

$$\begin{cases} Y_1, \dots, Y_n \sim N(\mu, \sigma^2), \sigma^2 \text{ is known.} \\ \mu \sim N(\mu_0, \frac{\sigma^2}{k_0}) \quad (\text{Prior}) \end{cases}$$

↑ prior parameters

## Posterior Bayes Estimators

$$\rightarrow P(\mu | Y_1, \dots, Y_n, \sigma^2) \sim N(\mu_n, \sigma_n^2)$$

$$\mu_n = w \bar{y} + (1-w) \mu_0, \quad w = \frac{n}{n+k_0}$$

$$\sigma_n^2 = \frac{\sigma^2}{n+k_0} = w \frac{\sigma^2}{n}$$

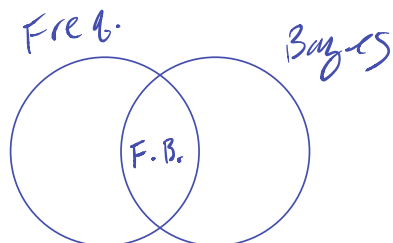
$$E[\mu | y_1, \dots, y_n, \sigma^2] = w\bar{y} + (1-w)\mu_0 \quad (\text{estimate})$$

lowercase

What are the frequentist properties

of  $E[\mu | y_1, \dots, y_n, \sigma^2]$  estimator!!

capital



$$w\bar{y} + (1-w)\mu_0$$

# Estimators: Bayes / Frequentist Unification


- Bayesian inference provides a straightforward procedure for producing estimators given your prior beliefs.



1. Compute posterior distribution

2. Summarize the posterior distribution with a point estimator (e.g. posterior mean or posterior mode) and a probability interval

- Frequentists provide tools for evaluating the sampling properties of an estimator.

- 
- Bias, variance and MSE of an estimator
  - Well-calibrated probability intervals

- Both are useful!

# The Bias-Variance Tradeoff

Reminder: an estimator is a random variable, an estimate is a constant

- *Bias*: systematic sampling error of the estimator
- *Variance*: variance of the estimator (from sampling & measurement error)
- Often we evaluate an estimator in terms of mean square error:  
$$\text{MSE}(\hat{\theta}) = E_Y(\hat{\theta} - \theta)^2$$
- The Bias-Variance tradeoff:  $\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2$

# The Bias-Variance Tradeoff

- Variance of an estimator comes sampling from a population
  - If you were to repeatedly draw new samples of the same size how much would your estimates vary?
  - e.g. if  $y_i \sim N(\mu, \sigma^2)$  then  $\text{Var}(\bar{Y}) = \sigma^2/n$

$$\hat{\mu}_{MLE} = \bar{Y}$$

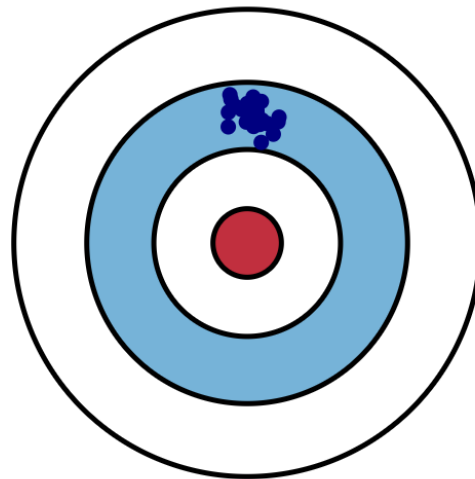
$$E[\bar{Y}] = \mu \quad (\text{unbiased})$$

$$\begin{aligned} \text{Var}(\bar{Y}) &= \text{Var}\left(\frac{1}{n} \sum Y_i\right) = \frac{1}{n^2} \sum \text{Var}(Y_i) \\ &= \frac{\sigma^2}{n} \end{aligned}$$

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

# Bias

The expected difference between the estimate and the response

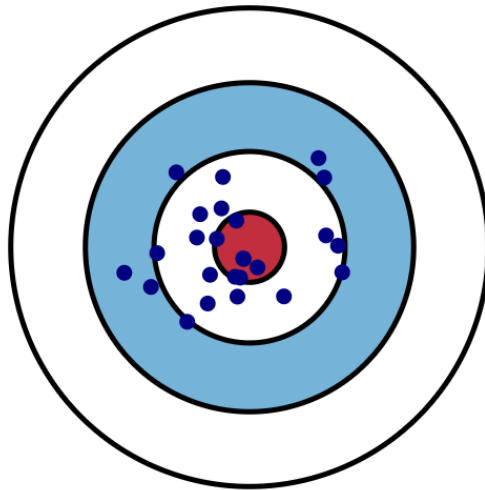


Statistical definition of bias:

$$E_Y[\hat{\theta} - \theta]$$

# Variance

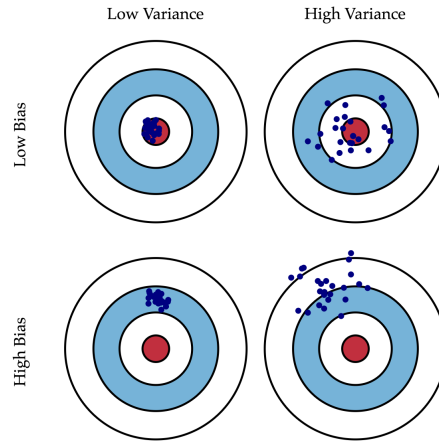
How variable is the prediction about its mean?



Statistical definition of variance:

$$E_Y[\hat{\theta} - E_Y[\hat{\theta}]]^2$$

# Bias and Variance



$$\underbrace{\text{MSE}(\hat{\theta})}_{\text{accuracy}} = \underbrace{\text{Var}(\hat{\theta})}_{\text{variance}} + \underbrace{\left( E[\hat{\theta} - \theta] \right)^2}_{\text{bias squared}}$$



# 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.

$$E[\mu | Y_1, \dots, Y_n, \sigma^2] = \hat{\mu}_{PM} = w \bar{Y} + (1-w)\mu_0$$

Bias:  $E[w\bar{Y} + (1-w)\mu_0 - \mu] = wE[\bar{Y}] + (1-w)E[\mu_0] - \mu$   
 $= w\mu + (1-w)\mu_0 - \mu = (1-w)(\mu_0 - \mu)$

Var:  $Var(\mu | Y_1, \dots, Y_n, \sigma^2) = Var(w\bar{Y} + (1-w)\mu_0) =$   
 $= w^2 Var(\bar{Y}) + (1-w)^2 Var(\mu_0)$   
 $= w^2 Var(\bar{Y}) \leq Var(\bar{Y}) = Var(\hat{\mu}_{MLE})$

*(Note:  $Var(\mu_0) = 0$  (const) is crossed out in the original image)*

## Example: IQ scores

- Scoring on **IQ tests** is designed to yield a  $N(100, 15)$  distribution for the general population  $sd$   
"
- 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  $\mu_0 = 100$   $\hookrightarrow I.V.$
- Suppose the true parameters for this town are  $\mu = 112$  and  $\sigma = 13$ 
  - The town is smarter on average than the general population

How does  $\bar{Y}$  compare to

$$\hat{\mu}_{PM} = w\bar{Y} + (1-w)100 \quad ??$$

## 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?

$$\text{Bias}^2$$

$$\frac{169}{n} \stackrel{?}{>} w^2 \frac{169}{n} + (1-w)^2 144$$

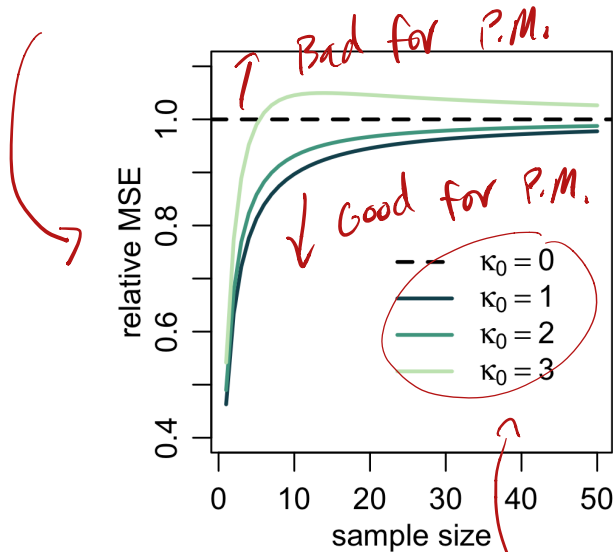
$$\begin{aligned} \text{Bias} &= (1-w)(100-112) \\ \Rightarrow \text{Bias}^2 &= (1-w)^2 144 \end{aligned}$$

$$\text{MSE}(\hat{\mu}_{MLE})$$

$$\text{MSE}(\hat{\mu}_{PM})$$

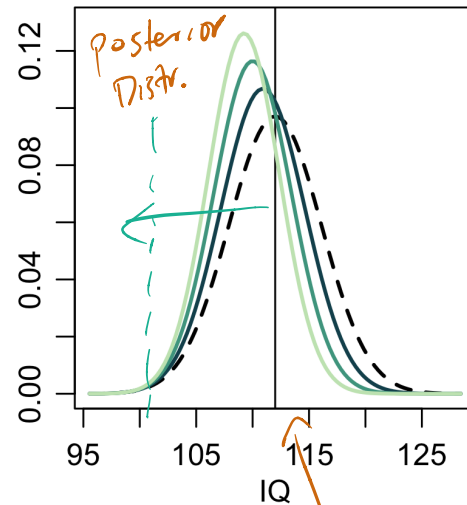
# Example: IQ scores

$$\frac{MSE(\hat{u}_{PM})}{MSE(\hat{u}_{MLE})}$$



$$w = \frac{n}{n + \kappa_0}$$

$n \rightarrow \infty, w \rightarrow 1$   
e.g.  $\hat{u}_{PM} \rightarrow \hat{u}_{MLE}$



$MSE_{PM}$ :

$$w \frac{169}{n} + (1-w)^2 144$$

Strength of  
prior knowledge  
(# of pseudo-obs,  
determines  $w$ )

Truth.

# Decision Theory

# 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  $\theta$

$$L(\hat{\theta}, \theta)$$

# Bayes Estimators

- The *loss function*:  $L(\hat{\theta}, \theta)$

penalize the difference  
between estimate & truth.

- 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:

$$E_{\theta|y}[L(\hat{\theta}, \theta)] = \int L(\hat{\theta}, \theta) p(\theta | y) d\theta$$

Monte Carlo

1. Sample  $\theta^s$   
from  $P(\theta | y)$

2.  $L(\hat{\theta}, \theta^s)$

Average all losses.

- Choose an estimator of  $\theta$  based on minimizing the Bayes risk.
- An estimator  $\hat{\theta}$  is said to be a **Bayes estimator** if it minimizes the Bayes risk among all estimators.

**Squared error loss :** *What  $\hat{\theta}$  should I choose?*

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

Differentiate with respect to  $\hat{\theta}$  and set equal to zero:

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

*✓ This*

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

$$\int 2(\hat{\theta} - \theta) P(\theta|y) d\theta = 0$$

$$\cancel{2} \int \hat{\theta} P(\theta|y) d\theta = \cancel{2} \int \theta P(\theta|y) d\theta$$



$$\hat{\theta} \underbrace{\int P(\theta|y) d\theta}_1 = \int \theta P(\theta|y) d\theta$$

$$\tilde{\theta} = E[\theta|y]$$

✓ Posterior  
Mean  
Definition.

# Absolute loss

$$\min_{\hat{\theta}} E_{\theta|y} |\hat{\theta} - \theta| = \min_{\hat{\theta}} \int |\hat{\theta} - \theta| p(\theta | y) d\theta$$

Aside

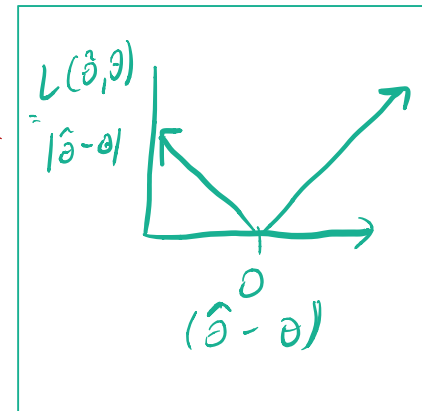
Differentiate with respect to  $\hat{\theta}$  and set equal to zero:

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

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

$$\int \text{sign}(\hat{\theta} - \theta) p(\theta | y) d\theta =$$

$\hookrightarrow -1 \text{ if } \hat{\theta} \leq \theta$   
 $+1 \text{ if } \hat{\theta} \geq \theta$



$$\text{sign}(x) = +1 \text{ if } x \geq 0$$

$$-1 \text{ if } x < 0$$

$$\int_{\theta = -\infty}^{\theta = \hat{\theta}} +1 P(\theta|y) d\theta + \int_{\theta = \hat{\theta}}^{\theta = +\infty} -1 P(\theta|y) d\theta$$

This is just a CDF!!

$$P(\theta < \hat{\theta} | y) - (1 - P(\theta < \hat{\theta} | y)) = 0$$

$$2P(\theta < \hat{\theta} | y) = 1$$

$$P(\theta < \hat{\theta} | y) = 1/2$$

$\hat{\theta}$  is the posterior median

Bayes Estimator for absolute loss.

# Loss functions in practice

- Squared error and absolute error are good default loss functions
  - Motivated largely by mathematical considerations
- In practice we should define a loss function specific to our problem
- Loss in dollars? Loss in "quality of life"?

# Decision making: flu example

- The CDC produces estimates of the expected prevalence and severity of flu during flu season
- Assume  $\theta$  represents severity of the flu
- $p(\theta | y)$  is CDC posterior distribution based on initial data about the upcoming flu season
- $\hat{\theta}$  determines how much flu vaccine to make. How do we determine  $L(\hat{\theta}, \theta)$ ?

e.g.  
Poisson  
rate

+ Shelf life of vaccine  
+ Too little = death  
+ Too much is wasted money