
Prior and Hierarchical Modeling

*Prof. Nicholas Zabaras
University of Notre Dame
Notre Dame, IN, USA*

*Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>*

September 11, 2017

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)



References

- C P Robert, [The Bayesian Choice: From Decision-Theoretic Motivations to Computational Implementation](#), Springer-Verlag, NY, 2001 ([online resource](#))
- A Gelman, JB Carlin, HS Stern and DB Rubin, [Bayesian Data Analysis](#), Chapman and Hall CRC Press, 2nd Edition, 2003.
- J M Marin and C P Robert, [The Bayesian Core](#), Spring Verlag, 2007 ([online resource](#))
- D. Sivia and J Skilling, [Data Analysis: A Bayesian Tutorial](#), Oxford University Press, 2006.
- Bayesian Statistics for Engineering, [Online Course at Georgia Tech](#), B. Vidakovic.
- Kevin Murphy, [Machine Learning, A probabilistic Perspective](#), Chapter 5.



Contents

- Prior modeling, Conjugate priors , Exponential families
- Linearity of the Posterior Mean, Example: Gaussian with unknown mean and variance
- Extension to Multivariate Gaussians, Poisson with unknown mean
- Mixture of conjugate priors, Limitations of Conjugate Priors
- MaxEnt priors, Non-informative priors
- Translation and Scale invariance
- Improper priors, Jeffrey's prior
- Pros and Cons of improper priors
- Lack of robustness of the normal prior
- Hierarchical Bayesian Models, Empirical Bayes



Selection of Prior Distribution

- Once the prior distribution is selected, Bayesian inference can be performed almost mechanically.
- A critical point of Bayesian statistics is the choice of the prior.
- Seldom there is enough “subjective information” to lead to an ‘exact’ determination of the prior distribution.
- Selection of prior includes subjectivity
 - ✓ Subjectivity does not imply being unscientific – one can use scientific information to guide the specification of priors.
 - ✓ We will review some of the work on uninformative and robust priors.



Prior Selection

The prior distribution is a key to Bayesian inference

- The available prior information is seldom precise enough to lead to determination of *the* prior distribution

Strategies for prior determination

- Use a partition of Θ in intervals and determine the probability of each interval. Then approach π by an histogram
- Alternatively, select significant elements of Θ , evaluate their respective likelihoods and deduce a likelihood curve proportional to π
- Both of these approaches don't work for unbounded Θ (quite difficult to subjectively estimate probabilities at the tails of the distribution).
- When no information is available on θ , use the marginal distribution of x to derive information on π ,

$$m(x) = \int_{\Theta} f(x | \theta) \pi(\theta) d\theta$$

- This perspective is the core of empirical & hierarchical Bayes models.



Informative Priors

- The prior is a tool summarizing the available information on a phenomenon of interest, as well as the uncertainty related with this information.
- Informative priors convey specific and definite information about parameters θ associated with the random phenomenon.
- Pre-existing evidence which has already been taken into account is part of the informative priors. This information can be based on historical data, insight or personal beliefs.
- Typical subgroups of informative priors
 - conjugate, non-conjugate
 - exponential families
 - maximum entropy priors



Conjugate Priors

- Consider a class of probability distribution P . For every prior $\pi(\theta) \in P$, if the posterior distribution $\pi(\theta|x)$ belongs to P and the likelihood $f(x|\theta)$ to a family F , then the P class is **conjugate** for F .
- Conjugate priors are analytically tractable. Finding the posterior reduces to an updating of the corresponding parameters of the prior.
- Consider a coin flipping example:

- Let θ the probability that the coin will draw heads
- Prior $\theta \sim \mathcal{Be}(a,b)$
- Data: the coin flipped n times with n_H of those were heads (binomial)
- Posterior:

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int_0^1 f(x|\theta)\pi(\theta)d\theta} = \frac{\theta^{a+n_H-1}(1-\theta)^{b+n-n_H-1}}{\text{beta}(a+n_H, b+n-n_H)} = \mathcal{Be}(a+n_H, b+n-n_H)$$

- The role of conjugate priors is generally to provide a first approximation to the adequate prior distribution which should be followed by a robustness analysis.



Conjugate Priors

A family \mathcal{F} of probability distributions on Θ is conjugate for a likelihood function $f(x|\theta)$ if, for every $\pi \in \mathcal{F}$, the posterior distribution $\pi(\theta|x)$ also belongs to \mathcal{F}

Of interest is the case when \mathcal{F} is parameterized : switching from prior to posterior distribution results in updating of the corresponding parameters.

C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter 3](#) (full text available)



Conjugate Priors: Justifications

- Limited/finite information conveyed by x
- Preservation of the structure of $\pi(\theta)$
- Exchangeability motivations
- Device of virtual past observations (prior plays the role of past virtual observations)
- Linearity of some estimators
- Tractability and simplicity
- Conjugate priors can be used as first approximations to adequate priors, backed up by robustness analysis
- Their restrictive nature can be attenuated by using hyperpriors on the hyperparameters themselves

Exponential Families

- Conjugate prior distributions are usually associated with **Exponential Families**, a class of probability distributions sharing a certain form as specified below.
- Suppose \mathbf{x} are observations from the **Exponential Family**

$$f(\mathbf{x} | \theta) = C(\theta)h(\mathbf{x})\exp\{R(\theta) \cdot T(\mathbf{x})\}$$

We call this an **exponential family**.

- When $\Theta \subset \mathbb{R}^k, X \subset \mathbb{R}^k$ and

$$f(\mathbf{x} | \theta) = h(\mathbf{x})\exp\{\theta \cdot \mathbf{x} - \psi(\theta)\}$$

the family is called **natural family** of dimension k .

$T(\mathbf{x})$ are here sufficient statistics.



Exponential Families: Example

- Consider the likelihood function

$$f(x|\theta) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-\theta)^2}{2}\right\}$$

- This is a normal distribution (unknown mean, unit variance). For this case note that:

$$f(\mathbf{x}|\theta) = h(\mathbf{x}) \exp\{\theta \cdot \mathbf{x} - \psi(\theta)\} \quad R(\theta) = \theta ; \quad T(x) = x ; \quad \psi(\theta) = \frac{\theta^2}{2} ; \quad h(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$$

- Consider the normal distribution (unknown mean, unknown variance)

$$f(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

define $\theta = (\mu, \sigma)$, we can then see that

$$f(\mathbf{x}|\theta) = C(\theta)h(\mathbf{x}) \exp\{R(\theta) \cdot T(\mathbf{x})\}$$

$$f(\mathbf{x}|\theta) = h(\mathbf{x}) \exp\{R(\theta) \cdot T(\mathbf{x}) - \psi(\theta)\}$$

$$R(\theta) = \left(\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2}\right)^T ; \quad T(x) = \left(x, -\frac{x^2}{2}\right)^T ; \quad C(\theta) = \frac{1}{\sigma} e^{-\frac{\mu^2}{2\sigma^2}} ;$$

$$\psi(\theta) = \frac{\mu^2}{2\sigma^2} - \log \frac{1}{\sigma} ; \quad h(x) = \frac{1}{\sqrt{2\pi}} .$$



Exponential Families

□ Conjugate distributions for exponential families

➤ Likelihood

$$f(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \exp \{ R(\boldsymbol{\theta}) \cdot T(\mathbf{x}) - \psi(\boldsymbol{\theta}) \}$$

➤ Conjugate Prior

$$\pi(\boldsymbol{\theta} | \mu, \lambda) \propto \exp \{ R(\boldsymbol{\theta}) \cdot \mu - \lambda \psi(\boldsymbol{\theta}) \}, \lambda > 0$$

Hyper Parameters

➤ Posterior

$$\pi(\boldsymbol{\theta} | \mathbf{x}) \propto \exp \{ R(\boldsymbol{\theta}) \cdot [\mu + T(\mathbf{x})] - (\lambda + 1) \psi(\boldsymbol{\theta}) \}$$

$$i.e. \quad \pi(\boldsymbol{\theta} | \mathbf{x}) = \pi(\boldsymbol{\theta} | \mu + T(\mathbf{x}), \lambda + 1)$$



Exponential Families: Example

- Normal distribution (unknown mean, known variance)

Likelihood : $x_1 | \theta \sim \mathcal{N}(\theta, \sigma^2), \sigma^2 = \text{known}, x_1 \in \mathbb{R}$

$$f(x_1 | \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_1-\theta)^2}{2\sigma^2}}$$

- Conjugate prior

$$\theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

- Posterior

$$\theta | x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2), \sigma_1^{-2} = \sigma_0^{-2} + \sigma^{-2}, \mu_1 = \underbrace{\frac{\sigma_0^{-2}\mu_0 + \sigma^{-2}x_1}{\sigma_0^{-2} + \sigma^{-2}}}_{\text{weighted average of the observation } x_1 \text{ and the prior mean}}$$

- Posterior predictive (proof given in an earlier lecture):

$$\pi(x | x_1) = \int \pi(x | \theta) \pi(\theta | x_1) d\theta \sim \int e^{-\frac{(x-\theta)^2}{2\sigma^2}} e^{-\frac{(\theta-\mu_1)^2}{2\sigma_1^2}} d\theta \sim \mathcal{N}(\mu_1, \sigma^2 + \sigma_1^2)$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004

Gaussian With Multiple Observations - Unknown Mean

- Assume we have observations $X_i | \mu \sim \mathcal{N}(\mu, \sigma^2)$ and $\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$
- The posterior is then:

$$\mu | x_1, x_2, \dots, x_n \sim \mathcal{N}(\mu_n, \sigma_n^2),$$

$$\frac{1}{\sigma_n^2} = \frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \Rightarrow \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n \sigma_0^2 + \sigma^2} = \frac{\sigma^2}{n + \frac{\sigma^2}{\sigma_0^2}}$$

$$\mu_n = \sigma_n^2 \left(\frac{\sum_{i=1}^n x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) = \sigma_n^2 \left(\frac{\sum_{i=1}^n x_i + \mu_0 (\sigma^2 / \sigma_0^2)}{\sigma^2} \right)$$

- One can think of the prior as n_0 virtual observations with $n_0 = \frac{\sigma^2}{\sigma_0^2}$ and

$$\sigma_n^2 = \frac{\sigma^2}{n + n_0}, \mu_n = \frac{\sum_{i=1}^n x_i + n_0 \mu_0}{n + n_0}$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Standard Exponential Families

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Normal $N(\theta, \sigma^2)$	Normal $\mathcal{N}(\mu, \tau^2)$	$\mathcal{N}(\rho(\sigma^2\mu + \tau^2x), \rho\sigma^2\tau^2)$ $\rho^{-1} = \sigma^2 + \tau^2$
Poisson $\mathcal{P}(\theta)$	Gamma $\mathcal{P}(\theta)\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + x, \beta + 1)$
Gamma $\mathcal{G}(v, \theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + v, \beta + x)$
Binomial $\mathcal{B}(n, \theta)$	Beta $\mathcal{B}\text{e}(\alpha, \beta)$	$\mathcal{B}\text{e}(\alpha + x, \beta + n - x)$
Negative Binomial $\mathcal{N}\text{eg}(m, \theta)$	Beta $\mathcal{B}\text{e}(\alpha, \beta)$	$\mathcal{B}\text{e}(\alpha + m, \beta + x)$
Multinomial $\mathcal{M}_k(\theta_1, \dots, \theta_k)$	Dirichlet $\mathcal{D}(\alpha_1, \dots, \alpha_k)$	$\mathcal{D}(\alpha_1 + x_1, \dots, \alpha_k + x_k)$
Normal $\mathcal{N}(\mu, 1/\theta)$	Gamma $\mathcal{G}\text{a}(\alpha, \beta)$	$\mathcal{G}(\alpha + 0.5, \beta + (\mu - x)^2/2)$

Linearity of the Posterior Mean

- Consider the natural prior

$$\pi(\theta | \mu, \lambda) \propto \exp\{\theta \cdot \mu - \lambda \psi(\theta)\}, \lambda > 0$$

- With $\mu \in X$, then

$$\mathbb{E}^\pi [\nabla \psi(\theta)] = \frac{\mu}{\lambda}, \text{ where } \nabla \psi(\theta) = \left(\frac{\partial \psi(\theta)}{\partial \theta_1}, \dots, \frac{\partial \psi(\theta)}{\partial \theta_k} \right)$$

- Consider that i.i.d. data x_1, x_2, \dots, x_n . Then with likelihood $f(x, \theta)$

$$\mathbb{E}^\pi [\nabla \psi(\theta) | x_1, x_2, \dots, x_n] = \frac{\mu + \lambda \bar{x}}{\lambda + n}$$

Gaussian with Unknown Mean and Variance

- Assume we have observations $x_i | \mu \sim \mathcal{N}(\mu, \sigma^2)$. The likelihood has the form:

$$f(x|\mu, \sigma^2) \propto \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) = \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \right]\right) \Rightarrow$$

$$f(x|\mu, \sigma^2) \propto \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} [(n-1)s^2 + n(\bar{x} - \mu)^2]\right)$$

- Here

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i , \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- The form of the likelihood and the subsequent discussion shows that the conjugate prior density must also have the product form $\pi(\sigma^2)\pi(\mu|\sigma^2)$ where the marginal distribution of σ^2 is scaled inverse- χ^2 and the conditional distribution $\pi(\mu|\sigma^2)$ is normal (so that marginally μ has a Student-t distribution)

$$\mu|\sigma^2 \sim \mathcal{N}(\mu_0, \frac{\sigma^2}{K_0}), \quad \sigma^2 \sim \text{Inv-}\chi^2(v_0, \sigma_0^2)$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Gamma and Inverse Gamma Distributions

- Let us review some of the needed distributions:

Gamma

$$\theta \sim \text{Gamma}(\alpha, \beta)$$

shape $\alpha > 0$

$$p(\theta) = \text{Gamma}(\theta | \alpha, \beta)$$

inverse scale $\beta > 0$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad \theta > 0$$

$$E(\theta) = \frac{\alpha}{\beta}$$

$$\text{var}(\theta) = \frac{\alpha}{\beta^2}$$

$$\text{mode}(\theta) = \frac{\alpha-1}{\beta}, \text{ for } \alpha \geq 1$$

Inverse-gamma

$$\theta \sim \text{Inv-gamma}(\alpha, \beta)$$

shape $\alpha > 0$

$$p(\theta) = \text{Inv-gamma}(\theta | \alpha, \beta)$$

scale $\beta > 0$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$$

$$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$$

$$\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$$

$$\text{mode}(\theta) = \frac{\beta}{\alpha+1}$$

Bayesian Data Analysis, A. Gelman, J. Carlin, H. Stern and D. Rubin, 2004



Chi-Square χ^2 Distribution

Chi-square	$\theta \sim \chi_{\nu}^2$ $p(\theta) = \chi_{\nu}^2(\theta)$	degrees of freedom $\nu > 0$
------------	--	------------------------------

$p(\theta) = \frac{2^{-\nu/2}}{\Gamma(\nu/2)} \theta^{\nu/2-1} e^{-\theta/2}, \quad \theta > 0$ same as Gamma($\alpha = \frac{\nu}{2}$, $\beta = \frac{1}{2}$)	$E(\theta) = \nu$, $\text{var}(\theta) = 2\nu$ $\text{mode}(\theta) = \nu - 2$, for $\nu \geq 2$
--	---

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Scaled Inverse χ^2 Distribution

Inverse-chi-square

$$\theta \sim \text{Inv-}\chi_{\nu}^2$$
$$p(\theta) = \text{Inv-}\chi_{\nu}^2(\theta)$$

degrees of freedom $\nu > 0$

$$p(\theta) = \frac{2^{-\nu/2}}{\Gamma(\nu/2)} \theta^{-(\nu/2+1)} e^{-1/(2\theta)}, \quad \theta > 0$$

same as Inv-gamma($\alpha = \frac{\nu}{2}$, $\beta = \frac{1}{2}$)

$$E(\theta) = \frac{1}{\nu-2}, \text{ for } \nu > 2$$

$$\text{var}(\theta) = \frac{2}{(\nu-2)^2(\nu-4)}, \nu > 4$$

$$\text{mode}(\theta) = \frac{1}{\nu+2}$$

Scaled inverse-chi-square

$$\theta \sim \text{Inv-}\chi^2(\nu, s^2)$$
$$p(\theta) = \text{Inv-}\chi^2(\theta|\nu, s^2)$$

degrees of freedom $\nu > 0$

scale $s > 0$

$$p(\theta) = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} s^{\nu} \theta^{-(\nu/2+1)} e^{-\nu s^2/(2\theta)}, \quad \theta > 0$$

same as Inv-gamma($\alpha = \frac{\nu}{2}$, $\beta = \frac{\nu}{2}s^2$)

$$E(\theta) = \frac{\nu}{\nu-2} s^2$$

$$\text{var}(\theta) = \frac{2\nu^2}{(\nu-2)^2(\nu-4)} s^4$$

$$\text{mode}(\theta) = \frac{\nu}{\nu+2} s^2$$

Bayesian Data Analysis, A. Gelman, J. Carlin, H. Stern and D. Rubin, 2004



Gaussian with Unknown Mean and Variance

- The joint prior density corresponds to:

$$\pi(\mu, \sigma^2) \propto \underbrace{\left(\sigma^2\right)^{-1/2} \exp\left(-\frac{1}{2\sigma^2} \kappa_0 (\mu_0 - \mu)^2\right)}_{\mathcal{N}(\mu_0, \frac{\sigma^2}{\kappa_0})} \underbrace{\left(\sigma^2\right)^{-(v_0/2+1)} \exp\left(-\frac{1}{2\sigma^2} v_0 \sigma_0^2\right)}_{\text{Inv-}\chi^2(v_0, \sigma_0^2)}$$
$$\propto \underbrace{\sigma^{-1} \left(\sigma^2\right)^{-(v_0/2+1)} \exp\left(-\frac{1}{2\sigma^2} \left[v_0 \sigma_0^2 + \kappa_0 (\mu_0 - \mu)^2\right]\right)}_{\mathcal{N-Inv-}\chi^2(\mu_0, \sigma_0^2 / \kappa_0; v_0, \sigma_0^2)}$$

- This is labeled the $\mathcal{N-Inv-\chi^2}(\mu_0, \sigma_0^2 / \kappa_0; v_0, \sigma_0^2)$ distribution. Has 4 parameters – the location and scale of μ and the degrees of freedom and scale of σ^2 .

Gaussian with Unknown Mean and Variance

- The joint prior density corresponds to:

$$\pi(\mu, \sigma^2) \propto \underbrace{\sigma^{-1} (\sigma^2)^{-(\nu_0/2+1)} \exp\left(-\frac{1}{2\sigma^2} [\nu_0 \sigma_0^2 + \kappa_0 (\mu_0 - \mu)^2]\right)}_{\mathcal{N}\text{-}Inv-\chi^2(\mu_0, \sigma_0^2/\kappa_0; \nu_0, \sigma_0^2)}$$

- The appearance of σ^2 in the conditional distribution of $\mu | \sigma^2$ means that μ and σ^2 are dependent in their joint conjugate prior density.
- For example, if σ^2 is large, then a high-variance prior distribution is induced on μ .
- It makes sense for the prior variance of the mean to be tied to σ^2 , which is the sampling variance of the observation x . This way, prior belief about μ , is calibrated by the scale of measurement of x and is equivalent to κ_0 prior measurements on this scale ([see earlier analysis](#)).

$$\mu | \sigma^2 \sim \mathcal{N}(\mu_0, \frac{\sigma^2}{\kappa_0}), \quad \sigma^2 \sim Inv-\chi^2(\nu_0, \sigma_0^2)$$



Gaussian with Unknown Mean and Variance

- Multiplying the likelihood & prior distributions gives the joint posterior as:

$$\pi(\mu, \sigma^2 | x) \propto \sigma^{-1} (\sigma^2)^{-(\nu_0/2+1)} \exp\left(-\frac{1}{2\sigma^2} [\nu_0 \sigma_0^2 + \kappa_0 (\mu_0 - \mu)^2]\right)$$

$$\times \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} [(n-1)s^2 + n(\bar{x} - \mu)^2]\right) =$$

$$\propto \sigma^{-1} (\sigma^2)^{-(\nu_n/2+1)} \exp\left(-\frac{1}{2\sigma^2} [\nu_n \sigma_n^2 + \kappa_n (\mu_n - \mu)^2]\right)$$

$$N - Inv - \chi^2(\mu_n, \sigma_n^2 / \kappa_n; \nu_n, \sigma_n^2)$$

where:

$$\mu_n = \frac{\kappa_0}{\kappa_0 + n} \mu_0 + \frac{n}{\kappa_0 + n} \bar{x}$$

$$\kappa_n = \kappa_0 + n, \quad \nu_n = \nu_0 + n$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{x} - \mu_0)^2$$

- In the derivation above, we used the following:

$$n(\bar{x} - \mu)^2 + \kappa_0(\mu_0 - \mu)^2 = \frac{\kappa_0 n}{\kappa_0 + n} (\bar{x} - \mu_0)^2 + (\kappa_0 + n) \left(\mu - \frac{\kappa_0}{\kappa_0 + n} \mu_0 - \frac{n\bar{x}}{\kappa_0 + n} \right)^2$$

Gaussian with Unknown Mean and Variance

- Using the joint posterior for fixed σ^2 , the conditional posterior distribution for μ is:

$$\mu | \sigma^2, x \sim \mathcal{N}(\mu_n, \sigma^2 / \kappa_n) = \mathcal{N}\left(\frac{\kappa_0 \mu_0 + n \bar{x}}{\kappa_0 + n}, \frac{\sigma^2}{\kappa_0 + n}\right)$$

- Note that this result confirms the posterior obtained in the earlier example of a Gaussian with known variance.
- Similarly the marginal posterior of σ^2 is given as:

$$\pi(\sigma^2 | x) \propto \int_{-\infty}^{\infty} (\sigma^2)^{-(\nu_n/2+3/2)} \exp\left(-\frac{\nu_n \sigma_n^2}{2\sigma^2} \left[1 + \frac{\kappa_n (\mu_n - \mu)^2}{\nu_n \sigma_n^2}\right]\right) d\mu =$$

Use: $\int_{-\infty}^{\infty} e^{-a(1+b(\mu-\mu_0))^2} d\mu = \frac{\sqrt{\pi}}{2\sqrt{ab}} e^{-a} \operatorname{erf}\left(\sqrt{ab}x - \frac{ab\mu_0}{\sqrt{ab}}\right)|_{-\infty}^{+\infty} = \frac{\sqrt{\pi}}{\sqrt{ab}} e^{-a}$

$$\propto (\sigma^2)^{-(\nu_n/2+3/2)} \exp\left(-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right) \sigma = (\sigma^2)^{-(\nu_n/2+1)} \exp\left(-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right)$$

Scaled
inverse-
chi-square

$$\theta \sim \operatorname{Inv-}\chi^2(\nu, s^2)$$
$$p(\theta) = \operatorname{Inv-}\chi^2(\theta | \nu, s^2)$$

degrees of freedom $\nu > 0$
scale $s > 0$

$$p(\theta) = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} s^\nu \theta^{-(\nu/2+1)} e^{-\nu s^2/(2\theta)}, \quad \theta > 0$$

same as Inv-gamma($\alpha = \frac{\nu}{2}$, $\beta = \frac{\nu}{2}s^2$)

$$E(\theta) = \frac{\nu}{\nu-2} s^2$$
$$\operatorname{var}(\theta) = \frac{2\nu^2}{(\nu-2)^2(\nu-4)} s^4$$
$$\operatorname{mode}(\theta) = \frac{\nu}{\nu+2} s^2$$

$$\sigma^2 | x \propto \operatorname{Inv-}\chi^2(\nu_n, \sigma_n^2)$$

Gaussian with Unknown Mean and Variance

- The marginal posterior of μ can be computed as:

$$\pi(\mu | x) \propto \int_0^{\infty} (\sigma^2)^{-(\nu_n/2+3/2)} \exp \left\{ -\frac{1}{2\sigma^2} \underbrace{\left[\nu_n \sigma_n^2 + \kappa_n (\mu_n - \mu)^2 \right]}_z \right\} d\sigma^2 =$$

$\mathcal{N}\text{-}Inv\text{-}\chi^2(\mu_n, \sigma_n^2/\kappa_n; \nu_n, \sigma_n^2)$

$$\begin{aligned} \text{Set } z &= \frac{A}{2\sigma^2} \Rightarrow \\ dz &= -\frac{A}{2\sigma^4} d\sigma^2 = \\ &= -\frac{A}{2A^2} 4z^2 d\sigma^2 = -\frac{2z^2}{A} d\sigma^2 \end{aligned}$$

$$\propto \int_0^{\infty} \left(\frac{A}{2z} \right)^{-(\nu_n/2+3/2)} e^{-z} \left(-\frac{A}{2z^2} \right) dz \propto A^{-(\nu_n/2+1/2)} \int_0^{\infty} z^{(\nu_n-1)/2} e^{-z} dz$$

$$\propto \left[\nu_n \sigma_n^2 + \kappa_n (\mu_n - \mu)^2 \right]^{-(\nu_n/2+1/2)} \propto \left[1 + \frac{\kappa_n}{\nu_n \sigma_n^2} (\mu_n - \mu)^2 \right]^{-(\nu_n/2+1/2)}$$

$$\pi(\mu | x) \propto \underbrace{\left[1 + \frac{\kappa_n (\mu - \mu_n)^2}{\nu_n \sigma_n^2} \right]^{-(\nu_n+1)/2}}_{t_{\nu_n}(\mu | \mu_n, \sigma_n^2 / \kappa_n)}$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)



Student-t Distribution

Student- <i>t</i>	$\theta \sim t_\nu(\mu, \sigma^2)$	degrees of freedom $\nu > 0$
	$p(\theta) = t_\nu(\theta \mu, \sigma^2)$	location μ
	t_ν is short for $t_\nu(0, 1)$	scale $\sigma > 0$

$p(\theta) = \frac{\Gamma((\nu+1)/2)}{\Gamma(\nu/2)\sqrt{\nu\pi}\sigma} (1 + \frac{1}{\nu}(\frac{\theta-\mu}{\sigma})^2)^{-(\nu+1)/2}$	$E(\theta) = \mu$, for $\nu > 1$
	$\text{var}(\theta) = \frac{\nu}{\nu-2}\sigma^2$, for $\nu > 2$
	$\text{mode}(\theta) = \mu$

- Compare this with our result:

$$\pi(\mu | x) \propto \underbrace{\left[1 + \frac{\kappa_n (\mu - \mu_n)^2}{\nu_n \sigma_n^2} \right]^{-(\nu_n + 1)/2}}_{t_{\nu_n}(\mu | \mu_n, \sigma_n^2 / \kappa_n)}$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Generalizing to Multivariate Gaussians

- Suppose you have $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \sim (\text{i.i.d}) \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- We do not know $\boldsymbol{\mu}$ or $\boldsymbol{\Sigma}$
- When both sets of parameters are unknown, a conjugate family of priors is one in which

$$\boldsymbol{\Sigma} \sim \mathcal{IW}(\nu, \boldsymbol{\Lambda}^{-1})$$

and

$$\boldsymbol{\mu} | \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\eta}, \boldsymbol{\Sigma} / \kappa)$$

- The Wishart distribution is a multivariate analog of the Gamma distribution. If matrix \mathbf{U} has the Wishart distribution, then \mathbf{U}^{-1} has the inverse-Wishart distribution.
- The quantity ν is a positive scalar, while $\boldsymbol{\Lambda}$ is a positive definite matrix. They play roles analogous to those played by α and β , respectively, in the Gamma distribution.
- The other parameters of the prior are the mean vector $\boldsymbol{\eta}$ and κ , the latter of which represents the ``a priori number of observations''.



Wishart Distribution

Wishart	$W \sim \text{Wishart}_\nu(S)$ $p(W) = \text{Wishart}_\nu(W S)$ (implicit dimension $k \times k$)	degrees of freedom ν symmetric, pos. definite $k \times k$ scale matrix S
---------	--	---

$$p(W) = \left(2^{\nu k/2} \pi^{k(k-1)/4} \prod_{i=1}^k \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1} \times |S|^{-\nu/2} |W|^{(\nu-k-1)/2} \times \exp\left(-\frac{1}{2}\text{tr}(S^{-1}W)\right), W \text{ pos. definite}$$

Gamma	$\theta \sim \text{Gamma}(\alpha, \beta)$ $p(\theta) = \text{Gamma}(\theta \alpha, \beta)$	shape $\alpha > 0$ inverse scale $\beta > 0$
-------	---	---

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}, \quad \theta > 0$$
$$\begin{aligned} \text{E}(\theta) &= \frac{\alpha}{\beta} \\ \text{var}(\theta) &= \frac{\alpha}{\beta^2} \\ \text{mode}(\theta) &= \frac{\alpha-1}{\beta}, \text{ for } \alpha \geq 1 \end{aligned}$$

Bayesian Data Analysis, A. Gelman, J. Carlin, H. Stern and D. Rubin, 2004



Inverse Wishart Distribution

Inverse-Wishart

$W \sim \text{Inv-Wishart}_\nu(S^{-1})$ degrees of freedom ν
 $p(W) = \text{Inv-Wishart}_\nu(W|S^{-1})$ symmetric, pos. definite
(implicit dimension $k \times k$) $k \times k$ scale matrix S

$$p(W) = \left(2^{\nu k/2} \pi^{k(k-1)/4} \prod_{i=1}^k \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1} \times |S|^{\nu/2} |W|^{-(\nu+k+1)/2} \times \exp\left(-\frac{1}{2} \text{tr}(SW^{-1})\right), W \text{ pos. definite}$$
$$\text{E}(W) = (\nu - k - 1)^{-1} S$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Generalizing to Multivariate Gaussians

- The prior has the following form:

$$\pi(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(\nu+d)/2+1}$$

$$\times \exp\left(-\frac{1}{2}\text{tr}(\boldsymbol{\Lambda}\boldsymbol{\Sigma}^{-1}) - \frac{\kappa}{2}(\boldsymbol{\mu} - \boldsymbol{\eta})^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \boldsymbol{\eta})\right)$$

- The posterior has the same form but with different parameters
 - ν , $\boldsymbol{\Lambda}$, $\boldsymbol{\eta}$, and κ , now become ν_n , $\boldsymbol{\Lambda}_n$, $\boldsymbol{\mu}_n$, and κ_n , respectively where:

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$\boldsymbol{\mu}_n = \frac{\kappa\boldsymbol{\eta} + n\bar{\mathbf{x}}}{\kappa + n}$$

$$\nu_n = \nu + n$$

$$\kappa_n = \kappa + n$$

$$\mathbf{S}^2 = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T$$

$$\boldsymbol{\Lambda}_n = \boldsymbol{\Lambda} + n\mathbf{S}^2 + \frac{(\bar{\mathbf{x}} - \boldsymbol{\eta})(\bar{\mathbf{x}} - \boldsymbol{\eta})^T}{1/\kappa + 1/n}$$

Generalizing to Multivariate Gaussians

- We thus have the following properties of component distributions of the posterior
- The conditional distribution of μ given Σ and the data is a Gaussian:

$$\mu \mid \Sigma, \mathbf{x} \sim \mathcal{N}(\mu_n, \Sigma / \kappa_n)$$

- The marginal posterior of μ is a multivariate Student t.
- The marginal posterior of Σ is Inverse-Wishart:

$$\Sigma \sim \mathcal{IW}(v_n, \Lambda_n^{-1})$$



Multivariate Student t Distribution

Multivariate
Student-*t*

$$\begin{aligned}\theta &\sim t_\nu(\mu, \Sigma) \\ p(\theta) &= t_\nu(\theta|\mu, \Sigma) \\ (\text{implicit dimension } d)\end{aligned}$$

degrees of freedom $\nu > 0$
location $\mu = (\mu_1, \dots, \mu_d)$
symmetric, pos. definite
 $d \times d$ scale matrix Σ

$$p(\theta) = \frac{\Gamma((\nu+d)/2)}{\Gamma(\nu/2)\nu^{d/2}\pi^{d/2}} |\Sigma|^{-1/2} \times (1 + \frac{1}{\nu}(\theta - \mu)^T \Sigma^{-1} (\theta - \mu))^{-(\nu+d)/2}$$

$$\begin{aligned}\mathbb{E}(\theta) &= \mu, \text{ for } \nu > 1 \\ \text{var}(\theta) &= \frac{\nu}{\nu-2} \Sigma, \text{ for } \nu > 2 \\ \text{mode}(\theta) &= \mu\end{aligned}$$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Generalizing to Multivariate Gaussians

- The Jeffreys noninformative prior (to be introduced later in this lecture) when μ and Σ are unknown is:

$$\pi(\mu, \Sigma) \propto |\Sigma|^{-(d+1)/2}$$

- This is the limit of the previous conjugate prior as $\kappa \rightarrow 0, v \rightarrow -1$ and $|\Lambda| \rightarrow 0$.
- The resulting posterior is proper given by:

$$\Sigma | x \sim \mathcal{IW}\left(n-1, (nS^2)^{-1}\right)$$

and

$$\mu | \Sigma, x \sim \mathcal{N}\left(\bar{x}, \Sigma / n\right)$$

- More results can be found on our multivariate Gaussian lecture.



Example: Poisson with Unknown Mean

- Assume we have some counting observations $X_i \stackrel{i.i.d.}{\sim} \mathcal{P}(\theta)$ i.e.

$$f(x_i | \theta) = e^{-\theta} \frac{\theta^{x_i}}{x_i!}, x_i = 0, 1, 2, \dots$$

- Assume a Gamma prior for θ , i.e.

$$\pi(\theta) = \text{Ga}(\theta; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}$$

- The posterior is then:

$$\pi(\theta | x_1, x_2, \dots, x_n) = \text{Ga}\left(\theta; \alpha + \sum_{i=1}^n x_i, \beta + n\right)$$

- We can now think of the prior as having `` β virtual observations who sum to α ''.

Mixture of Conjugate Priors

- Robust priors are useful, but can be computationally expensive to use.
- Conjugate priors simplify the computation, but are often not robust, and not flexible enough to encode our prior knowledge.
- A mixture of conjugate priors is also conjugate and can approximate any kind of prior. Thus such priors provide a good compromise between computational convenience and flexibility.
- Example: to model coin tosses, we can take a prior which is a mixture of two beta distributions to model coin tosses.

$$p(\theta) = 0.5 \text{Beta}(\theta | 20, 20) + 0.5 \text{Beta}(\theta | 30, 10)$$

- If θ comes from the first distribution, the coin is fair, but if it comes from the second, it is biased towards heads.



Mixture of Conjugate Priors

- If we have a prior distribution which is a mixture of conjugate distributions to a given likelihood, then the posterior is in closed form and is a mixture of conjugate distributions, i.e. with

$$\pi(\theta) = \sum_{i=1}^K w_i \pi_i(\theta)$$

the following posterior is obtained:

$$\pi(\theta | x) = \frac{\sum_{i=1}^K w_i \pi_i(\theta) f(x | \theta)}{\underbrace{\sum_{i=1}^K w_i \int \pi_i(\theta) f(x | \theta) d\theta}_A} = \sum_{i=1}^K \frac{w_i}{A} \pi_i(\theta) f(x | \theta) = \sum_{i=1}^K w'_i \frac{\pi_i(\theta) f(x | \theta)}{\int \pi_i(\theta) f(x | \theta) d\theta} = \sum_{i=1}^K w'_i \pi_i(\theta | x)$$

where:

$$w'_i = \frac{w_i \int \pi_i(\theta) f(x | \theta) d\theta}{\sum_{i=1}^K w_i \int \pi_i(\theta) f(x | \theta) d\theta}, \quad \sum_{i=1}^K w'_i = 1.$$

- One can approximate arbitrary closely any prior distribution by a mixture of conjugate distributions ([Brown, 1986](#))



Mixture of Conjugate Priors

- As an example, suppose we use the mixture prior

$$p(\theta) = 0.5 \text{Beta}(\theta | a_1, b_1) + 0.5 \text{Beta}(\theta | a_2, b_2)$$

$a_1 = b_1 = 20, a_2 = b_2 = 10$, we observe N_1 heads, N_0 tails

- The posterior becomes $p(\theta | \mathcal{D}) = p(Z=1 | \mathcal{D}) \text{Beta}(\theta | a_1 + N_1, b_1 + N_0)$
 $+ p(Z=2 | \mathcal{D}) \text{Beta}(\theta | a_2 + N_1, b_2 + N_0)$

- The posterior mixing weights are given as:

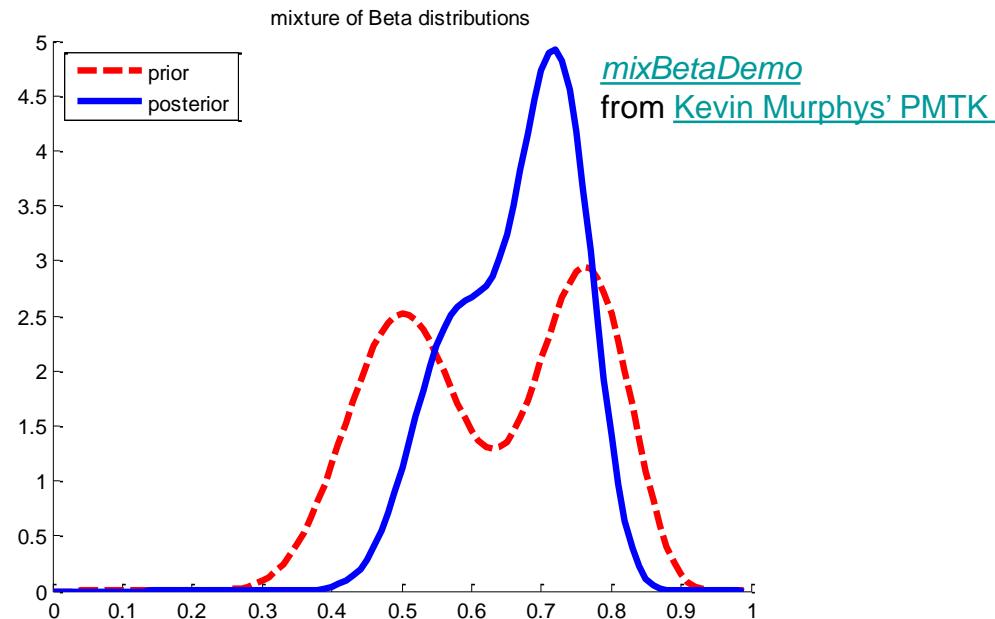
$$p(Z=k | \mathcal{D}) = \frac{p(Z=k)p(\mathcal{D} | Z=k)}{\sum_{k'} p(Z=k')p(\mathcal{D} | Z=k')} = \frac{p(Z=k)p(\mathcal{D} | Z=k)}{p(\mathcal{D})}$$

- If $N_1 = 20$ heads and
 $N_0 = 10$ tails, then, using

$$p(\mathcal{D} | Z=1) = \binom{N}{N_1} \frac{B(a_1 + N_1, b_1 + N_0)}{B(a_1, b_1)}$$

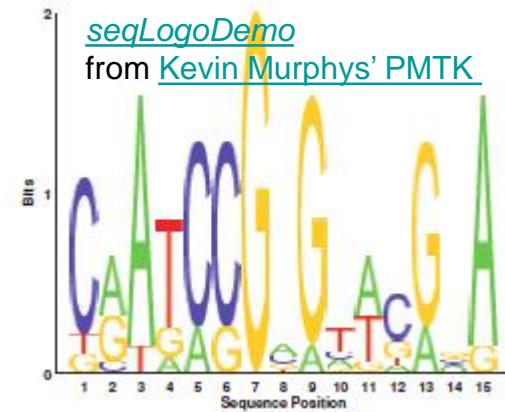
- The posterior finally becomes

$$\begin{aligned} p(\theta | \mathcal{D}) &= \\ &0.346 \text{Beta}(\theta | 40, 30) + \\ &0.654 \text{Beta}(\theta | 30, 20) \end{aligned}$$



Mixture of Conjugate Priors

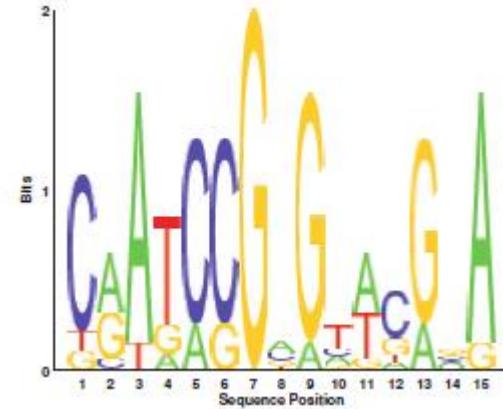
- Dirichlet-multinomial models are widely used in biosequence analysis. Consider the sequence logo problem.
- Suppose we want to find locations which represent **coding regions of the genome**. Such locations often **have the same letter across all sequences** (mostly all A's, or all T's, or all C's, or all G's).
- We believe adjacent locations are conserved together. We let $Z_t = 1$ if location t is conserved, and let $Z_t = 0$ otherwise. We add a dependence between adjacent Z_t variables using a Markov chain.
- To define a likelihood model, $p(\mathbf{N}_t | Z_t)$, where \mathbf{N}_t is the vector of (A,C,G,T) counts for column t . We make this a multinomial distribution with parameter $\boldsymbol{\theta}_t$.
- Since each column has a different distribution, we will want to integrate out $\boldsymbol{\theta}_t$ and thus compute the marginal likelihood $p(\mathbf{N}_t | Z_t)$.



Mixture of Conjugate Priors

$$p(\mathbf{N}_t | Z_t) = \int p(\mathbf{N}_t | \boldsymbol{\theta}_t) p(\boldsymbol{\theta}_t | Z_t) d\boldsymbol{\theta}_t$$

- But what prior should we use for $\boldsymbol{\theta}_t$?
- When $Z_t = 0$ we can use a uniform prior, $p(\boldsymbol{\theta} | Z_t = 0) = \text{Dir}(1, 1, 1, 1)$, but what should we use if $Z_t = 1$?
- If the column is conserved, $Z_t = 1$, it could be a nearly pure column of A's, C's, G's, or T's. A natural approach is to use a mixture of Dirichlet priors, each tilted towards the appropriate corner of the 4-d simplex,



$$p(\boldsymbol{\theta} | Z_t = 1) = 1/4 \text{Dir}(\boldsymbol{\theta} | (10, 1, 1, 1)) + \dots + 1/4 \text{Dir}(\boldsymbol{\theta} | (1, 1, 1, 10))$$

- Since this is conjugate, we can easily compute $p(\mathbf{N}_t | Z_t)$ from the Eq. on the top of the slide ([Brown et al. 1993](#))

Limitations of Conjugate Priors

- The conjugate prior can have a strange shape or be difficult to handle without simulation techniques.^a
- Consider a logistic regression example. The indicator variable $y \in \{0,1\}$ and $x \in \mathbb{R}^k$.

$$\Pr(y=1 | \theta, x) = 1 - \Pr(y=0 | \theta, x) = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}, \quad \Pr(y=0 | \theta, x) = \frac{1}{1 + \exp(\theta^T x)}$$

- For a sample $(y_1, x_1), \dots, (y_n, x_n)$ from the above distribution, the likelihood for n observations is exponential conditional upon x_i 's as

$$f(y_1, \dots, y_n | x_1, \dots, x_n, \theta) = \exp\left(\theta^T \sum_{i=1}^n y_i x_i\right) \prod_{i=1}^n \left(1 + \exp(\theta^T x_i)\right)^{-1}$$

and the resulting rather complicated conjugate prior is:

$$\pi(\theta | \mu, \lambda) \propto \exp(\theta^T \mu) \prod_{i=1}^n \left(1 + \exp(\theta^T x_i)\right)^{-\lambda}$$

^a C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter](#) 3 (full text available)



Summary: Conjugate Priors

PROS.

- Simple to handle, can be interpreted through imaginary observations.
- Considered as the least informative ones.

CONS.

- Not applicable to all likelihood functions.
- Not flexible, cannot account for constraints e.g. $\theta > 0$.
- Approximation by mixtures while feasible is very tedious and thus not used in practice.



Maximum Entropy Priors

- If nothing is known about a distribution except that it belongs to a certain class, then the distribution with the largest entropy should be chosen as the default.^a
- The entropy is defined as

➤ discrete case $\mathbb{H}(\pi) = -\sum_k \pi(\theta_k) \log(\pi(\theta_k))$

- When some statistics (moments) of the prior distribution are known,

$$\mathbb{E}_\pi [g_k(\theta)] = w_k, k = 1, \dots, K$$

the maximum entropy distribution is of the form:

$$\pi(\theta_i) = \frac{\exp\left(\sum_{k=1}^K \lambda_k g_k(\theta_i)\right)}{\sum_j \exp\left(\sum_{k=1}^K \lambda_k g_k(\theta_j)\right)}, \lambda_k = \text{Lagrange multipliers}$$

- However, the constraints may not be compatible, e.g. $\mathbb{E}(\theta^2) \geq \mathbb{E}^2(\theta)$.

^a C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter 3](#) (full text available)



Maximum Entropy Priors

- For the continuous case, we define the entropy as the Kullback-Leibler divergence between π and some invariant non-informative prior for the problem π_0 , i.e.

$$\mathbb{H}(\pi) = - \int \pi_0(\theta) \log \left(\frac{\pi(\theta)}{\pi_0(\theta)} \right) d\theta$$

- As for the discrete case, the maximum entropy distribution is of the form:

$$\pi(\theta) = \frac{\exp \left(\sum_{k=1}^K \lambda_k g_k(\theta) \right) \pi_0(\theta)}{\int \exp \left(\sum_{k=1}^K \lambda_k g_k(\theta) \right) \pi_0(\theta) d\theta}, \quad \lambda_k = \text{Lagrange multipliers}$$

- The selection of π_0 is not obvious or easy.

Maximum Entropy Priors

□ Examples

- **discrete case** $\theta \in \{1, 2, \dots, n\}$

maximum entropy distribution $\pi(\theta) = \frac{1}{n}$

- **continuous** $\mathbb{E}_\pi(\theta) = \mu$

Maximum entropy distribution $\pi(\theta) \propto e^{\lambda\theta}$ (bad improper prior)

- **continuous case** $\mathbb{E}(\theta) = \mu, \text{var}(\theta) = \sigma^2$

maximum entropy distribution $\pi(\theta) = \mathcal{N}(\theta; \mu, \sigma^2)$

- The prior maximizing the entropy is, in a sense, minimizing the prior information brought through π about θ .

Maximum Entropy Priors

- Consider $\Theta = \{0, 1, 2, \dots, n\}$. Suppose $\mathbb{E}_\pi(\theta) = 5$, then the MaxEnt distribution is:

$$\pi(\theta) = \frac{e^{\lambda_1 \theta}}{\sum_{\theta=0}^{\infty} e^{\lambda_1 \theta}} = \frac{e^{\lambda_1 \theta}}{1 + (e^{\lambda_1}) + (e^{\lambda_1})^2 + \dots} \Rightarrow \pi(\theta) = (1 - e^{\lambda_1}) e^{\lambda_1 \theta}$$

- We need to compute λ_1 to satisfy the constraint.

$$\mathbb{E}_\pi(\theta) = 5 \Rightarrow \sum_{\theta=0}^{\infty} \theta (1 - e^{\lambda_1}) e^{\lambda_1 \theta} = 5 \Rightarrow (1 - e^{\lambda_1}) \sum_{\theta=0}^{\infty} \theta e^{\lambda_1 \theta} = (1 - e^{\lambda_1}) \frac{d}{d\lambda_1} \sum_{\theta=0}^{\infty} e^{\lambda_1 \theta} = 5 \Rightarrow$$

$$(1 - e^{\lambda_1}) \frac{d}{d\lambda_1} \left(\frac{1}{1 - e^{\lambda_1}} \right) = 5 \Rightarrow \frac{e^{\lambda_1}}{1 - e^{\lambda_1}} = 5 \Rightarrow e^{\lambda_1} = \frac{1}{6}$$

- The maximum entropy distribution is thus: $\pi(\theta) = p(1 - p)^\theta$, which is the Geometric density Geo(5/6) with parameter $p = 1 - e^{\lambda_1} = 5/6$.

Noninformative Priors

- The motivation for noninformative priors
 - when prior information about the model is too vague or unreliable, it is usually impossible to justify the choice of prior distributions on a subjective basis.
 - “Objectivity” requirements which force us to provide prior distributions with as little subjective input as possible, in order to base inference on the sampling model alone.
- An intrinsic and acceptable notion of noninformative priors should satisfy invariance under reparametrization.

Noninformative Priors

- Noninformative priors are intended to have as little influence on the posterior as possible i.e. ‘letting the data speak for themselves’.
- Assume a distribution $p(x|\lambda)$ governed by a parameter λ , and a prior $p(\lambda) = \text{const}$ e.g. if λ is a discrete variable with K states, this simply amounts to setting the prior probability of each state to $1/K$.
- In the case of continuous λ there are two difficulties with this approach. If the domain of λ is unbounded, this prior distribution cannot be correctly normalized (**improper prior**).
- Improper priors can often be used provided the corresponding posterior distribution is proper.
 - For example, if we put a uniform prior distribution over the mean of a Gaussian, then the posterior distribution for the mean, once we have observed at least one data point, will be proper.



Noninformative Priors

- If we don't have strong beliefs about what θ should be, it is common to use an uninformative prior, and to let the data speak for itself.
- Consider as an example a Bernoulli parameter, $\theta \in [0,1]$.
- An uninformative prior would be the uniform distribution, $\text{Beta}(1,1)$. In this case, the posterior mean and MLE are:

$$\mathbb{E}[\theta | \mathcal{D}] = \frac{N_1 + 1}{N_1 + N_0 + 2}$$

$$\bar{\theta} = \frac{N_1}{N_1 + N_0}$$

- One could argue that the prior wasn't completely uninformative after all.

$$\text{Beta}(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \quad \mathbb{E}[x] = \frac{\alpha}{\alpha + \beta}$$



Noninformative Priors

- By the above argument, the most non-informative prior is

$$\lim_{c \rightarrow 0} \text{Beta}(c, c) = \text{Beta}(0, 0)$$

- This prior is a mixture of two equal point masses at 0 and 1.
- It is called *the Haldane prior*.
- Note that the Haldane prior is **an improper prior**, meaning it does not integrate to 1. However, as long as we see at least one head and at least one tail, the posterior will be proper.
- We will see shortly that the **right uninformative prior is:**

$$\text{Beta}(1/2, 1/2)$$

$$\text{Beta}(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \quad \mathbb{E}[x] = \frac{\alpha}{\alpha + \beta}$$



Noninformative Priors

- A second difficulty arises from the transformation behavior of a probability density under a nonlinear change of variables.
- If a function $h(\lambda)$ is constant, and we change variables to $\lambda = \eta^2$, then $h(\eta) = h(\eta^2)$ will also be constant. However, if we choose the density $p_\lambda(\lambda)$ to be constant, then the density of η will be given by

$$p_\eta(\eta) = p_\lambda(\lambda) \left| \frac{d\lambda}{d\eta} \right| = p_\lambda(\eta^2) 2\eta \propto \eta$$

and so the density over η will not be constant.

- This issue does not arise when we use maximum likelihood, because the likelihood function $p(x|\lambda)$ is a simple function of λ and so we are free to use any convenient parameterization.
- If, however, we are to choose a prior distribution that is constant, we must take care to use an appropriate representation for the parameters.



Translation Invariant Prior

- Translation Invariant: If the likelihood is of the form

$$p(x | \mu) = f(x - \mu)$$

then $f(\cdot)$ is translation invariant and μ is a location parameter.

- Note that if we shift x by a constant to give $\bar{x} = x + c$ then

$$p(\bar{x} | \bar{\mu}) = f(\bar{x} - \bar{\mu}), \text{ where } \bar{\mu} = \mu + c$$

- Thus the form of the density remains the same.
- We would like to find a prior that satisfies this translational invariance – a density independent of the origin.

Translation Invariant Prior

- We want a prior that assigns equal probability to the interval $A \leq \mu \leq B$ as to the interval $A-c \leq \mu \leq B-c$.

$$\int_A^B p(\mu) d\mu = \int_{A-c}^{B-c} p(t) dt = \int_{t=\mu-c}^B p(\mu-c) d\mu$$

- A translation invariance requirement is thus that the prior distribution should satisfy:

$$p(\mu) = p(\mu - c) \text{ for every } c \in \mathbb{R} \Rightarrow$$

$p(\mu) = \text{constant (improper prior)}$

- This flat prior is improper – but the resulting posterior is proper assuming

$$\int f(x - \theta) d\theta < \infty$$

Having seen $N \geq 1$ data points will satisfy this. One data point is enough to fix the location

- Example of a location parameter is the mean μ of a Gaussian. The noninformative prior is obtained from the conjugate prior

$$\mathcal{N}(\mu | \mu_0, \sigma_0^2) \text{ with } \sigma_0^2 \rightarrow \infty.$$



Scale Invariant Prior

- Scale Invariant: If the likelihood is of the form

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

then $f(\cdot)$ is scale invariant and σ is the scale parameter.

- Note that if we change the scale by a constant to give $\bar{x} = cx$ then

$$p(\bar{x} | \bar{\sigma}) = \frac{1}{\bar{\sigma}} f\left(\frac{\bar{x}}{\bar{\sigma}}\right), \text{ where } \bar{\sigma} = c\sigma$$

- Thus the form of the density remains the same.
- We would like to find a prior that satisfies this scale invariance – a density independent of the scaling used.

Scale Invariant Prior

- We want a prior that assigns equal probability to the interval $A \leq \sigma \leq B$ as to the interval $A/c \leq \sigma \leq B/c$.

$$\int_A^B p(\sigma) d\sigma = \int_{A/c}^{B/c} p(t) dt = \int_{t=\frac{\sigma}{c}}^{\frac{B}{c}} p\left(\frac{\sigma}{c}\right) \frac{1}{c} d\sigma$$

- A translation invariance requirement is thus that the prior distribution should satisfy:

$$p(\sigma) = p\left(\frac{\sigma}{c}\right) \frac{1}{c} \text{ for every } c \in \mathbb{R} \Rightarrow$$

$$p(\sigma) \propto \frac{1}{\sigma} \text{ (improper prior)} \Leftrightarrow p(\ln \sigma) = \text{const}$$

- We can approximate this with a $p(\sigma) = \text{Gamma}(\sigma | 0, 0)$. *This improper prior leads to a proper posterior if we observe $N \geq 2$ data* (we need at least 2 data points to estimate a variance)



Scale Invariant Prior

- Example of a scale parameter is the std σ of a Gaussian after we account for the location parameter:

$$\mathcal{N}(x|\mu, \sigma^2) \propto \frac{1}{\sigma} e^{-\left(\frac{\tilde{x}}{\sigma}\right)^2}, \quad \tilde{x} = x - \mu$$

- We can express this in terms of the precision $\lambda = 1/\sigma^2$ rather than σ itself.
- A distribution $p(\sigma) \propto 1/\sigma$ corresponds to a distribution over λ of the form $p(\lambda) \propto 1/\lambda$.
- We have seen that the conjugate prior for λ was $\text{Gamma}(\lambda | a_0, b_0)$.
The noninformative prior is obtained with $a_0 = b_0 = 0$. In this case, the posterior depends only from the data and not from the prior.

$$p(\lambda | X, \mu) = \prod_{n=1}^N f(x_n | \mu) \text{Gamma}(\lambda | a_0, b_0) \propto \lambda^{N/2+a_0-1} \exp\left(-b_0\lambda - \frac{1}{2}\lambda \sum_{n=1}^N (x_n - \mu)^2\right)$$

Jeffrey's Noninformative Priors

- Jeffrey's proposes a more intrinsic approach which avoids the need to take the invariance structure into account.
- Given a likelihood $f(x | \theta)$, Jeffrey's noninformative prior distributions are based on **Fisher information**, given by

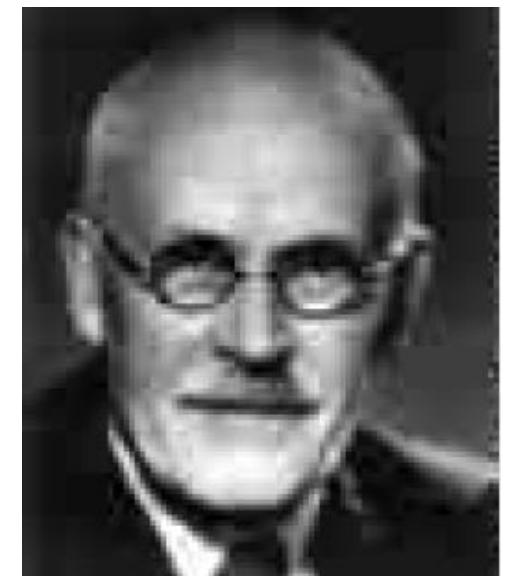
$$I(\theta) = \mathbb{E}_{X|\theta} \left(\frac{\partial \log f(X | \theta)}{\partial \theta} \frac{\partial \log f(X | \theta)^T}{\partial \theta} \right) = -\mathbb{E}_{X|\theta} \left(\frac{\partial^2 \log f(X | \theta)}{\partial \theta^2} \right)$$

the corresponding prior distribution is

$$\pi(\theta) \propto |I(\theta)|^{-1/2}$$

Determinant of I

Sir Harold Jeffreys
(1891–1989)



Jeffreys Noninformative Priors

□ Jeffreys Invariance Principle:

- Any rule for defining the prior distribution on θ should lead to an equivalent result when using a transformed parameterization
- Let $\phi = h(\theta)$ and h be an invertible function with inverse function $\theta = g(\phi)$, then

$$\pi(\phi) = \pi(g(\phi)) \left| \frac{dg(\phi)}{d\phi} \right| = \pi(\theta) \left| \frac{d\theta}{d\phi} \right|$$

- Jeffreys noninformative priors $\pi(\phi) \propto |I(\phi)|^{1/2}$ satisfy this invariant reparameterization requirement.

$$I(\phi) = -\mathbb{E}_{X|\phi} \left(\frac{\partial^2 \log f(X|\phi)}{\partial \phi^2} \right) = -\mathbb{E}_{X|\theta} \left(\frac{\partial^2 \log f(X|\phi)}{\partial \theta^2} \left| \frac{d\theta}{d\phi} \right|^2 \right) = I(\theta) \left| \frac{d\theta}{d\phi} \right|^2$$



Jeffrey's Noninformative Priors

- For example, consider normally distributed data with unknown mean
- Likelihood

$$x_i | \theta \sim \mathcal{N}(\theta, \sigma^2) \text{ (known } \sigma)$$

i.e.

$$f(x_{1:n} | \theta) \propto \exp\left(-\frac{n(\bar{x} - \theta)^2}{2\sigma^2}\right), \text{ where } \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

- Then:

$$\frac{\partial^2 \log f(x_{1:n} | \theta)}{\partial \theta^2} = -\frac{n}{\sigma^2} \Rightarrow \pi(\theta) \propto 1$$

Jeffrey's Noninformative Priors

- Consider normally distributed data with unknown variance

➤ Likelihood $X_i | \theta \sim \mathcal{N}(\mu, \theta)$ (*known* μ)

i.e.

$$f(x_{1:n}|\theta) \propto \theta^{-n/2} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2 + n(\bar{x} - \mu)^2}{2\theta}\right)$$

Then: $\frac{\partial^2 \log f(x_{1:n}|\theta)}{\partial \theta^2} = \frac{n}{2\theta^2} - \frac{\sum_{i=1}^n (x_i - \mu)^2 + n(\bar{x} - \mu)^2}{\theta^3} \Rightarrow$

$$\begin{aligned} I(\theta) &= -\mathbb{E}_{X|\theta} \left(\frac{n}{2\theta^2} - \frac{\sum_{i=1}^n (x_i - \mu)^2 + n(\bar{x} - \mu)^2}{\theta^3} \right) = -\frac{n}{2\theta^2} + \mathbb{E}_{X|\theta} \left(\frac{\sum_{i=1}^n (x_i - \mu)^2}{\theta^3} \right) \\ &= -\frac{n}{2\theta^2} + \frac{n}{\theta^2} = \frac{n}{2\theta^2} \end{aligned}$$

➤ Jeffrey's prior $\pi(\theta = \sigma^2) \propto \frac{1}{\theta} = \frac{1}{\sigma^2}$ (favors small variance)

➤ Note that $\pi(\phi = \log \theta) \propto \frac{1}{\theta} \left| \frac{d\theta}{d\phi} \right| = \frac{1}{\theta} \theta = 1$



Jeffrey's Noninformative Priors

- Consider data following a binomial distribution (mean $n\theta$)
 - Likelihood

$$f(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

Then:

$$\frac{\partial^2 \log f(x|\theta)}{\partial \theta^2} = -\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2} \Rightarrow I(\theta) = -\mathbb{E}_{X|\theta} \left(-\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2} \right) = \frac{n\theta}{\theta^2} + \frac{n-n\theta}{(1-\theta)^2} = \frac{n}{\theta(1-\theta)}$$

- The Jeffrey's prior is:

$$\pi(\theta) \propto [\theta(1-\theta)]^{-1/2} = \text{Beta}\left(\theta; \frac{1}{2}, \frac{1}{2}\right)$$

Beta Distribution: $p(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$

$$\begin{aligned} E(\theta) &= \frac{\alpha}{\alpha+\beta} \\ \text{var}(\theta) &= \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} \\ \text{mode}(\theta) &= \frac{\alpha-1}{\alpha+\beta-2} \end{aligned}$$

- For a multinoulli random variable with K states, one can show that the Jeffreys' prior is:

$$\pi(\theta) = \text{Dir}\left(\frac{1}{2}, \dots, \frac{1}{2}\right)$$

- Note that this is not any of the expected answers:

$$\pi(\theta) = \text{Dir}\left(\frac{1}{K}, \dots, \frac{1}{K}\right) \text{ or } \pi(\theta) = \text{Dir}(1, \dots, 1)$$

Pros and Cons of Jeffrey's Priors

- It can lead to incoherencies; i.e. the Jeffrey's prior for Gaussian data and $\theta = (\mu, \sigma)$ unknown is $\pi(\theta) \propto \sigma^{-2}$. Indeed using: $\ln f(x|\theta) = \ln \frac{1}{(2\pi)^{1/2}} - \ln \sigma - \frac{1}{2\sigma^2}(x-\mu)^2$

$$I(\theta) = \mathbb{E}_{X|\theta} \begin{bmatrix} \frac{1}{\sigma^2} & \frac{2(x-\mu)}{\sigma^3} \\ \frac{2(x-\mu)}{\sigma^3} & \frac{3(\mu-x)^2}{\sigma^4} - \frac{1}{\sigma^2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{bmatrix} \Rightarrow \pi(\theta) \propto \frac{1}{\sigma^2}$$

- However if these parameters are assumed a priori independent (using the results [derived earlier](#)) then $\pi(\theta) \propto \sigma^{-1}$.
- Automated procedure that however cannot incorporate any “physical” information.
- It does NOT satisfy the likelihood principle. The Fisher information can differ for two experiments providing proportional likelihoods. An example is given next for the [Binomial](#) and [Negative Binomial](#) distributions.

C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter](#) 3 (full text available)



Jeffrey's Priors and the Likelihood Principle

- The binomial and negative binomial models could lead to the same likelihood.
- If $x \sim \mathcal{B}(n, \theta)$, the non-informative prior is $\mathcal{Be}\left(\theta; \frac{1}{2}, \frac{1}{2}\right)$ ([recall proof](#)).
- Consider $n \sim \mathcal{NB}(x, \theta)$ ([negative Binomial](#)) (n trials until we have x successes with probability θ),
- Using the formula for the mean (number of failures) given below, the mean for the number of trials is then:

$$\mathbb{E}(n) = \mathbb{E}(n - x) + x = \frac{x(1-\theta)}{\theta} + x = \frac{x}{\theta}$$

Note that the interpretation of symbols for the \mathcal{NB} tables given below: θ =number of failures until you observe α successes (so in the notation of our problem above, $x \leftarrow \alpha$ and $n \leftarrow \alpha + \theta$, $\theta \leftarrow \beta/(\beta+1)$ =probability of success, $1-\theta \leftarrow 1/(\beta+1)$ =probability of failure). Note that α/β (mean of # of failures) in our problem becomes $x(1-\theta)/\theta$.

Negative binomial	$\theta \sim \text{Neg-bin}(\alpha, \beta)$ $p(\theta) = \text{Neg-bin}(\theta \alpha, \beta)$	shape $\alpha > 0$ inverse scale $\beta > 0$
	$p(\theta) = \binom{\theta+\alpha-1}{\alpha-1} \left(\frac{\beta}{\beta+1}\right)^\alpha \left(\frac{1}{\beta+1}\right)^\theta$ $\theta = 0, 1, 2, \dots$	$\mathbb{E}(\theta) = \frac{\alpha}{\beta}$ $\text{var}(\theta) = \frac{\alpha}{\beta^2}(\beta + 1)$

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Jeffrey's Priors and the Likelihood Principle

- Thus for $n \sim \mathcal{NB}(x, \theta)$ (negative Binomial), the likelihood is of the form

$$f(n | \theta) = \binom{x + (n-x)-1}{x-1} \theta^x (1-\theta)^{n-x} \Rightarrow$$

- The Fisher information is then (using the earlier result, $\mathbb{E}(n) = \frac{x}{\theta}$):

$$I(\theta) = \mathbb{E}_{n|\theta} \left(\frac{x}{\theta^2} + \frac{n-x}{(1-\theta)^2} \right) = \frac{x}{\theta^2} + \frac{\left(\mathbb{E}_{n|\theta}(n) - x \right)}{(1-\theta)^2} = \frac{x}{\theta^2(1-\theta)}$$

- The corresponding prior is then given as (improper and different from the one obtained using [the Binomial likelihood](#)):

$$\pi(\theta) \propto \theta^{-1} (1-\theta)^{-1/2} \neq \mathcal{Be}\left(\theta; \frac{1}{2}, \frac{1}{2}\right)$$

Negative binomial	$\theta \sim \text{Neg-bin}(\alpha, \beta)$ $p(\theta) = \text{Neg-bin}(\theta \alpha, \beta)$	shape $\alpha > 0$ inverse scale $\beta > 0$
	$p(\theta) = \binom{\theta+\alpha-1}{\alpha-1} \left(\frac{\beta}{\beta+1}\right)^\alpha \left(\frac{1}{\beta+1}\right)^\theta$ $\theta = 0, 1, 2, \dots$	$E(\theta) = \frac{\alpha}{\beta}$ $\text{var}(\theta) = \frac{\alpha}{\beta^2} (\beta + 1)$

C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter](#) 3 (full text available for Cornell students)



Improper Prior Distributions

- In most cases, the prior distribution is determined on a subjective or theoretical basis which provides a prior π such that

$$\int_{\Theta} \pi(\theta) d\theta = +\infty$$

In such cases, the prior distribution is said to be improper (or generalized).

- The fact that the prior distribution is improper weakens the symmetry between the observations and the parameters, but as long as the posterior distribution is defined, Bayesian methods apply as well.
- The usual convention is to take the posterior distribution $\pi(\theta|y)$ associated with an improper π as given by Bayes' formula

$$\pi(\theta | y) = \frac{\pi(y | \theta)\pi(\theta)}{\int_{\Theta} \pi(y | \theta)\pi(\theta)d\theta}$$

when the marginal distribution $\int_{\Theta} \pi(y | \theta)\pi(\theta)d\theta$ is well defined.



Improper Priors

- Extension of the posterior distribution $\pi(\theta|x)$ associated with an improper prior π given by Bayes's formula

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int_{\Theta} f(x|\theta)\pi(\theta)d\theta}$$

when

$$\int_{\Theta} f(x|\theta)\pi(\theta)d\theta < \infty$$

- Example (Normal + improper)

If $x \sim \mathcal{N}(\theta, 1)$ and $\pi(\theta) = w$, constant, the pseudo-marginal distribution is

$$m(x) = w \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-\theta)^2}{2}\right\} d\theta$$

and the posterior distribution of θ is

$$\pi(\theta|x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-\theta)^2}{2}\right\}$$

i.e., corresponds to $\mathcal{N}(x, 1)$ [independently of w]



Improper Prior Distribution

- Extension from a prior distribution to a prior σ -finite measure π such that

$$\int_{\Theta} \pi(\theta) d\theta = +\infty$$

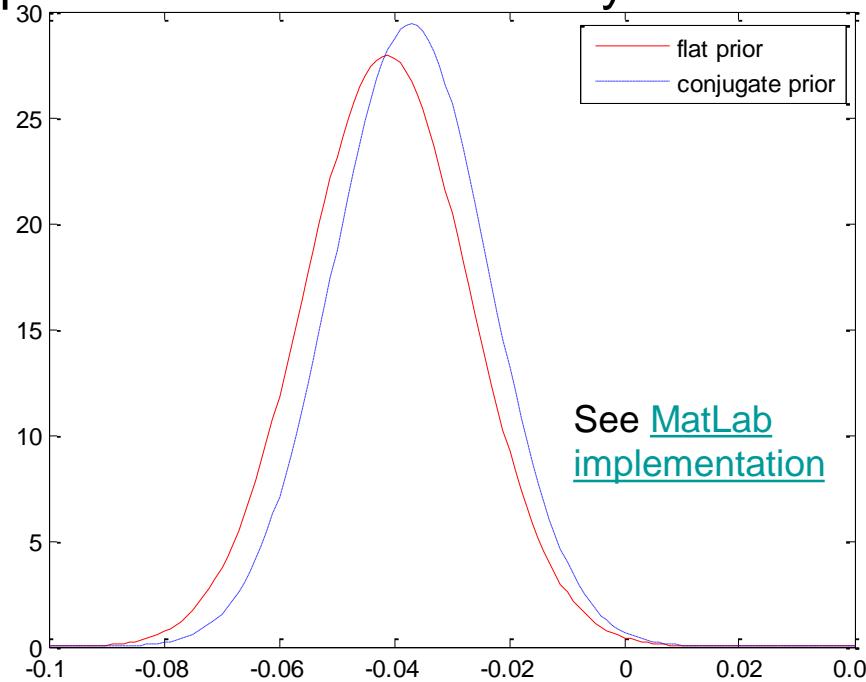
Formal extension: π cannot be interpreted as a probability any longer

- Justifications
 - Often only way to derive a prior in non-informative/automatic settings
 - Performances of associated estimators usually good
 - Often occur as limits of proper distributions
 - More robust answer against possible misspecifications of the prior
 - Improper priors preferable to vague proper priors such as a normal distribution $\mathcal{N}(0, 100^2)$



Lack of Robustness of the Normal Prior

- Comparison of two posterior distributions corresponding to the flat prior (plain) and a conjugate prior (dotted) $\mathcal{N}(0, 0.1\bar{\sigma}^2)$ (where the variance $\bar{\sigma}^2$ refers here to the empirical variance of the sample). We use the data [normaldata](#). This shows the lack of robustness of the normal prior.
- When the [hyperparameters in the prior](#) vary, both the range and location of the posterior are not limited by the data.



J.-M. Marin & C. P. Robert, [The Bayesian Core](#), Springer, 2nd edition, [chapter 2](#) (full text available)



Robust Priors: Priors with Heavy Tails

- In many cases, we are not very confident in our prior, so we want to make sure it does not have an undue influence on the result.
- This can be done by using *robust priors, which typically have heavy tails, which avoids forcing things to be too close to the prior mean.*
- As an example, consider $x \sim \mathcal{N}(\theta, 1)$. We observe that $x = 5$ and we want to estimate θ . The MLE is $\hat{\theta} = 5$, which seems reasonable. The posterior mean under a uniform prior is also $\mathbb{E}[\theta | x = 5] = 5$
- Suppose we know that the prior median is 0, and the prior quantiles are at -1 and 1, so $p(\theta \leq -1) = p(-1 < \theta \leq 0) = p(0 < \theta \leq 1) = p(1 < \theta) = 0.25$. Let us also assume the prior is smooth and unimodal.
- Using the prior $\mathcal{N}(\theta | 0, 2.19^2)$ satisfies these prior constraints. But in this case the posterior mean is 3.43, which is not very satisfactory.
- Use Cauchy prior $\mathcal{T}(\theta | 0, 1, 1)$. This also satisfies the prior constraints of our example. But this time we find that the posterior mean is about 4.6, which seems much more reasonable.

[robustPriorDemo](#)
from [Kevin Murphys' PMTK](#)



Prior Distributions

- There is no such a thing as ‘the exact prior’. In most applications, however, there is “true” prior.
- Although conjugate priors are limited, they remain the most widely used class of priors for convenience and simple interpretability.
- There is a whole literature on the subject: reference & objective priors.
- Empirical Bayes: the prior is constructed from the data.
- In all cases, you should play around with different priors and do a sensitivity analysis.

Hierarchical Bayesian Models

- It often helps to decompose prior knowledge into several levels particularly when the available data is hierarchical.
- The **hierarchical Bayes method** is a powerful tool for expressing rich statistical models that more fully reflect a given problem than a simpler model could.
- Often the prior on θ depends in turn on other parameters ϕ that are not mentioned in the likelihood. So, the prior $\pi(\theta)$ must be replaced by a prior $\pi(\theta|\phi)$, and a prior $\pi(\phi)$ on the newly introduced parameters ϕ is required, resulting in a posterior probability $\pi(\theta,\phi|x)$.

$$\pi(\theta, \phi | x) \sim \underbrace{\pi(x | \theta, \phi)}_{\pi(x|\theta)} \pi(\theta, \phi) \sim \pi(x | \theta) \underbrace{\pi(\theta | \phi) \pi(\phi)}_{\pi(\theta, \phi)}$$

- This is the simplest example of a *hierarchical Bayes model*.
- The process may be repeated, e.g., ϕ may depend on parameters ψ , which will require their own prior. **Eventually the process must terminate**, with priors that do not depend on any other parameters.



Hierarchical Bayesian Models

- Consider m -level hierarchical Bayesian model

$$\pi(\theta) = \int_{\Theta_1 \times \Theta_2 \times \dots \times \Theta_m} \pi(\theta | \theta_1) \pi(\theta_1 | \theta_2) \dots \pi(\theta_{m-1} | \theta_m) \pi(\theta_m) d\theta_1 \dots d\theta_m$$

- Two level hierarchical modeling gives:

✓ Full posterior: $\pi(\theta, \theta_1 | x) \sim \underbrace{\pi(x | \theta) \pi(\theta | \theta_1)}_{\pi(\theta | \theta_1, x)} \pi(\theta_1)$

✓ Conditional posterior: $\pi(\theta | \theta_1, x) \sim \pi(x | \theta) \pi(\theta | \theta_1)$

✓ Marginal Posterior: $\pi(\theta | x) = \int \pi(\theta, \theta_1 | x) d\theta_1$



Hierarchical Bayes

- A key requirement for computing the posterior $p(\theta|\mathcal{D})$ is the specification of a prior $p(\theta|\eta)$, where η are the hyper-parameters.
- What if we don't know how to set η ?
- In some cases, we can use uninformative priors as discussed earlier.
- A more Bayesian approach is to put a prior on our priors! *In terms of graphical models (showing explicitly dependence relations)*, we can represent the situation as follows:

$$\eta \rightarrow \theta \rightarrow \mathcal{D}$$

- This is an example of a hierarchical Bayesian model, also called a **multi-level model**, since there are multiple levels of unknown quantities.



Hierarchical Bayes: Modeling Cancer Rates

- Consider the problem of predicting cancer rates in various cities.
- We measure the people in various cities, N_i , and the people who died of cancer in these cities, x_i . We assume $x_i \sim \text{Bin}(N_i, \theta_i)$ and we estimate the cancer rates θ_i .
- We can estimate them all separately, but this will suffer from the sparse data problem (underestimation of the rate of cancer due to small N_i).
- We can assume all the θ_i are the same (*parameter tying*). But the assumption that all the cities have the same rate is a rather strong one.
- As a compromise we assume that the θ_i are similar, but that there may be city-specific variations. This can be modeled by assuming $\theta_i \sim \text{Beta}(a, b)$. The full joint distribution can be written as

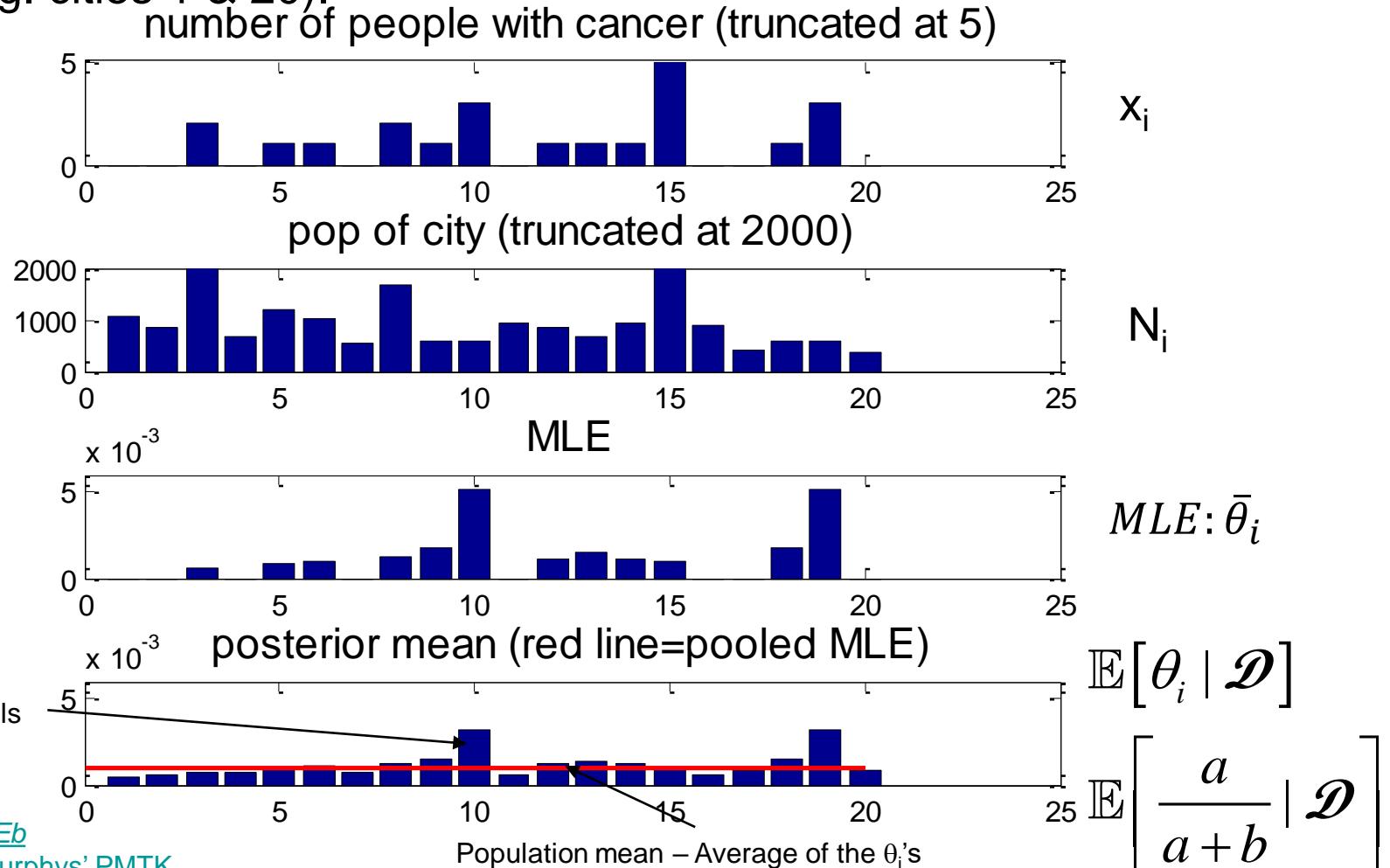
$$p(\mathcal{D}, \theta, \eta) = p(\eta) \prod_{i=1}^N \text{Bin}(x_i | N_i, \theta_i) \text{Beta}(\theta_i | \eta), \eta = (a, b)$$

- By treating η as an unknown (hidden variable), we *allow the data-poor cities to borrow statistical strength from data-rich ones*.



Hierarchical Bayes: Modeling Cancer Rates

- Compute $p(\eta, \theta | \mathcal{D})$, then the marginal $p(\theta | \mathcal{D})$. The posterior mean is shrunk towards the pooled estimate more strongly for cities with small N_i (e.g. cities 1 & 20).



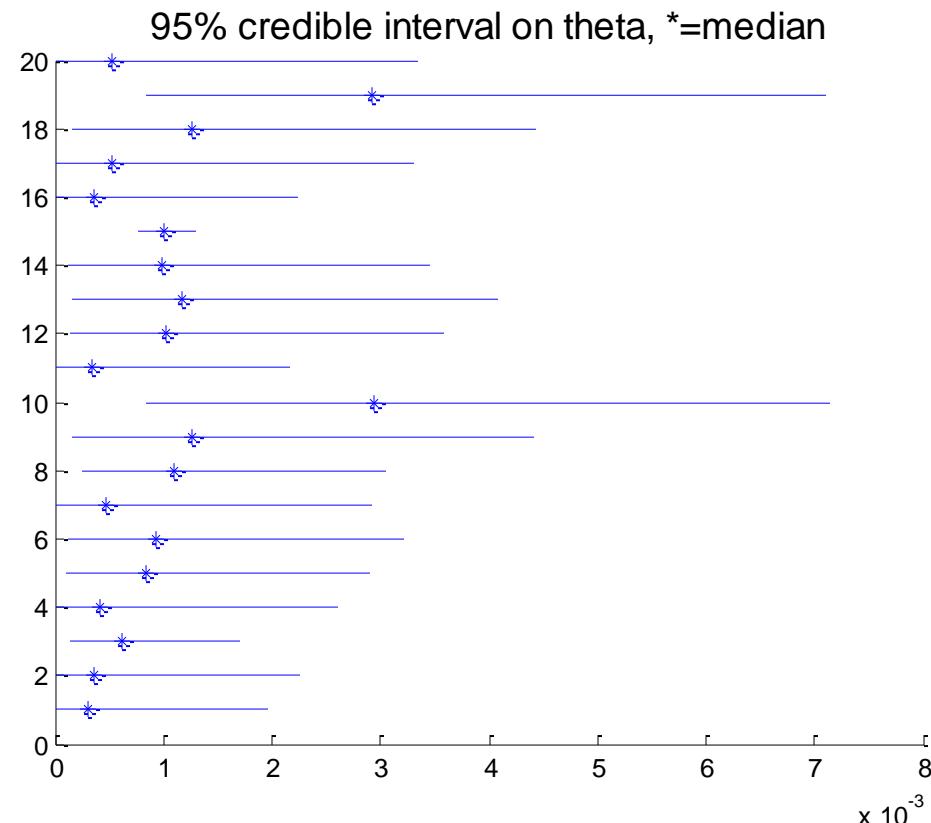
[cancerRatesEb](#)
from [Kevin Murphy's PMTK](#)



Hierarchical Bayes: Modeling Cancer Rates

- 95% posterior credible intervals for θ_i .
- City 15, which has a very large population, has small posterior uncertainty. It has the largest impact on the posterior of η which in turn impacts the estimate of the cancer rates for other cities.
- Cities 10 and 19, which have the highest MLE, also have the highest posterior uncertainty, reflecting the fact that such a high estimate is in conflict with the prior (which is estimated from all the other cities).

[cancerRatesEb](#)
[from Kevin Murphys' PMTK](#)



Empirical Bayes - Evidence Approximation

- In hierarchical Bayesian models, we need to compute the posterior on multiple levels of latent variables. For example, in a two-level model, we need to compute

$$p(\boldsymbol{\eta}, \boldsymbol{\theta} | \mathcal{D}) \propto p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\eta}) p(\boldsymbol{\eta})$$

- In some cases, we can analytically marginalize out $\boldsymbol{\theta}$; this leaves us with the simpler problem of just computing $p(\boldsymbol{\eta} | \mathcal{D})$.
- As a computational shortcut, we can *approximate the posterior on the hyper-parameters with a point-estimate*,

$$p(\boldsymbol{\eta} | \mathcal{D}) \approx \delta_{\bar{\boldsymbol{\eta}}}(\boldsymbol{\eta}), \bar{\boldsymbol{\eta}} = \operatorname{argmax} p(\boldsymbol{\eta} | \mathcal{D}) = \operatorname{argmax} \left[\int p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\eta}) d\boldsymbol{\theta} \right]$$

- Since $\boldsymbol{\eta}$ is typically much smaller than $\boldsymbol{\theta}$ in dimensionality, it is less prone to overfitting, so we can safely use a uniform prior on $\boldsymbol{\eta}$.
- The quantity inside the brackets is *the marginal or integrated likelihood, often called the evidence*. This overall approach is called *empirical Bayes (EB) or type-II maximum likelihood*. In machine learning, it is sometimes called the *evidence procedure*.



Empirical Bayes

- Empirical Bayes violates the principle that the prior should be chosen independently of the data.
- We can just view it as a cheap approximation to inference in a hierarchical Bayesian model, just as we viewed MAP estimation as an approximation to inference in the one level model $\theta \rightarrow \mathcal{D}$.
- We can construct a hierarchy in which the more integrals one performs, the “more Bayesian” one becomes:

Method	Definition
Maximum likelihood	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)$
MAP estimation	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)p(\theta \eta)$
ML-II (Empirical Bayes)	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)$
MAP-II	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)p(\eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)p(\eta)$
Full Bayes	$p(\theta, \eta \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta \eta)p(\eta)$

Empirical Bayes

- Let us return to [the cancer rates model](#). We can analytically integrate out θ_i , and write down the marginal likelihood directly, as follows:

$$p(\mathcal{D} | a, b) = \prod_i \int \mathcal{Bin}(x_i | N_i, \theta_i) \mathcal{Beta}(\theta_i | a, b) d\theta_i = \prod_i \binom{N_i}{x_i} \frac{B(a + x_i, b + N_i - x_i)}{B(a, b)}$$

- Various ways of maximizing this wrt a and b are discussed in Minka.
- Having estimated a and b, we can plug in the hyper-parameters to compute the posterior $p(\theta_i | \mathcal{D}, \bar{a}, \bar{b})$ in the usual way, using conjugate analysis.
- It can be shown that *the posterior mean of each θ_i is a weighted average of its local MLE and the prior means, which depends on $\eta = (a, b)$.*
- Since η is estimated using all the data, *each θ_i is influenced by all data.*

Minka, T. (2000e). [Estimating a Dirichlet distribution](#), Technical Report.



Empirical Bayes: Gaussian-Gaussian Model

- We now consider an example where the data is real-valued. We use a *Gaussian likelihood and a Gaussian prior*.
- Suppose we have data from multiple related groups, e.g. x_{ij} is the test score for **student i** in **school j**, $j = 1:D$, $i = 1:N_j$. *We want to estimate the mean score for each school, θ_j .*
- Since N_j may be small for some schools, we regularize the problem by using a hierarchical Bayesian model, where θ_j comes from a common prior, $\mathcal{N}(\mu, \tau^2)$.
- The joint distribution has the following form:

$$p(\boldsymbol{\theta}, \mathcal{D} | \boldsymbol{\eta}, \sigma^2) = \prod_{j=1}^D \mathcal{N}(\theta_j | \mu, \tau^2) \prod_{i=1}^{N_j} \mathcal{N}(x_{ij} | \theta_j, \sigma^2), \boldsymbol{\eta} = (\mu, \tau)$$

- We assume for simplicity that σ^2 is known.

Empirical Bayes: Gaussian-Gaussian Model

- We rewrite the joint distribution exploiting the fact that *N_j Gaussian measurements with values x_{ij} and variance σ² are equivalent to one measurement* $\bar{x}_j = \frac{1}{N_j} \sum_{i=1:N_j} x_{ij}$ *with variance* $\sigma_j^2 = \sigma^2 / N_j$

$$\text{measurement } \bar{x}_j = \frac{1}{N_j} \sum_{i=1:N_j} x_{ij} \quad \text{with variance } \sigma_j^2 = \sigma^2 / N_j$$

- This yields

$$p(\theta, \mathcal{D} | \hat{\eta}, \sigma^2) = \prod_{j=1}^D \mathcal{N}(\theta_j | \hat{\mu}, \hat{\tau}^2) \mathcal{N}(\bar{x}_j | \theta_j, \sigma_j^2)$$

where: $\hat{\mu} = \frac{1}{D} \sum_{j=1}^D \bar{x}_j \equiv \bar{x}$, $\sigma^2 + \hat{\tau}^2 = \frac{1}{D} \sum_{j=1}^D (\bar{x}_j - \bar{x})^2$

- From this, closing the square on θ_j , it follows that the posteriors are:

$$p(\theta_j | D, \hat{\mu}, \hat{\tau}^2) = \mathcal{N}(\theta_j | \hat{B}_j \hat{\mu} + (1 - \hat{B}_j) \bar{x}_j, (1 - \hat{B}_j) \sigma_j^2),$$

$$\hat{\mu} = \frac{1}{D} \sum_{j=1}^D \bar{x}_j \equiv \bar{x}, \sigma^2 + \hat{\tau}^2 = \frac{1}{D} \sum_{j=1}^D (\bar{x}_j - \bar{x})^2, \hat{B}_j = \frac{\sigma_j^2}{\sigma_j^2 + \hat{\tau}^2}$$

Empirical Bayes: Gaussian-Gaussian Model

- Note that for constant σ_j^2

$$\int \mathcal{N}(\theta_j | \mu, \tau^2) \mathcal{N}(\bar{x}_j | \theta_j, \sigma_j^2) d\theta_j = \mathcal{N}(\bar{x}_j | \mu, \sigma^2 + \tau^2) \Rightarrow$$

$$p(\mathcal{D} | \mu, \tau^2, \sigma^2) = \prod_{i=1}^D \mathcal{N}(\bar{x}_j | \mu, \sigma^2 + \tau^2)$$

- We can now derive the previously shown estimates using MLE:

$$\hat{\mu} = \frac{1}{D} \sum_{j=1}^D \bar{x}_j \equiv \bar{x}, \quad \sigma^2 + \hat{\tau}^2 = \frac{1}{D} \sum_{j=1}^D (\bar{x}_j - \bar{x})^2$$

- For non-constant σ_j^2 , you need to use Expectation-Maximization to derive the EB estimate.

James Stein Estimator

$$p(\theta_j | D, \hat{\mu}, \tau^2) = N(\theta_j | \hat{B}_j \hat{\mu} + (1 - \hat{B}_j) \bar{x}_j, (1 - \hat{B}_j) \sigma_j^2),$$

$$\hat{\mu} = \frac{1}{D} \sum_{j=1}^D \bar{x}_j \equiv \bar{x}, \quad \sigma^2 + \hat{\tau}^2 = \frac{1}{D} \sum_{j=1}^D (\bar{x}_j - \bar{x})^2, \quad \hat{B}_j = \frac{\sigma_j^2}{\sigma_j^2 + \hat{\tau}^2}$$

- The quantity $0 \leq \hat{B}_j \leq 1$ controls the degree of *shrinkage towards the overall mean*, $\hat{\mu}$. If the data is reliable for group j , then σ_j^2 will be small relative to $\hat{\tau}^2$; hence \hat{B}_j will be small, and we will put more weight on \bar{x}_j when we estimate θ_j . However, groups with small N_j will get regularized (shrunk towards the overall mean $\hat{\mu}$) more heavily.
- For σ_j constant across j , the posterior mean becomes (**James Stein estimator**):

$$\hat{\theta}_j = \hat{B} \bar{x} + (1 - \hat{B}) \bar{x}_j = \bar{x} + (1 - \hat{B})(\bar{x}_j - \bar{x}), \quad \hat{B} = \frac{\sigma^2}{\sigma^2 + \hat{\tau}^2}$$

Predicting Baseball Scores

- This is an example of shrinkage applied to baseball batting averages.
- We observe the number of hits for $D = 18$ players during the first $T = 45$ games. Let the number of hits b_i and assume $b_j \sim \mathcal{B}(T, \theta_j)$, where θ_j is the “true” batting average for player j . The goal is to estimate the θ_j .
- The MLE is $\hat{\theta}_j = x_j$, $x_j = b_j/T$ being the empirical batting average. One can use an Empirical Bayes approach to do better.
- To apply the Gaussian shrinkage approach described above, we require that the likelihood be Gaussian, $x_j \sim \mathcal{N}(\theta_j, \sigma^2)$ for known σ^2 . (We drop the i subscript since we assume $N_j = 1$, since x_j already represents the average for player j .)
- However, in this example we have a binomial likelihood. While this has the right mean, $\mathbb{E}[x_j] = \theta_j$, the variance is not constant:

$$\text{var}[x_j] = \frac{1}{T^2} \text{var}[b_j] = \frac{T\theta_j(1-\theta_j)}{T^2}$$

- Efron, B. and C. Morris (1975). Data analysis using stein's estimator and its generalizations. J. of the Am. Stat. Assoc. 70(350), 311–319.



Predicting Baseball Scores

- So we apply a variance stabilizing transform to x_j to better match the Gaussian assumption.

$$y_j = f(x_j) = \sqrt{T} \arcsin(2x_j - 1)$$

- Now we have approximately $y_j \sim \mathcal{N}(f(\theta_j), 1) = \mathcal{N}(\mu_j, 1)$. We use Gaussian shrinkage to estimate the μ_j using

$$\hat{\mu}_j = \hat{B}\bar{y} + (1 - \hat{B})\bar{y}_j = \bar{y} + (1 - \hat{B})(\bar{y}_j - \bar{y})$$

with $\sigma^2 = 1$, and we then transform back to get

$$\hat{\theta}_j = 0.5 \sin\left(\frac{\hat{\mu}_j}{\sqrt{T}} + 1\right)$$

- The results are shown next.

Consider a transform $Y = f(X)$ where $\mathbb{E}[X] = \mu$, $\text{var}[X] = \sigma^2$ s.t.

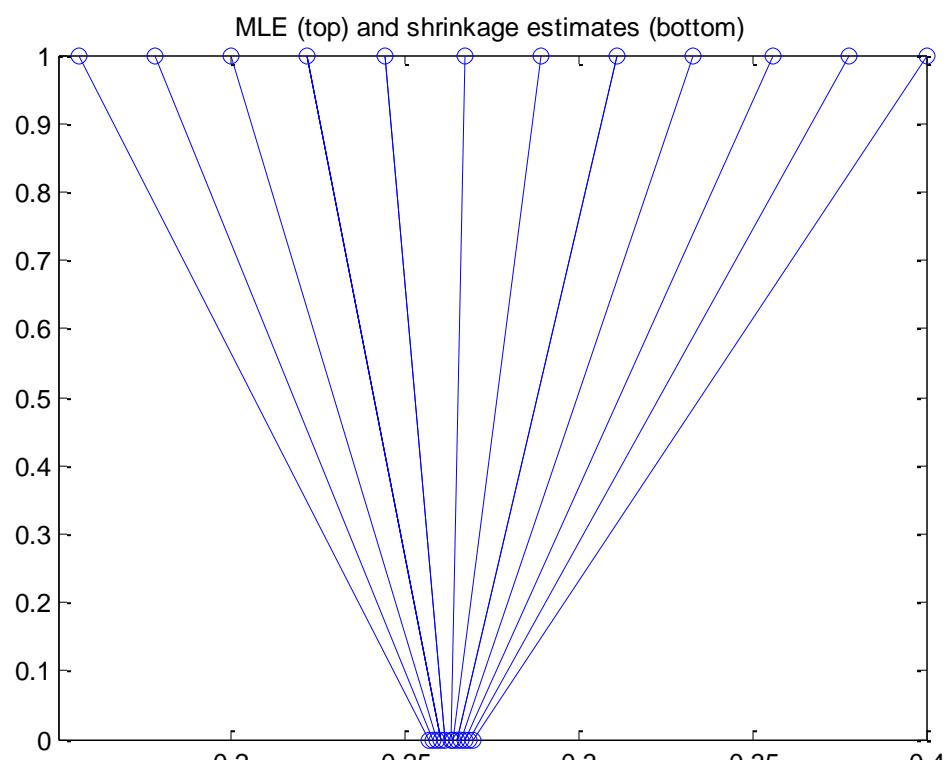
$Y = f(X) \approx f(\mu) + f'(\mu)(X - \mu)$ with $\text{var}[Y] = f'(\mu)^2 \sigma^2(\mu)$

If $f'(\mu)^2 \sigma^2(\mu)$ is independent of μ , we call $f(X)$ a variance stabilizing transform

$$\text{Here : } f'(\mu)^2 \sigma^2(\mu) = \frac{4T}{1 - (2x_j - 1)^2} \Bigg|_{\mu = \mathbb{E}[x_j] = \theta_j} \frac{T\theta_j(1 - \theta_j)}{T^2} = 1$$



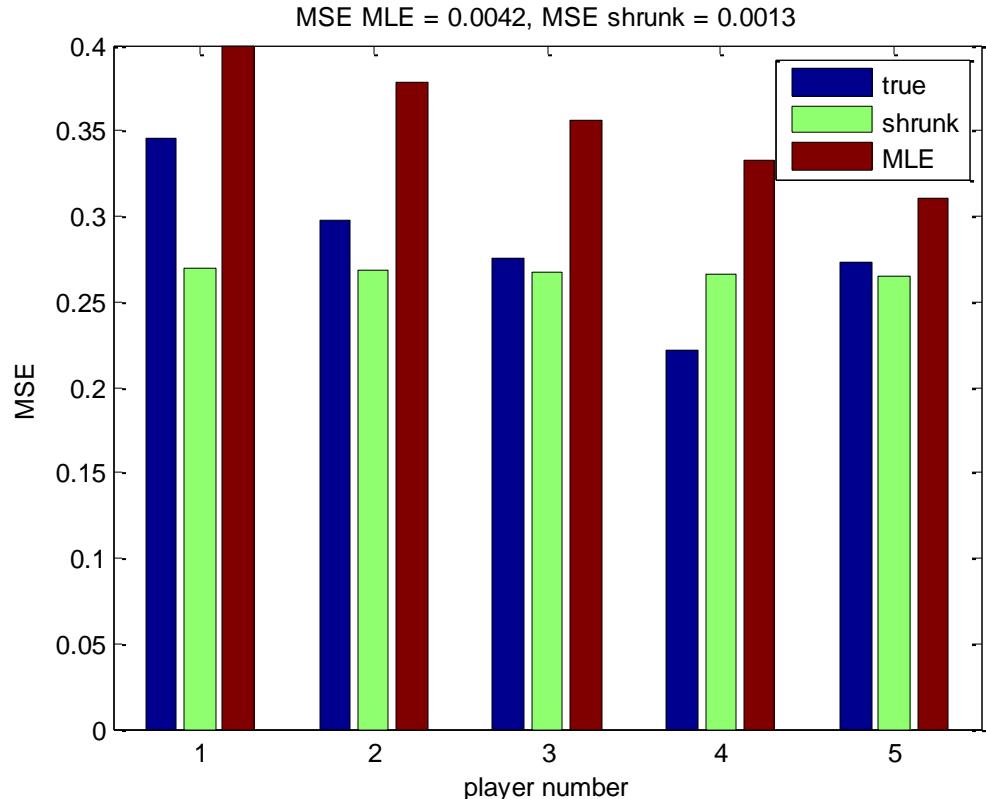
Predicting Baseball Scores



- We plot the MLE $\hat{\theta}_j$.
- All the estimates have shrunk towards the global mean, 0.265.

[ShrinkageDemoBaseBall](#)
from Kevin Murphys' PMTK

Predicting Baseball Scores



[ShrinkageDemoBaseBall](#)
from Kevin Murphys' PMTK

$$MSE = \frac{1}{N} \sum_{j=1}^D (\theta_j - \bar{\theta}_j)^2$$

- We plot the true value θ_j , the MLE $\hat{\theta}_j$ and the posterior mean $\bar{\theta}_j$
- The “true” values of θ_j are estimated from a large number of independent games.) On average, the shrunken estimate is much closer to the true parameters than the MLE is.
- The mean squared error is over three times smaller using the $\bar{\theta}_j$ shrinkage estimates than using the MLEs $\hat{\theta}_j$