\theta}_0)\} d\boldsymbol{x}.$

If  $\boldsymbol{x} = \{x_1, \dots, x_n\} \in \mathcal{X}^n$  is a random sample from  $p(\boldsymbol{x} | \boldsymbol{\theta})$ , then

$$k(\boldsymbol{\theta}_0 | \boldsymbol{\theta}) = n \int_{\mathcal{X}} p(\boldsymbol{x} | \boldsymbol{\theta}) \log \frac{p(\boldsymbol{x} | \boldsymbol{\theta})}{p(\boldsymbol{x} | \boldsymbol{\theta}_0)} d\boldsymbol{x}.$$

□ For objective results, exclusively based on model assumptions and data, the *reference* posterior distribution  $\pi^*(\boldsymbol{\theta} | \boldsymbol{x})$  should be used.

□ Hence, *reject if (and only if) the expected reference posterior intrinsic discrepancy  $d(\boldsymbol{\theta}_0 | \boldsymbol{x})$  is too large,*

$$d(\boldsymbol{\theta}_0 | \boldsymbol{x}) = \int_{\Theta} \delta(\boldsymbol{\theta}_0, \boldsymbol{\theta}) \pi^*(\boldsymbol{\theta} | \boldsymbol{x}) d\boldsymbol{\theta} > d^*, \text{ for some } d^* > 0.$$

This is the *Bayesian reference criterion (BRC)*.

□ The *reference test statistic*  $d(\boldsymbol{\theta}_0 | \boldsymbol{x})$  is nonnegative, it is invariant both under reparametrization and under sufficient transformation of the data, and it is a measure, in natural information units (nits) of the expected discrepancy between  $p(\boldsymbol{x} | \boldsymbol{\theta}_0)$  and the true model.

- *Calibration of the BRC*

- The reference test statistic  $d(\theta_0 | \mathbf{x})$  is the posterior expected value of the intrinsic discrepancy between  $p(\mathbf{x} | \theta_0)$  and  $p(\mathbf{x} | \theta)$ .
  - A reference test statistic value  $d(\theta_0 | \mathbf{x}) \approx 1$  suggests that data are clearly compatible with the Hypothesis that  $\theta = \theta_0$ .
  - A test statistic value  $d(\theta_0 | \mathbf{x}) \log(10) = 2.303$  nits implies that, given data  $\mathbf{x}$ , the *average* value of the likelihood ratio *against* the hypothesis,  $p(\mathbf{x} | \theta) / p(\mathbf{x} | \theta_0)$ , is expected to be about 10: *mild evidence* against  $\theta_0$ .
  - Similarly,  $d(\theta_0 | \mathbf{x}) \approx \log(100) = 4.605$  (expected likelihood ratio against  $\theta_0$  about 100), indicates *strong evidence* against  $\theta_0$ , and  $\log(1000) = 6.908$ , *conclusive evidence* against  $\theta_0$ .
- Strong connections between BRC and intrinsic estimation:
  - The *intrinsic estimator* is the value of  $\theta$  with minimizes the reference test statistic:  $\theta^* = \arg \inf_{\theta \in \Theta} d(\theta | \mathbf{x})$ .
  - The regions defined by  $\{\theta; d(\theta | \mathbf{x}) \leq d^*\}$  are invariant *reference posterior  $q(d^*)$ -credible regions* for  $\theta$ . For regular problems and large samples,  $q(\log(10)) \approx 0.95$  and  $q(\log(100)) \approx 0.995$ .

- *A canonical example: Testing a value for the Normal mean*

□ In the simplest case where the variance  $\sigma^2$  is known,

$$\delta(\mu_0, \mu) = n(\mu - \mu_0)^2 / (2\sigma^2), \quad \pi^*(\mu | \mathbf{x}) = \mathbf{N}(\mu | \bar{x}, \sigma / \sqrt{n}),$$

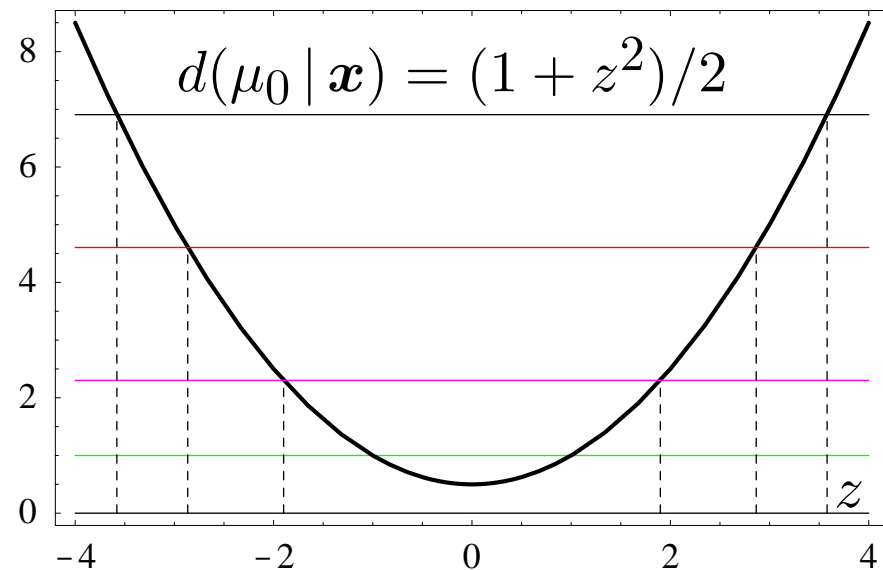
$$d(\mu_0 | \mathbf{x}) = \frac{1}{2}(1 + z^2), \quad z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

Thus rejecting  $\mu = \mu_0$  if  $d(\mu_0 | \mathbf{x}) > d^*$  is equivalent to rejecting if  $|z| > \sqrt{2d^* - 1}$  and, hence, to a conventional two-sided frequentist test with significance level  $\alpha = 2(1 - \Phi(|z|))$ .

$d^*$	$ z $	$\alpha$
$\log(10)$	1.8987	0.0576
$\log(100)$	2.8654	0.0042
$\log(1000)$	3.5799	0.0003

□ The expected value of  $d(\mu_0 | \mathbf{x})$  if the hypothesis is **true** is

$$\int_{-\infty}^{\infty} \frac{1}{2}(1 + z^2) \mathbf{N}(z | 0, 1) dz = 1$$



- Fisher's tasting tea lady*

- Data  $\mathbf{x} = \{x_1, \dots, x_n\}$ , random from  $p(x | \theta) = \theta^x (1 - \theta)^{1-x}$ ,  
 $r = \sum x_j$ . Intrinsic discrepancy  $\delta(\theta_0, \theta) = n \min\{k(\theta_0 | \theta), k(\theta | \theta_0)\}$ ,  
 $k(\theta_1 | \theta_2) = \theta_2 \log \frac{\theta_2}{\theta_1} + (1 - \theta_2) \log \frac{1-\theta_2}{1-\theta_1}$ ,  $\pi^*(\theta | r, n) = \text{Be}(\theta | r + \frac{1}{2}, n - r + \frac{1}{2})$

Intrinsic test statistic

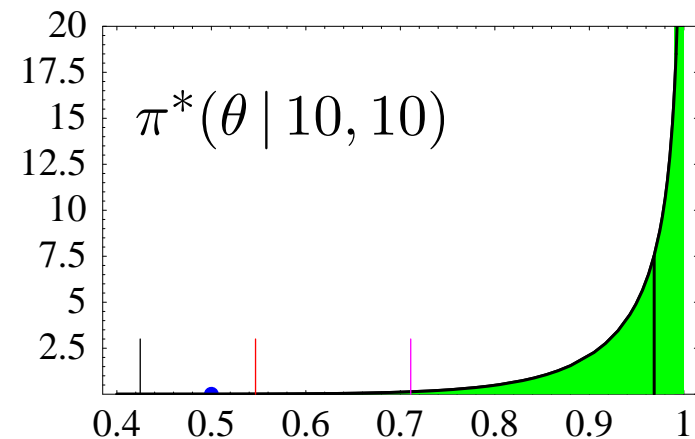
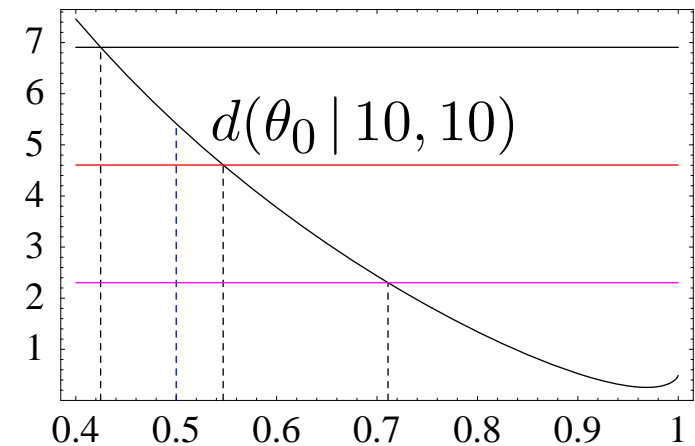
$$d(\theta_0 | r, n) = \int_0^1 \delta(\tilde{\theta}, \theta) \pi^*(\theta | r, n) d\theta$$

- Fisher's example:  $\mathbf{x} = \{10, 10\}$ ,  
 Test  $\theta_0 = 1/2$ ,  $\theta^*(\mathbf{x}) = 0.9686$   
 $d(\theta_0 | 10, 10) = 5.414 = \log[224]$

Using  $d^* = \log[100] = 4.61$ ,  
 the value  $\theta_0 = 1/2$  is **rejected**.

$$\Pr[\theta < 0.5 | \mathbf{x}] = 0.00016$$

$d(\theta^*   \mathbf{x})$	$\theta^*$	$\Pr[\theta < \theta^*   \mathbf{x}]$
$\log[10]$	0.711	0.00815
$\log[100]$	0.547	0.00043
$\log[1000]$	0.425	0.00003



- *Asymptotic approximation*

- For large samples, the posterior approaches  $N(\theta | \hat{\theta}, (nF(\hat{\theta}))^{-1/2})$ , where  $F(\theta)$  is Fisher's function. Changing variables, the posterior distribution of  $\phi = \phi(\theta) = \int F^{1/2}(\theta) d\theta = 2 \arcsin \sqrt{\theta}$  is approximately normal  $N(\phi | \hat{\phi}, n^{-1/2})$ . Since  $d(\theta, \mathbf{x})$  is invariant,  $d(\theta_0, \mathbf{x}) \approx \frac{1}{2}[1 + n\{\phi(\theta_0) - \phi(\hat{\theta})\}^2]$ .

- *Testing for a majority ( $\theta_0 = 1/2$ )*

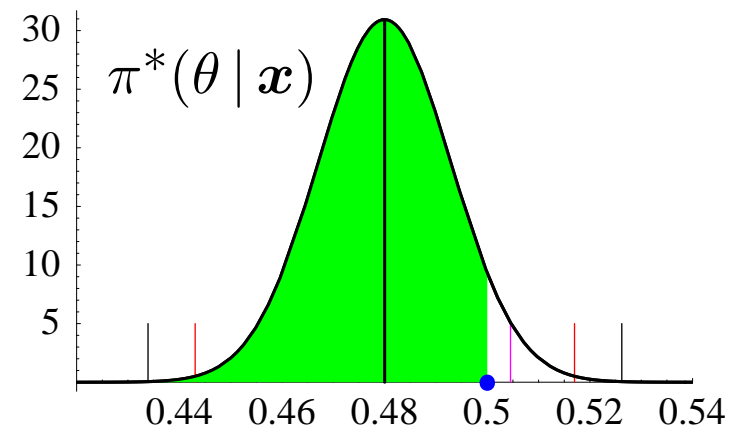
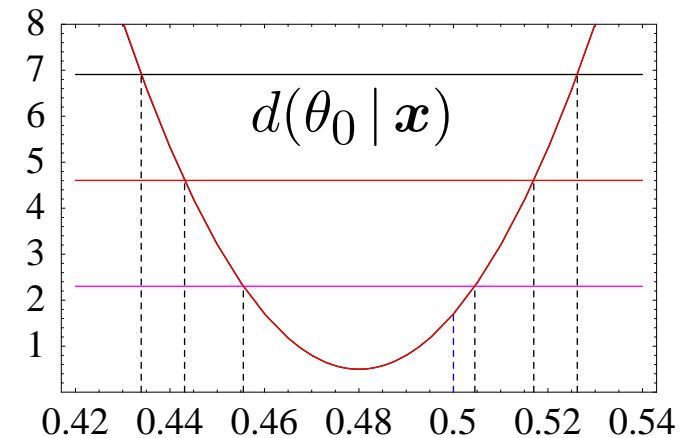
$$\mathbf{x} = \{720, 1500\}, \quad \theta^*(\mathbf{x}) = 0.4800$$

$d(\theta^*   \mathbf{x})$	$R = (\theta_0^*, \theta_1^*)$	$\Pr[\theta \in R   \mathbf{x}]$
$\log[10]$	(0.456, 0.505)	0.9427
$\log[100]$	(0.443, 0.517)	0.9959
$\log[1000]$	(0.434, 0.526)	0.9997

Very mild evidence against  $\theta = 0.5$ :

$$d(0.5 | 720, 1500) = 1.67$$

$$\Pr(\theta < 0.5 | 720, 1500) = 0.9393$$



# Basic References

Many available on line at [www.uv.es/bernardo](http://www.uv.es/bernardo)

- *Introductions*

Bernardo, J. M. and Ramón, J. M. (1998).

An introduction to Bayesian reference analysis. *The Statistician* **47**, 1–35.

Bernardo, J. M. (2003). Bayesian Statistics.

*Encyclopedia of Life Support Systems (EOLSS):*

*Probability and Statistics*, (R. Viertl, ed). Oxford, UK: UNESCO.

Bernardo, J. M. (2005). Reference Analysis.

*Handbook of Statistics* **25** (D. K. Dey and C. R. Rao eds.)

Amsterdam: Elsevier, 17–90

- *Textbooks*

Gelman, A., Carlin, J. B., Stern, H. and Rubin, D. B. (2003).

*Bayesian Data Analysis* (2nd ed.) New York: CRC Press.

Bernardo, J. M. and Smith, A. F. M. (1994).

*Bayesian Theory*. Chichester: Wiley.

2nd ed. to appear in June 2006

- *Research papers on reference analysis (cronological order)*

Bernardo, J. M. (1979). Reference posterior distributions for Bayesian inference. *J. Roy. Statist. Soc. B* **41**, 113—147, (with discussion).

Reprinted in *Bayesian Inference* (N. G. Polson and G. C. Tiao, eds.), Brookfield, VT: Edward Elgar, (1995), 229—263.

Berger, J. O. and Bernardo, J. M. (1992). On the development of reference priors. *Bayesian Statistics 4* (J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.) Oxford: University Press, 35–60 (with discussion).

Bernardo, J. M. (1997) . Noninformative priors do not exist. *J. Statist. Planning and Inference* **65**, 159—189, (with discussion).

Bernardo, J. M. and Rueda, R. (2002). Bayesian hypothesis testing: A reference approach. *Internat. Statist. Rev.* **70**, 351–372.

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.

Bernardo, J. M. (2005). Intrinsic credible regions: An objective Bayesian approach to interval estimation. *Test* **14**, 317-384 (invited paper, with discussion).

- **Valencia International Meetings on Bayesian Statistics**

Sponsored by the University of Valencia. Held every four years in Spain. World forums on research and applications of Bayesian analysis.

## **8th Valencia International Meeting on Bayesian Statistics**

**Benidorm (Alicante), June 1st – 6th 2006**

**[www.uv.es/valenciameeting](http://www.uv.es/valenciameeting)**

- **Valencia Mailing List**

The **Valencia Mailing List** contains about 2,000 entries of people interested in **Bayesian Statistics**. It sends information about the Valencia Meetings and other material of interest to the Bayesian community.

If you do not belong to the list and want to be included, please send your e-mail to **<[valenciameeting@uv.es](mailto:valenciameeting@uv.es)>**

- **José-Miguel Bernardo contact data**

**<[jose.m.bernardo@uv.es](mailto:jose.m.bernardo@uv.es)>**

**[www.uv.es/bernardo](http://www.uv.es/bernardo)**