

# Bayesian Essentials



**Feng Li**

feng.li@cufe.edu.cn

**School of Statistics and Mathematics  
Central University of Finance and Economics**

**NOTE: Most of the contents are from Bayesian Statistics course  
taught by Professor Mattias Villani**

**<http://www.mattiasvillani.com/teaching/bayesian-statistics/>**

- Likelihood
- Bayesian inference
- The Bernoulli model
- The Normal model

- Bernoulli trials:

$$x_1, \dots, x_n | \theta \stackrel{iid}{\sim} \text{Bern}(\theta).$$

- Likelihood:

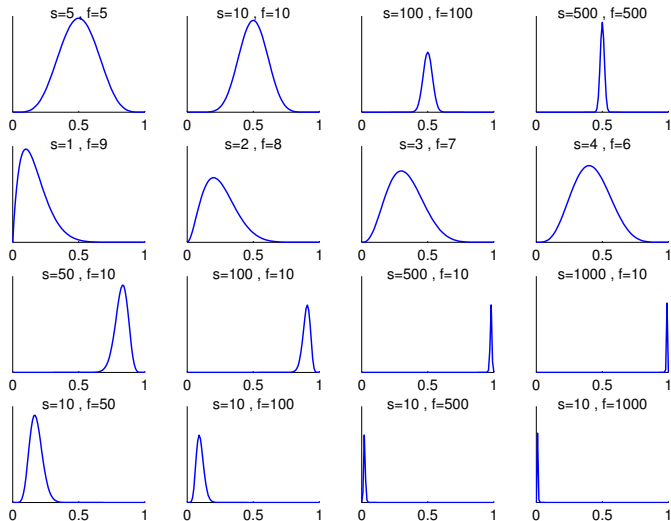
$$\begin{aligned} p(x_1, \dots, x_n | \theta) &= p(x_1 | \theta) \cdots p(x_n | \theta) \\ &= \theta^s (1 - \theta)^f, \end{aligned}$$

where  $s = \sum_{i=1}^n x_i$  is the number of successes in the Bernoulli trials and  $f = n - s$  is the number of failures.

- Given the data  $x_1, \dots, x_n$ , we may plot  $p(x_1, \dots, x_n | \theta)$  as a function of  $\theta$ .

# The likelihood function from Bernoulli trials

Likelihood function of the Bernoulli model for different data



# Uncertainty and subjective probability

- Will the likelihood give us an idea of which values of  $\theta$  that should be regarded as probable (in some sense)? Kind of, but ... No!
- In order to say that one value of  $\theta$  is more probable than another we clearly must think of  $\theta$  as random. But  $\theta$  may be something that we know is non-random, e.g. a fixed natural constant.
- **Bayesian: doesn't matter if  $\theta$  is fixed or random.** What matters is whether or not You know the value of  $\theta$ . If  $\theta$  is uncertainty to You, then You can assign a probability distribution to  $\theta$  which reflects Your knowledge about  $\theta$ . **Subjective probability.**

## Learning from data - Bayes' theorem

- Given that you have formulated a distribution for  $\theta$ ,  $p(\theta)$ , how can we learn from data? That is, how do we make the transition from  $p(\theta) \rightarrow p(\theta|Data)$ ? Bayes' theorem is the key.
- One form of Bayes' theorem reads ( $A$  and  $B$  are events)

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}.$$

So that Bayes' theorem 'reverses the conditioning', i.e. takes us from  $p(B|A)$  to  $p(A|B)$ .

- Let  $A = \theta$  and  $B = Data$

$$p(\theta|Data) = \frac{p(Data|\theta)p(\theta)}{p(Data)}.$$

- Interpreting the likelihood function as a probability density for  $\theta$  is just as wrong as ignoring the factor  $p(A)/p(B)$  in Bayes' theorem.

# Generalized Bayes' theorem

- From your basic statistics textbook:

$$p(A_i|B) = \frac{p(B|A_i)p(A_i)}{p(B)} = \frac{p(B|A_i)p(A_i)}{\sum_{i=1}^k p(B|A_i)p(A_i)}.$$

- Let  $\theta_1, \dots, \theta_k$  be  $k$  different values on a parameter  $\theta$ . Bayes' Theorem:

$$p(\theta_i|Data) = \frac{p(Data|\theta_i)p(\theta_i)}{p(Data)} = \frac{p(Data|\theta_i)p(\theta_i)}{\sum_{i=1}^k p(Data|\theta_i)p(\theta_i)}.$$

- If  $\theta$  takes on a continuum of values

$$p(\theta|Data) = \frac{p(Data|\theta)p(\theta)}{\int_{\theta} p(Data|\theta)p(\theta)d\theta}.$$

# The joy of ignoring a normalizing constant

- When *Data* is known,  $p(Data)$  in Bayes' theorem is just a constant that makes  $p(\theta|Data)$  integrate to one. Example:  $x \sim N(\mu, \sigma^2)$

$$p(x) = (2\pi\sigma^2)^{-1/2} \exp \left[ -\frac{1}{2\sigma^2} (x - \mu)^2 \right].$$

- We may write

$$p(x) \propto \exp \left[ -\frac{1}{2\sigma^2} (x - \mu)^2 \right].$$

- Short form of Bayes' theorem

$$p(\theta|Data) \propto p(Data|\theta)p(\theta)$$

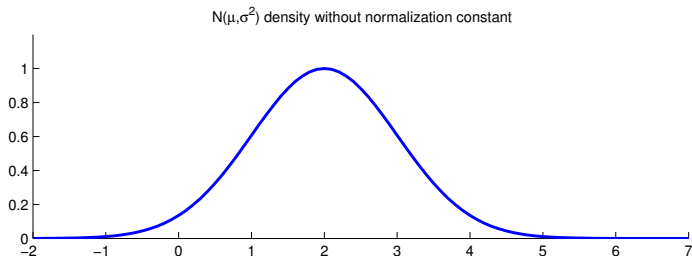
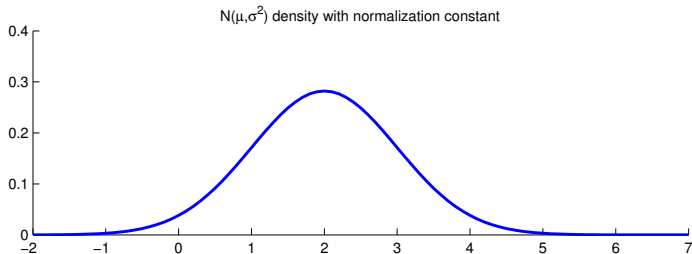
or

$$\text{Posterior} \propto \text{Likelihood} \cdot \text{Prior}$$



# Normalization constant is not important

Illustration that the normalization constant is unimportant



# Bayesian updating

- Suppose: you already have  $x_1, x_2, \dots, x_n$  data points, and the corresponding posterior  $p(\theta|x_1, \dots, x_n)$
- Now, a fresh additional data point  $x_{n+1}$  arrives.
- The posterior based on all available data is

$$p(\theta|x_1, \dots, x_{n+1}) \propto p(x_{n+1}|\theta, x_1, \dots, x_n)p(\theta|x_1, \dots, x_n).$$

- The following are therefore equivalent:
  - Analyzing the likelihood of all data  $x_1, \dots, x_{n+1}$  with the prior based on no data  $p(\theta)$
  - Analyzing the likelihood of the fresh data point  $x_{n+1}$  with the 'prior' equal to the posterior based on the old data  $p(\theta|x_1, \dots, x_n)$ .
- Yesterday's posterior is today's prior.

# Bernoulli trials - Beta prior

- Model:

$$x_1, \dots, x_n | \theta \stackrel{iid}{\sim} \text{Bern}(\theta)$$

- Prior:

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

$$p(y) = \frac{\Gamma(\alpha, \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha-1} (1-y)^{\beta-1} \quad \text{for } 0 \leq y \leq 1.$$

- Posterior

$$\begin{aligned} p(\theta | x_1, \dots, x_n) &\propto p(x_1, \dots, x_n | \theta) p(\theta) \\ &= \theta^s (1-\theta)^f \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &= \theta^{s+\alpha-1} (1-\theta)^{f+\beta-1}. \end{aligned}$$

- But this is recognized as proportional to the  $\text{Beta}(\alpha + s, \beta + f)$  density. That is, the prior-to-posterior mapping reads

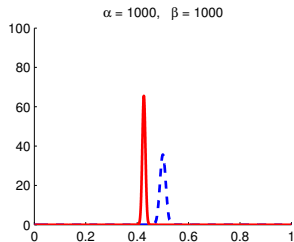
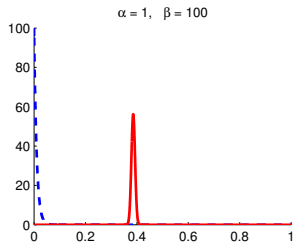
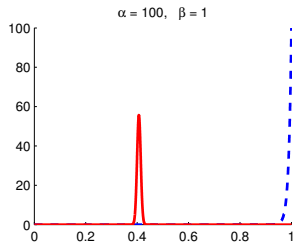
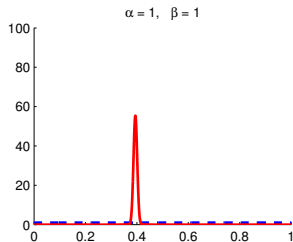
$$\theta \sim \text{Beta}(\alpha, \beta) \xrightarrow{x_1, \dots, x_n} \theta | x_1, \dots, x_n \sim \text{Beta}(\alpha + s, \beta + f).$$

## Bernoulli example: spam emails

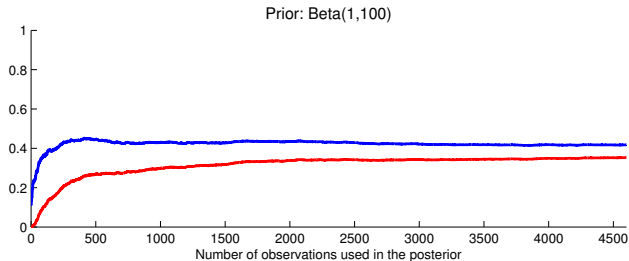
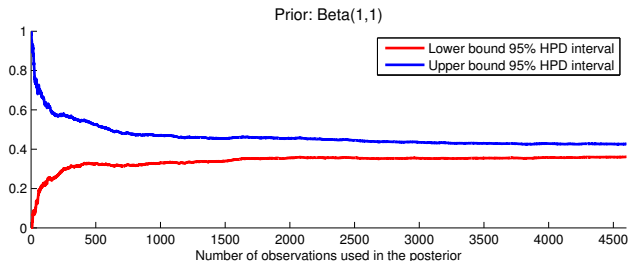
- George has gone through his collection of 4601 e-mails. He classified 1813 of them to be spam.
- Let  $x_i = 1$  if  $i$ :th email is spam. Assume  $x_i|\theta \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$  and  $\theta \sim \text{Beta}(\alpha, \beta)$ .
- Posterior

$$\theta|x \sim \text{Beta}(\alpha + 1813, \beta + 2788)$$

# Spam data: The effect of different priors



# Spam data: Posterior convergence



# Normal data with known variance - uniform prior

- Model:

$$x_1, \dots, x_n | \theta, \sigma^2 \stackrel{iid}{\sim} N(\theta, \sigma^2).$$

- Prior:

$$p(\theta) \propto c$$

- Likelihood (see Technical Appendix A):

$$\begin{aligned} p(x_1, \dots, x_n | \theta, \sigma^2) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp \left[ -\frac{1}{2\sigma^2} (x_i - \theta)^2 \right] \\ &\propto \exp \left[ -\frac{1}{2(\sigma^2/n)} (\theta - \bar{x})^2 \right]. \end{aligned}$$

- Posterior

$$\theta | x_1, \dots, x_n \sim N(\bar{x}, \sigma^2/n)$$

- Prior

$$\theta \sim N(\mu_0, \tau_0^2)$$

- Posterior (see Technical Appendix A)

$$\begin{aligned} p(\theta|x_1, \dots, x_n) &\propto p(x_1, \dots, x_n|\theta, \sigma^2)p(\theta) \\ &\propto N(\theta|\mu_n, \tau_n^2), \end{aligned}$$

where

$$\frac{1}{\tau_n^2} = \frac{n}{\sigma^2} + \frac{1}{\tau_0^2},$$

$$\mu_n = w\bar{x} + (1 - w)\mu_0,$$

and

$$w = \frac{\frac{n}{\sigma^2}}{\frac{n}{\sigma^2} + \frac{1}{\tau_0^2}}.$$



$$\theta \sim N(\mu_0, \tau_0^2) \xrightarrow{x_1, \dots, x_n} \theta|x \sim N(\mu_n, \tau_n^2).$$

Posterior precision = Data precision + Prior precision

Posterior mean =

$$\frac{\text{Data precision}}{\text{Posterior precision}}(\text{Data mean}) + \frac{\text{Prior precision}}{\text{Posterior precision}}(\text{Prior mean})$$

- Conjugate priors
- Poisson model
- 'Non-Informative' priors
- Jeffreys' prior

# Conjugate priors

- Normal likelihood: Normal prior  $\rightarrow$  Normal posterior. (posterior belongs to the same distribution family as prior)
- Binomial likelihood: Beta prior  $\rightarrow$  Beta posterior.
- *Conjugate priors*: Let  $\mathcal{F} = \{p(y|\theta), \theta \in \Theta\}$  be a class of sampling distributions. A family of distributions  $\mathcal{P}$  is conjugate for  $\mathcal{F}$  if

$$p(\theta) \in \mathcal{P} \Rightarrow p(\theta|x) \in \mathcal{P}$$

holds for all  $p(x|\theta) \in \mathcal{F}$ .

- *Natural conjugate prior*:  $p(\theta) = c \cdot p(y_1, \dots, y_n|\theta)$  for some constant  $c$ , i.e. the prior is of the same functional form as the likelihood.

- Likelihood from iid Poisson sample  $y = (y_1, \dots, y_n)$

$$p(y|\theta) = \left[ \prod_{i=1}^n p(y_i|\theta) \right] \propto \theta^{(\sum_{i=1}^n y_i)} \exp(-\theta n),$$

so that the sum of counts  $\sum_{i=1}^n y_i$  is a sufficient statistic for  $\theta$ .

- *Natural conjugate prior for Poisson parameter  $\theta$*

$$p(\theta) \propto \theta^{\alpha-1} \exp(-\theta\beta) \propto \text{Gamma}(\alpha, \beta)$$

which contains the info:  $\alpha - 1$  counts in  $\beta$  observations.

- *Posterior for Poisson parameter  $\theta$* . Multiplying the poisson likelihood and the Gamma prior gives the posterior

$$\begin{aligned} p(\theta|y_1, \dots, y_n) &\propto \left[ \prod_{i=1}^n p(y_i|\theta) \right] p(\theta) \\ &\propto \theta^{\sum_{i=1}^n y_i} \exp(-\theta n) \theta^{\alpha-1} \exp(-\theta \beta) \\ &= \theta^{\alpha + \sum_{i=1}^n y_i - 1} \exp[-\theta(\beta + n)], \end{aligned}$$

which is proportional to the  $\text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n)$  distribution.

- In summary

Model:  $y_1, \dots, y_n | \theta \stackrel{iid}{\sim} \text{Po}(\theta)$

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

Posterior:  $\theta | y_1, \dots, y_n \sim \text{Gamma}(\alpha + \sum_{i=1}^n y_i, \beta + n)$ .

## Poisson example - Number of bomb hits in London

$$n = 576, \sum_{i=1}^n y_i = 229 \cdot 0 + 211 \cdot 1 + 93 \cdot 2 + 35 \cdot 3 + 7 \cdot 4 + 1 \cdot 5 = 537.$$

Average number of hits per region  $= \bar{y} = 537/576 \approx 0.9323$ .

$$p(\theta|y) \propto \theta^{\alpha+537-1} \exp[-\theta(\beta + 576)]$$

$$E(\theta|y) = \frac{\alpha + \sum_{i=1}^n y_i}{\beta + n} \approx \bar{y} \approx 0.9323,$$

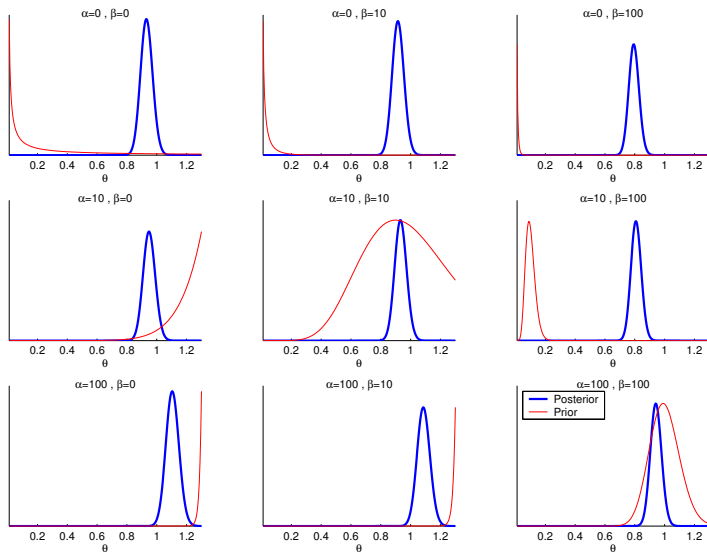
and

$$SD(\theta|y) = \left( \frac{\alpha + \sum_{i=1}^n y_i}{(\beta + n)^2} \right)^{1/2} = \frac{(\alpha + \sum_{i=1}^n y_i)^{1/2}}{(\beta + n)} \approx \frac{(537)^{1/2}}{576} \approx 0.0402.$$

if  $\alpha$  and  $\beta$  are small compared to  $\sum_{i=1}^n y_i$  and  $n$ .

# Poisson bomb hits in London

Analysis of bomb hits in regions of London – Poisson model with Gamma prior



- Bayesian 95% interval: the probability that the unknown parameter  $\theta$  lies in the interval is 0.95. What a relief!
- Approximate 95% credible interval for  $\theta$  (for small  $\alpha$  and  $\beta$ ):

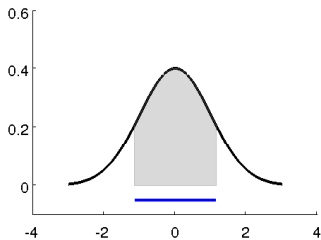
$$E(\theta|y) \pm 1.96 \cdot SD(\theta|y) = [0.8535; 1.0111]$$

- An exact 95% equal-tail interval is  $[0.8550; 1.0125]$  (assuming  $\alpha = \beta = 0$ )
- An exact Highest Posterior Density (HPD) interval is  $[0.8525; 1.0144]$ . Obtained numerically, assuming  $\alpha = \beta = 0$ .

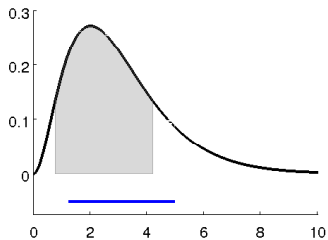


# Illustration of different interval types

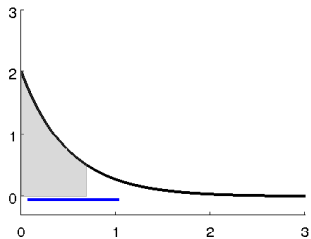
Symmetric distribution



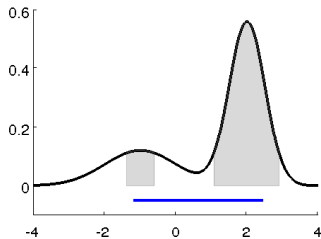
Skewed distribution



Skewed monotone distribution



Bimodal distribution



# Non-informative priors

- ... do not exist!
- ... may be improper and still lead to proper posterior
- Regularization priors
- Ideal communication. Present the posterior distributions for all possible priors.
- Practical communication - Reference priors.
- Cannot demand that users specify priors on high-dimensional in detail. Model the prior in terms of a few hyperparameters.
- Subjective consensus: when extreme priors give essentially the same posterior. This will happen, given enough data as

$$p(\theta|y) \rightarrow N[\hat{\theta}, I^{-1}] \text{ for all } p(\theta) \text{ as } n \rightarrow \infty.$$

- A common non-informative prior is Jeffreys' prior

$$p(\theta) = |I(\theta)|^{1/2},$$

where

$$J(\theta) = -E_{y|\theta} \left[ \frac{d^2 \ln p(y|\theta)}{d\theta^2} \right]$$

is the expected Fisher information.

$$y_1, \dots, y_n | \theta \stackrel{iid}{\sim} \text{Bern}(\theta).$$

$$\ln p(y|\theta) = s \ln \theta + f \ln(1 - \theta)$$

$$\frac{d \ln p(y|\theta)}{d\theta} = \frac{s}{\theta} - \frac{f}{(1 - \theta)}$$

$$\frac{d^2 \ln p(y|\theta)}{d\theta^2} = -\frac{s}{\theta^2} - \frac{f}{(1 - \theta)^2}$$

$$J(\theta) = \frac{E_{y|\theta}(s)}{\theta^2} + \frac{E_{y|\theta}(f)}{(1 - \theta)^2} = \frac{n\theta}{\theta^2} + \frac{n(1 - \theta)}{(1 - \theta)^2} = \frac{n}{\theta(1 - \theta)}$$

Thus, the Jeffreys' prior is

$$p(\theta) = |J(\theta)|^{1/2} \propto \theta^{-1/2}(1 - \theta)^{-1/2} \propto \text{Beta}(\theta|1/2, 1/2).$$

# Jeffreys' prior Binomial vs Negative binomial sampling

- Bernoulli experiment: Perform  $n$  independent trials with success probability  $\theta$  and count the number of successes. Here

$$y|\theta \sim \text{Bin}(\theta)$$

- Inverse Bernoulli experiment: Perform independent trials with success probability  $\theta$  until you have observed  $y$  successes. Here

$$y|\theta \sim \text{NegBin}(\theta)$$

- Exercise: Suppose you performed both of the two experiments and that in both cases you ended up doing  $n$  trials and observed  $y$  successes. Show that the likelihood function conveys the same information on  $\theta$  in both cases, but that Jeffreys prior is not the same in both models. Is this reasonable?

# Properties of Jeffreys prior

- Invariant to 1:1 transformations of  $\theta$ . Doesn't matter which parametrization we derive the prior, it always contains the same info.
- Two models with identical likelihood functions (up to constant) can yield different Jeffreys' prior. Jeffreys' prior does not respect the likelihood principle. The crux of the matter is the expectation with respect to the sampling distribution.
- Jeffreys' prior may be a very complicated (non-conjugate) distribution.
- Problematic in multivariate problems. Dubious results in many standard models.

- Multiparameter models
- Marginalization
- Normal model with unknown variance
- Bayesian analysis of multinomial data
- Bayesian analysis of multivariate normal data

# Marginalization

- Models usually contains several parameter  $\theta_1, \theta_2, \dots$ . Examples:  
 $x_i \stackrel{iid}{\sim} N(\theta, \sigma^2)$ ; multiple regression ...

- The Bayesian computes the joint posterior distribution

$$p(\theta_1, \theta_2, \dots, \theta_p | y) \propto p(y | \theta_1, \theta_2, \dots, \theta_p) p(\theta_1, \theta_2, \dots, \theta_p).$$

... or in vector form:

$$p(\theta) \propto p(y | \theta) p(\theta).$$

- Complicated to graph the joint posterior.
- Some of the parameters may not be of direct interest (nuisance parameters), but are nevertheless needed in the model.
- No problem: just integrate them out (marginalize with respect to, average over) all nuisance parameters.
- Example:  $\theta = (\theta_1, \theta_2)'$ , where  $\theta_2$  is a nuisance. We are interested in the marginal posterior of  $\theta_1$

$$p(\theta_1 | y) = \int p(\theta_1, \theta_2 | y) d\theta_2 = \int p(\theta_1 | \theta_2, y) p(\theta_2 | y) d\theta_2.$$



- Model:

$$y, \dots, y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

- Prior

$$p(\mu, \sigma^2) \propto (\sigma^2)^{-1}$$

- Posterior:

$$\begin{aligned}\mu | \sigma^2, y &\sim N\left(\bar{y}, \frac{\sigma^2}{n}\right) \\ \sigma^2 | y &\sim \text{Inv} - \chi^2(n-1, s^2),\end{aligned}$$

where

$$s^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

is the usual sample variance.

- Simulating the posterior of the normal model with non-informative prior:
  1. Draw  $X \sim \chi^2(n-1)$
  2. Compute  $\sigma^2 = \frac{(n-1)s^2}{X}$  (this a draw from  $\text{Inv-}\chi^2(n-1, s^2)$ )
  3. Draw a  $\mu$  from  $N\left(\bar{y}, \frac{\sigma^2}{n}\right)$  conditional on the previous draw  $\sigma^2$
  4. Repeat step 1-3 many times.
- The sampling is implemented in the R program `NormalNonInfoPrior.R`
- We may derive the marginal posterior analytically as

$$\mu|y \sim t_{n-1}\left(\bar{y}, \frac{s^2}{n}\right).$$

# Multinomial model with Dirichlet prior

- *Data*:  $y = (y_1, \dots, y_K)$ , where  $y_k$  counts the number of observations in the  $k$ th category.  $\sum_{k=1}^K y_k = n$ . Example: brand choices.
- Multinomial model:

$$p(y|\theta) \propto \prod_{k=1}^K \theta_k^{y_k}, \text{ where } \sum_{k=1}^K \theta_k = 1.$$

- *Conjugate prior*:  $\text{Dirichlet}(\alpha_1, \dots, \alpha_K)$

$$p(\theta) \propto \prod_{j=1}^K \theta_j^{\alpha_j - 1}.$$

- Moments of  $\theta = (\theta_1, \dots, \theta_K)' \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$

$$\begin{aligned} E(\theta_k) &= \frac{\alpha_k}{\sum_{j=1}^K \alpha_j} \\ V(\theta_k) &= \frac{E(\theta_k) [1 - E(\theta_k)]}{1 + \sum_{k=1}^K \alpha_k} \end{aligned}$$

- Note that  $\sum_{k=1}^K \alpha_k$  is the precision (inverse variance).

- 'Non-informative':  $\alpha_1 = \dots = \alpha_K = 1$  (uniform and proper).
- Simulating from the Dirichlet distribution:
  - Generate  $x_1 \sim \text{Gamma}(\alpha_1, \beta), \dots, x_K \sim \text{Gamma}(\alpha_K, \beta)$ , independently. Any  $\beta$  will do as long it is the same for all  $x_i$ .
  - Compute  $y_k = x_k / (\sum_{j=1}^K x_j)$ .
  - $y = (y_1, \dots, y_K)$  is a draw from the  $\text{Dirichlet}(\alpha_1, \dots, \alpha_K)$  distribution.
- *Prior-to-Posterior updating:*

*Model:*  $y = (y_1, \dots, y_K) \sim \text{Multin}(n; \theta_1, \dots, \theta_K)$

*Prior :*  $\theta = (\theta_1, \dots, \theta_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$

*Posterior :*  $\theta|y \sim \text{Dirichlet}(\alpha_1 + y_1, \dots, \alpha_K + y_K)$ .

# Multivariate normal - known covariance matrix

- Model:

$$y_1, \dots, y_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$$

where  $\Sigma$  is a known covariance matrix.

- Density

$$p(y|\mu, \Sigma) = |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(y - \mu)' \Sigma^{-1}(y - \mu)\right)$$

- Likelihood:

$$\begin{aligned} p(y_1, \dots, y_n|\mu, \Sigma) &\propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)' \Sigma^{-1}(y_i - \mu)\right) \\ &= |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{tr} \Sigma^{-1} S_\mu\right) \end{aligned}$$

where  $S_\mu = \sum_{i=1}^n (y_i - \mu)(y_i - \mu)'$ .

- Prior:

$$\mu \sim N_p(\mu_0, \Lambda_0)$$

- Posterior:

$$\mu \sim N(\mu_n, \Lambda_n)$$

where

$$\begin{aligned}\mu_n &= (\Lambda_0^{-1} + n\Sigma^{-1})^{-1}(\Lambda_0^{-1}\mu_0 + n\Sigma^{-1}\bar{y}) \\ \Lambda_n^{-1} &= \Lambda_0^{-1} + n\Sigma^{-1}\end{aligned}$$

- Note how the posterior mean is (matrix) weighted average of prior and data information.
- Noninformative prior: let the precision go to zero:  $\Lambda^{-1} \rightarrow 0$ .

# Multivariate normal - Conjugate prior

- Conjugate prior is Normal-IW( $\mu_0, \kappa_0, \Lambda_0, \nu_0$ )

$$\Sigma \sim \text{Inv} - \text{Wishart}(\Lambda_0, \nu_0)$$

$$\mu|\Sigma \sim N(\mu_0, \kappa_0^{-1}\Sigma)$$

- Density:

$$|\Sigma|^{-[(\nu_0+d)/2+1]} \exp\left(-\frac{1}{2}\text{tr}(\Lambda_0\Sigma^{-1}) - \frac{\kappa_0}{2}(\mu - \mu_0)'\Sigma^{-1}(\mu - \mu_0)\right)$$

- Posterior is also Normal IW

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

$$\kappa_n = \kappa_0 + n$$

$$\nu_n = \nu_0 + n$$

$$\Lambda_n = \Lambda_0 + S + \frac{\kappa_0 n}{\kappa_0 + n}(\bar{y} - \mu_0)(\bar{y} - \mu_0)'$$

- Bayesian prediction
- Decision theory



- We may use the estimated model for forecasting a future observation  $\tilde{y}$ .
- *Posterior predictive distribution* ( $y$  denotes available data at the time of forecasting)

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta, y) p(\theta|y) d\theta = \int_{\theta} p(\tilde{y}|\theta) p(\theta|y) d\theta$$

where the last step holds if  $p(\tilde{y}|\theta, y) = p(\tilde{y}|\theta)$ .

- The uncertainty that comes from not knowing  $\theta$  is represented in  $p(\tilde{y}|y)$  by averaging over  $p(\theta|y)$ .

- Let  $y = \sum_{i=1}^n y_i$  and  $\tilde{y}$  the outcome of the next trial

$$\begin{aligned} p(\tilde{y} = 1|y) &= \int_{\theta} p(\tilde{y} = 1|\theta) p(\theta|y) d\theta \\ &= \int_{\theta} \theta p(\theta|y) d\theta = E_{\theta|y}(\theta) = \frac{\alpha + y}{\alpha + \beta + n}. \end{aligned}$$

- Uniform prior ( $\alpha = \beta = 1$ )

$$p(\tilde{y} = 1|y) = \frac{y + 1}{n + 2}.$$

- Assume the uniform prior  $p(\theta) \propto c$ .

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta)p(\theta|y)d\theta$$

where

$$\begin{aligned}\theta|y &\sim N(\bar{y}, \sigma^2/n) \\ \tilde{y}|\theta &\sim N(\theta, \sigma^2)\end{aligned}$$

## Simulate from the predictive distribution - Normal model

- 1 Generate a posterior draw of  $\theta$  ( $\theta^{(1)}$ ) from  $N(\bar{y}, \sigma^2/n)$
- 2 Generate a draw of  $\tilde{y}$  ( $\tilde{y}^{(1)}$ ) from  $N(\theta^{(1)}, \sigma^2)$  (note the mean)
- 3 Repeat steps 1 and 2 a large number of times ( $N$ ) with the result:
  - Sequence of posterior draws:  $\theta^{(1)}, \dots, \theta^{(N)}$
  - Sequence of predictive draws:  $\tilde{y}^{(1)}, \dots, \tilde{y}^{(N)}$ .

- $\theta^{(1)} = \bar{y} + \varepsilon^{(1)}$ , where  $\varepsilon^{(1)} \sim N(0, \sigma^2/n)$ . (Step 1).
- $\tilde{y}^{(1)} = \theta^{(1)} + v^{(1)}$ , where  $v^{(1)} \sim N(0, \sigma^2)$ . (Step 2).
- $\tilde{y}^{(1)} = \bar{y} + \varepsilon^{(1)} + v^{(1)}$ .
- $\varepsilon^{(1)}$  and  $v^{(1)}$  are independent.
- The sum of two normal random variables follows a normal distribution, so  $\tilde{y}$  follows a normal distribution with

$$E(\tilde{y}|y) = E(\tilde{y}) = \bar{y}$$

$$V(\tilde{y}|y) = \frac{\sigma^2}{n} + \sigma^2 = \sigma^2 \left(1 + \frac{1}{n}\right).$$

- Note that the estimation uncertainty ( $\sigma^2/n$ ) is typically much less important than the intrinsic population uncertainty,  $\sigma^2$ .

- It easy to see that the predictive distribution is normal.
- The mean can be obtained from

$$E_{\tilde{y}|\theta}(\tilde{y}|\theta) = \theta$$

and then remove the conditioning on  $\theta$  by averaging over  $\theta$

$$E(\tilde{y}|y) = E_{\theta|y}(\theta) = \mu_n \text{ (Posterior mean of } \theta\text{)}.$$

- The predictive variance of  $\tilde{y}$  can be obtained from the conditional variance formula

$$\begin{aligned} V(\tilde{y}|y) &= E_{\theta|y}[V_{\tilde{y}|\theta}(\tilde{y}|\theta)] + V_{\theta|y}[E_{\tilde{y}|\theta}(\tilde{y}|\theta)] \\ &= E_{\theta|y}(\sigma^2) + V_{\theta|y}(\theta) \\ &= \sigma^2 + \tau_n^2 \\ &= (\text{Population variance} + \text{Posterior variance of } \theta). \end{aligned}$$

- In summary:

$$\tilde{y}|y \sim N(\mu_n, \sigma^2 + \tau_n^2).$$

- Let  $\theta$  be an unknown quantity. State of nature. Examples: Future inflation, Global temperature, Disease.
- Let  $a \in \mathcal{A}$  be an action. Ex: Interest rate, Energy tax, Operation.
- Choosing action  $a$  when state of nature turns out to be  $\theta$  gives utility

$$U(a, \theta)$$

- Alternatively loss  $L(a, \theta) = -U(a, \theta)$ .

- Loss table:

	$\theta_1$	$\theta_2$
$a_1$	$L(a_1, \theta_1)$	$L(a_1, \theta_2)$
$a_2$	$L(a_2, \theta_1)$	$L(a_2, \theta_2)$

- Example utility functions:

- Linear:  $L(a, \theta) = |a - \theta|$
- Quadratic:  $L(a, \theta) = (a - \theta)^2$
- Lin-Lin:

$$L(a, \theta) = \begin{cases} c_1 & \text{if } a \leq \theta \\ c_2 & \text{if } a > \theta \end{cases}$$

- Ad hoc decision rules:
  - *Minimax*. Choose the decision that minimizes the maximum loss.
  - *Minimax-regret*: Choose the decision rule that gives you least regret when you eventually find out the true value of  $\theta$ .
- Bayesian axiomatic theory gives you the rule: Choose the action that maximizes the (posterior) expected utility:

$$a_{\text{bayes}} = \operatorname{argmax}_{a \in \mathcal{A}} E_{p(\theta|y)}[L(a, \theta)],$$

where  $E_{p(\theta|y)}$  denotes the posterior expectation.

- Using simulated draws  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(N)}$  from  $p(\theta|y)$  :

$$E_{p(\theta|y)}[L(a, \theta)] \approx N^{-1} \sum_{i=1}^N L(a, \theta^{(i)})$$

- *Separation principle*: The analysis of uncertainty (i.e. the posterior of  $\theta$ ) is completely separated from the utilities of the choices.



- Choosing a point estimator is a decision problem.
- Which to choose: posterior median, mean or mode?
- It depends on your loss function:
  - Linear loss  $\rightarrow$  Posterior median is optimal
  - Quadratic loss  $\rightarrow$  Posterior mean is optimal
  - Lin-Lin loss  $\rightarrow c_1/(c_1 + c_2)$  quantile of the posterior is optimal
  - Zero-one loss  $\rightarrow$  Posterior mode is optimal
- Similar analysis can be used to select interval type: symmetric or HPD?

# Combining expert judgement

- Available:  $K$  unbiased expert forecast/judgement:  $\mathbf{y} = (y_1, y_2, \dots, y_K)'$  and

$$\mathbf{y}|\theta \sim N(\theta \cdot \mathbf{1}, \Sigma)$$

where  $\Sigma$  is assumed to be known.

- Assuming a uniform prior for  $\theta$ , the posterior distribution is of the form

$$\theta|\mathbf{y} \sim N(\mathbf{w}'\mathbf{y}, \psi^2)$$

where

$$\mathbf{w} = \frac{\mathbf{1}' \cdot \Sigma^{-1}}{\mathbf{1}' \cdot \Sigma^{-1} \cdot \mathbf{1}}$$
$$\psi = \frac{1}{\sqrt{\mathbf{1}' \cdot \Sigma^{-1} \cdot \mathbf{1}}}$$

- Weights can be negative, and it makes sense!
- Correlation between expert's forecasts are important. A poor expert can get a sizeable weight if she is negatively correlated with the rest.
- Estimate  $\Sigma$  with past forecast errors. Difficulty: estimate  $\Sigma$  precisely. A decent prior on  $\Sigma$  can help!