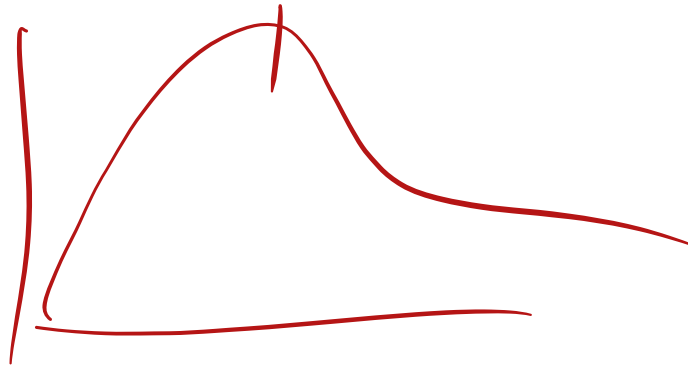


# Aside: Bayes Estimators

# Why the posterior mean?

- Often times we need to make a “decision” by providing a single estimate
- The posterior provides a full distribution over  $\theta$ , which can be summarized in infinitely many ways
- Specify a *loss function* which describes the cost of estimating  $\hat{\theta}$  when the truth is  $\theta$



# Bayes Estimators

- The *loss function*:  $L(\hat{\theta}, \theta)$ 
  - Squared error:  $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$
  - Absolute error:  $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$
- The *Bayes risk* is the posterior expected loss:  $E_{\theta|y}[L(\hat{\theta}, \theta)]$
- Summarize the posterior by minimizing the Bayes risk.
- An estimator  $\hat{\theta}$  is said to be a Bayes estimator if it minimizes the Bayes risk among all estimators.

$$\int_{\theta} L(\hat{\theta}, \theta) p(\theta|y) d\theta$$

# Examples Squared error loss

$$\hat{\theta}_{\text{B.E.}} = \min_{\hat{\theta}} E_{\theta|y}(\hat{\theta} - \theta)^2 = \min_{\hat{\theta}} \int (\hat{\theta} - \theta)^2 p(\theta | y) d\theta$$

$$= E[\theta | y]$$

Do THIS

$$\frac{d}{d\hat{\theta}} \int (\hat{\theta} - \theta)^2 p(\theta | y) d\theta = \int 2(\hat{\theta} - \theta) p(\theta | y) d\theta$$

# The Bias-Variance Tradeoff

- The prior distribution (usually) makes your estimator biased...
- But the prior distribution also (usually) reduces the variance!
- Example: compute the frequentist mean and variance of the posterior mean.

# Example: IQ scores

$$y_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma=13), \mu=112$$
$$\mu \sim N(\mu_0, \frac{\sigma}{\sqrt{n_0}})$$

- Scoring on **IQ tests** is designed to yield a  $N(100, 15)$  distribution for the general population
- We observe IQ scores for a sample of  $n$  individuals from a particular town and estimate  $\mu$ , the town-specific IQ score
- If we lacked knowledge about the town, a natural choice would be  $\mu_0 = 100$
- Suppose the true parameters for this town are  $\mu = 112$  and  $\sigma = 13$ 
  - The town is smarter on average than the general population

$$y_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma=1), \mu=1/2$$

$$\mu \sim N(\mu_0, \frac{\sigma}{\tau_{k_0}})$$

$$\underline{E[\mu | y_1, \dots, y_n] = w \bar{y} + (1-w)\mu_0}$$

$$w = \frac{n}{n + k_0}$$

Post. mean estimator:  $w \bar{Y} + (1-w)\mu_0$

$$k_0 = 0, \rightarrow \hat{\mu}_{MLE} = \bar{Y}$$

$$E_{y|\mu}[(\hat{\mu} - \mu)^2] = \text{Var}(\hat{\mu}) + \text{Bias}(\hat{\mu})^2$$

$$\text{Var}(\hat{\mu}_{MLE}) = \text{Var}(\bar{Y}) = \frac{\sigma^2}{n}$$

$$\text{Bias}(\hat{\mu}_{MLE}) = E[\bar{Y}] - \mu = 0$$

$$\text{MSE}(\hat{\mu}_{MLE}) = \frac{\sigma^2}{n}$$


---

$$\text{Var}(\hat{\mu}_{pm}) = \text{Var}(w\bar{Y} + (1-w)\mu_0) = w^2 \frac{\sigma^2}{n}$$

$$\begin{aligned} \text{Bias}(\hat{\mu}_{pm}) &= E[w\bar{Y} + (1-w)\mu_0] - \mu \\ &= w\mu + (1-w)\mu_0 - \mu \\ &= (1-w)(\mu_0 - \mu) \end{aligned}$$

$$\text{MSE}(\hat{\mu}_{pm}) = w^2 \frac{\sigma^2}{n} + (1-w)^2 (\mu_0 - \mu)^2$$

$$\text{MSE}(\hat{\mu}_{pm}) \begin{matrix} > \\ < \\ \sim \\ ?? \end{matrix} \text{MSE}(\hat{\mu}_{MLE})$$



# Example: IQ scores

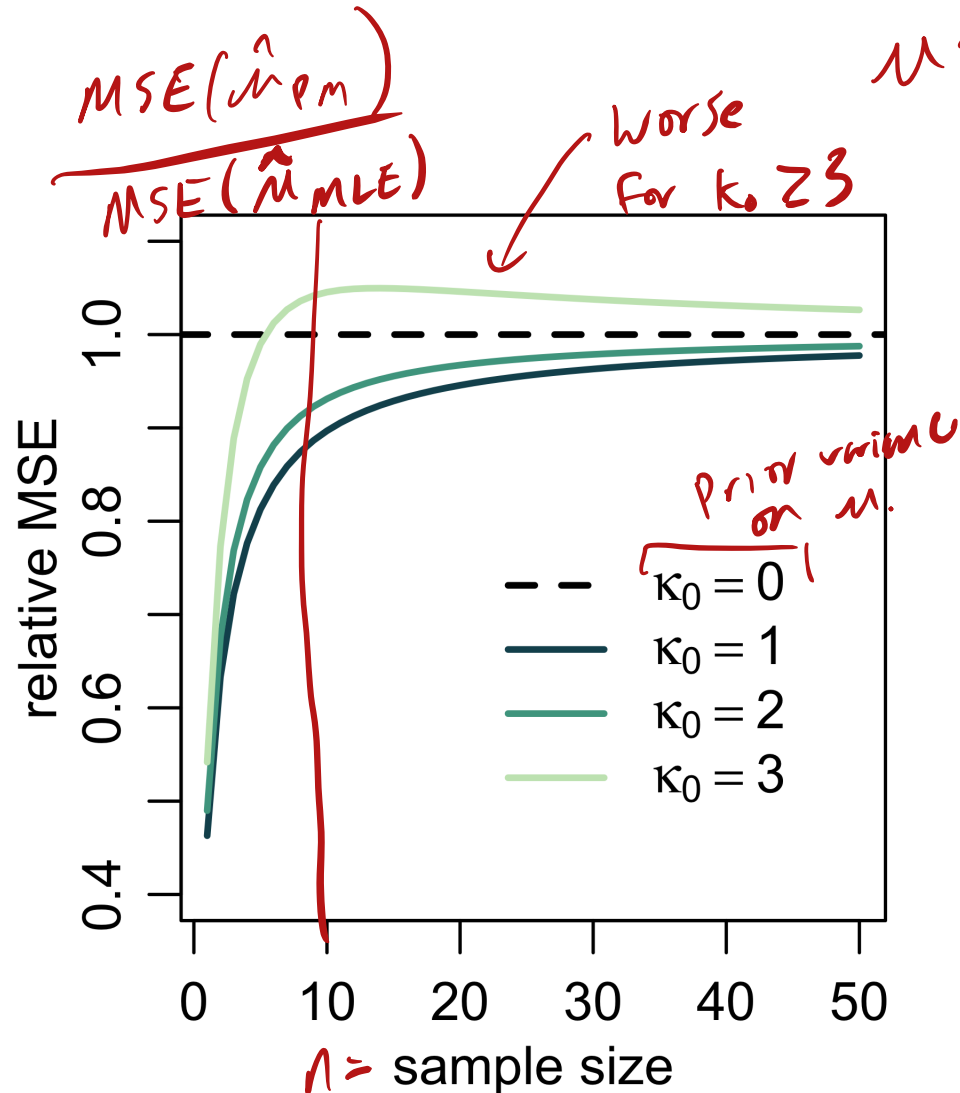
- What is the mean squared error of the MLE? MSE of the posterior mean?

$$\sigma = 13$$

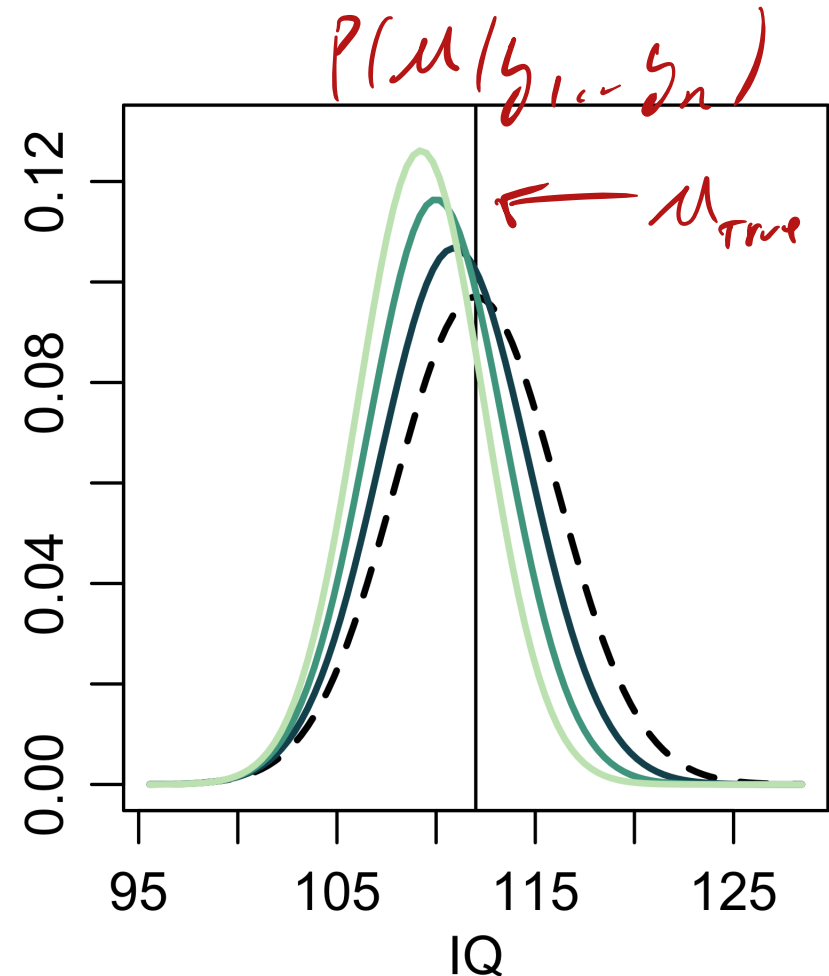
- $\text{MSE}[\hat{\mu}_{MLE}] = \text{Var}[\hat{\mu}_{MLE}] = \frac{\sigma^2}{n} = \frac{169}{n}$
- $\text{MSE}[\hat{\mu}_{PM}|\theta_0] = w^2 \frac{169}{n} + (1 - w)^2 144$
- Reminder:  $w = \frac{n}{\kappa_0 + n}$ . For what values of  $n$  and  $\kappa_0$  is the MSE smaller for the posterior mean estimator than the maximum likelihood?

$$\begin{matrix} \mu_0 & \mu_1 \\ 100 & 112 \end{matrix} \quad \left( \begin{matrix} 100 \\ 112 \end{matrix} \right)^2$$

# Example: IQ scores



$$\mu \sim N(\mu_0, \frac{\sigma^2}{k_0})$$



# The Multivariate Normal Distribution

$$Y_{p \times 1} = \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix} \sim N_p(\mu, \Sigma)$$

$$P(Y=y \mid \mu, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left[-\frac{1}{2}(y-\mu)' \Sigma^{-1} (y-\mu)\right]$$

$$\mu \in \mathbb{R}^p$$

$$\Sigma \in S_p^+ \text{ cone of positive def. matrices.}$$

$$\underbrace{a' \Sigma a > 0}_{\frac{p(p+1)}{2}} \text{ free parameters.}$$

$$P(\mu | \Sigma, y_1, \dots, y_n)$$

$$\mu \sim \mathcal{N}_p(\mu_0, \Lambda_0)$$

$$L(\mu) \propto \prod_{i=1}^n \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (y_i - \mu)' \Sigma^{-1} (y_i - \mu)\right]$$

$$\propto \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)' \Sigma^{-1} (y_i - \mu)\right)$$

$$\propto \exp\left(-\frac{1}{2} \mu^T A \mu + \mu^T b\right)$$

$$A = n \Sigma^{-1} \quad , \quad b = n \Sigma^{-1} \bar{y}$$

$$P(\mu) \propto \exp(\mu^T \Lambda_0^{-1} \mu + \mu^T \Lambda_0^{-1} \mu_0)$$

$$P(\mu | y_1, \dots, y_n) \propto \exp(\mu^T A_n \mu + \mu^T (b + \underbrace{\Lambda_0^{-1} \mu_0}_{n \Sigma^{-1} \bar{y} + \Lambda_0^{-1} \mu_0}))$$

$$A_n = \Sigma^{-1} + \Lambda_0^{-1}$$

$$n \Sigma^{-1} \bar{y} + \Lambda_0^{-1} \mu_0$$

$$\sim \mathcal{N}(\mu_n, \Sigma_n)$$

$$\Sigma_n = (n \Sigma^{-1} + \Lambda_0^{-1})^{-1}$$

$$\mu_n = (n \Sigma^{-1} + \Lambda_0^{-1})^{-1} (n \Sigma^{-1} \bar{y} + \Lambda_0^{-1} \mu_0)$$

Special  
Case

$$\Lambda_0^{-1} = k_0 \Sigma^{-1}$$

$$\mu_n = (n \Sigma^{-1} + k_0 \Sigma^{-1})^{-1} (n \Sigma^{-1} \bar{y} + k_0 \Sigma^{-1} \mu_0)$$

$$= \frac{1}{n+k_0} \Sigma (n \Sigma^{-1} \bar{y} + k_0 \Sigma^{-1} \mu_0)$$

$$= \frac{n}{n+k_0} \bar{y} + \frac{k_0}{n+k_0} \mu_0$$

$\frac{\sigma^2}{k_0}$   
(uninformative)  
prior

$$\mu \sim N(\mu_0, \tau^2)$$

$\Sigma$  unknown,  $\mu = 0$  known

$$\vec{y}_1, \dots, \vec{y}_n \sim N(0, \Sigma)$$

$$(\vec{y}_1 \quad \vec{y}_2 \quad \dots \quad \vec{y}_n)$$

$P \times n$

$$S = Y Y^T$$

sum of  
squares  
matrix

$$= \begin{pmatrix} \sum_{i=1}^n y_{i1}^2 & \left( \sum_{i=1}^n y_{i1} y_{ik} \right) & \dots \\ \vdots & \ddots & \ddots \\ \sum_{i=1}^n y_{ip}^2 & \dots & \dots \end{pmatrix}$$

ith row  
kth column

$$S \sim \text{Wishart}(n, \Sigma)$$

$$P(S) \propto |S|^{\frac{-(n-p-1)}{2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} S)}$$

$n > p \rightarrow Y Y^T$  is pos. def. w.p. 1.  
 $\rightarrow$  symmetric.

$$E[S] = n \Sigma$$

$$\frac{Y Y^T}{n} \rightarrow \Sigma$$

$$X \sim \text{Wish}, \quad X^{-1} \sim \text{Inv-Wish}.$$

$\Sigma \sim \text{Inv-Wish}$  is conjugate prior

$$L(\Sigma) \propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \underbrace{Y' \Sigma^{-1} Y}_{\text{scalar}}\right)$$

$$\begin{aligned} \text{tr}(ABC) &= \text{tr}(BCA) \\ &\propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{tr}(Y' \Sigma^{-1} Y)\right) \\ &= |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{tr}(\Sigma^{-1} \underbrace{Y Y'}_{\text{pp}})\right) \end{aligned}$$

$$\Sigma \sim \text{IW}(\nu, \Lambda_0) \rightarrow |\Sigma|^{-\frac{\nu+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} \Lambda_0)}$$

$$\rightarrow P(\Sigma | y_1, \dots, y_n) \sim \underline{\underline{\text{IW}(n+\nu, S+\Lambda_0)}}$$

$$\Sigma \sim \text{IW}(\nu, \Lambda_0) \quad \underline{\underline{\nu > p-1}}$$

$$\begin{aligned} \text{Default: } \Lambda_0 &= I \\ &\text{IW}(\text{small}, I) \end{aligned}$$

$$\tau_{eff} \propto |\varepsilon|^{-\frac{(p+1)}{2}}$$



# Dirichlet-Multinomial

## Metagenomics example

- Metagenomics is the study of genetic material recovered directly from environmental samples
- Map counts of genetic material to counts of microbial species
- Assume species are sampled with replacement
  - Observed sample is a multinomial distribution
- Total counts isn't meaningful (hard to control how much total sample)
- Relative counts are meaningful

# Multinomial Density

- Let  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$  with  $\sum_{i=1}^k \theta_i = 1$  # of species  $\sum \theta_i = 1$ ,  $\sum_{k=1}^K Y_k = n$

- If  $Y = (Y_1, \dots, Y_K) \sim \text{Mult}(n, \theta)$ , then:

- $Y_i \sim \text{Bin}(n, \theta_i)$

- $Y_i + Y_j \sim \text{Bin}(n, \theta_i + \theta_j)$

$\rightarrow P(Y|\theta) \propto \prod_{k=1}^K \theta_k^{Y_k}$

- What is  $\hat{\theta}_{MLE}$ ?

$$\hat{\theta}_{k,MLE} = \frac{Y_k}{\sum Y_k}$$

# Dirichlet Distribution

Generalizes Beta, conjugate for Multinomial.

$$\underline{\vec{\theta}_k \sim \text{Dir}(\alpha_1, \dots, \alpha_k)}$$

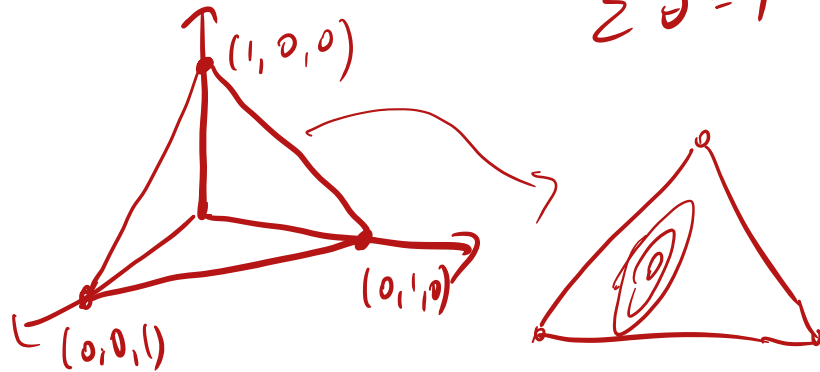
$$E[\theta_k] = \frac{\alpha_k}{\sum_{i=1}^K \alpha_i}$$

$$P(\theta | y) \propto \prod_k \theta_k^{(y_k + \alpha_k - 1)} \rightarrow \text{Dir}(y_1 + \alpha_1, \dots, y_k + \alpha_k)$$

$$K=3$$

$$(\theta_1, \theta_2, \theta_3)$$

$$\sum \theta = 1$$



$$z_i \sim \text{Gam}(\alpha_i, 1)$$

$$\left( \frac{z_1}{\sum_{k=1}^K z_k}, \dots, \frac{z_K}{\sum_{k=1}^K z_k} \right) \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$