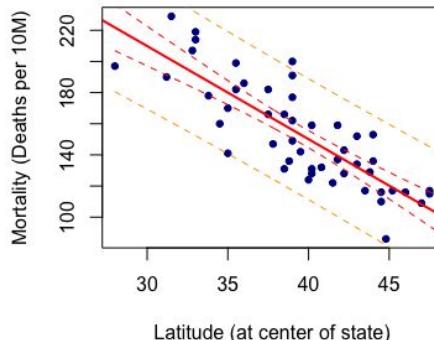


# Sampling distributions of $\hat{\beta}$

*14.4 of Rice*

08/05/2021

# In the previous lecture,



- Simple linear regression (SLR):
  - $(\hat{\beta}_0, \hat{\beta}_1)^T$  is a bivariate Normal random variable;
  - $(\hat{\beta}_0, \hat{\beta}_1)^T$  is independent of residuals;
  - $\text{RSS} = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \sim \chi_{n-2}^2$ .
- Confidence intervals & bands
  - CIs of  $\beta_0$  and  $\beta_1$ :  

$$\hat{\beta}_k \pm t_{n-2}(\alpha/2) \cdot \text{se}(\hat{\beta}_k)$$
  - CIs of the regression line  $\beta_0 + \beta_1 x$ :  

$$\hat{y} \pm t_{n-2}(\alpha/2) \cdot \text{se}(\hat{y})$$
  - PIs of the response:  

$$\hat{y} \pm t_{n-2}(\alpha/2) \cdot \text{se}(y - \hat{y})$$
- Multiple linear regression (MLR):
  - $Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \epsilon_i$ ,  $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ .
  - $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y}$  ←
  - To avoid collinearity, the design matrix  $\mathbf{X}$  must have rank  $p+1$ .  $\Rightarrow R^2$  and  $VIF_i$

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}$$

Design matrix

$$\text{rank}(\mathbf{X}) = p+1$$

$$\text{var}(\vec{\beta} \mid \vec{Y}) = A \Sigma A^T$$

↑  
Var( $\vec{Y}$ )

Sampling distribution of  $\vec{\beta}_{\text{hat}}$

**Proposition D'.** Under the MLR assumptions,  $\vec{\beta}_{\text{hat}}$  is multivariate Normal with:

$$E(\vec{\beta}_{\text{hat}}) = \vec{\beta}, \quad \text{var}(\vec{\beta}_{\text{hat}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

Proof.  $\vec{Y} \sim N(\mathbf{X}\vec{\beta}, \sigma^2 \mathbf{I})$

$$\vec{\beta}_{\text{hat}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y} \sim N(\vec{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}).$$

$$E \vec{\beta}_{\text{hat}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \vec{\beta} = \vec{\beta}.$$

$$\begin{aligned} \text{var}(\vec{\beta}_{\text{hat}}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \cancel{\mathbf{X}^T \mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \end{aligned}$$

## Sampling distribution of RSS

projection matrix of  $\text{span}(\mathbf{X})$ :

$$H = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$\vec{\beta}_{\text{hat}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_P \end{pmatrix}$$

**Lemma D'.** Under the MLR assumptions,

and  $\text{RSS} = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \dots - \hat{\beta}_p X_{ip})^2 \sim \chi_{n-p-1}^2$

$$\vec{\beta}_{\text{hat}} \perp\!\!\!\perp \overrightarrow{Y} - \mathbf{X} \vec{\beta}_{\text{hat}},$$

1. A matrix  $R$  is idempotent if  $R^2 = R$ .
2. If  $Z \sim N(0, I)$  and  $R$  is symmetric and idempotent of rank  $r$ , then  $Z^T R Z \sim \chi_r^2$ .

Proof\*.

$$\hat{e}_i = Y_i - \hat{\beta}_0 - \dots - \hat{\beta}_p X_{ip}$$

$$= Y_i - (\underbrace{1, X_{i1}, \dots, X_{ip}}_{\vec{\beta}_{\text{hat}}}) \vec{\beta}_{\text{hat}}$$

$$\begin{pmatrix} \hat{e}_1 \\ \vdots \\ \hat{e}_n \end{pmatrix} = \vec{Y} - \mathbf{X} \vec{\beta}_{\text{hat}}.$$

$$= \vec{Y} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y}$$

$$= (\mathbf{I} - H) \vec{Y}.$$

$$\text{cov}(\vec{\beta}_{\text{hat}}, \vec{Y} - \mathbf{X} \vec{\beta}_{\text{hat}})$$

$$= \text{cov}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{Y}, (\mathbf{I} - H) \vec{Y})$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{0}^2 \mathbf{I} (\mathbf{I} - H)$$

$$= 0 \quad (\underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{I} - H)}_{\mathbf{0}^2}).$$

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n \hat{e}_i^2 = \vec{Y}^T (\underbrace{\mathbf{I} - H}_{\text{rank } (\mathbf{I} - H)}) (\underbrace{\mathbf{I} - H}_{\text{rank } (\mathbf{I} - H)}) \vec{Y} \\ &= \vec{Y}^T (\mathbf{I} - H) \vec{Y} \sim \chi_{n-p-1}^2. \end{aligned}$$

$\text{rank } (\mathbf{I} - H)$

$= n - \text{rank } (\mathbf{I})$

$= n - p - 1.$

$$\text{Cl's of } \vec{\beta} \quad \left( \begin{array}{c} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{array} \right) = \vec{\beta}_{\text{nat}} \sim N(\vec{\beta}, \underbrace{\sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1}}_{\text{Var}})$$

$$\hat{\beta}_k \sim N(\hat{\beta}_k, \underbrace{\sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1}}_{\text{Var}})$$

**Theorem D'.** Under the MLR assumptions,

$$\frac{\hat{\beta}_k - \beta_k}{\text{se}(\hat{\beta}_k)} \sim t_{n-p-1}, \quad k = 0, 1, \dots, p.$$



95% exact confidence interval for  $\beta_k$  :

$$\hat{\beta}_k \pm t_{n-p-1}(\alpha/2) \cdot \text{se}(\hat{\beta}_k)$$



$$\text{RSS} = \vec{Y}^\top (I - H) \vec{Y}$$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n - p - 1}$$

$$\text{se}(\hat{\beta}_k) = \sqrt{\hat{\sigma}^2 [(\mathbf{X}^\top \mathbf{X})^{-1}]_{kk}}$$

# Cls of $\vec{\beta}$

**Example 3.** Interested in answering whether person's brain size and body size predictive of his or her intelligence, some researchers (Willerman, et al, 1991) collected the following data on a sample of  $n = 38$  college students:

- **Response ( $y$ ):** Performance IQ scores (PIQ) from the revised Wechsler Adult Intelligence Scale. Potential
- $x_1$ : Brain size based on the count obtained from MRI scans (given as count/10,000).
- $x_2$ : Height in inches.
- $x_3$ : Weight in pounds

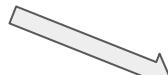
$$\vec{\beta}_{\text{hat}} = (\underline{\mathbf{X}^T \mathbf{X}})^{-1} \mathbf{X}^T \vec{\mathbf{Y}}$$

$$\vec{\mathbf{Y}} - \mathbf{X} \vec{\beta}_{\text{hat}}$$

```
coefs <- solve(t(X)%%X)%%t(X)%%y
Res <- y-X%*%coefs
RSS <- sum(Res^2)
s2_hat <- RSS/(n-ncol(X))
diag((X^T X)^{-1})
cov_diag_elements <- diag(solve(t(X)%%X))
se_beta_vec <- sqrt(s2_hat*cov_diag_elements)
alpha = 0.05
t_alpha <- qt(0.05/2, n-ncol(X), lower.tail = FALSE)
cbind(coefs - t_alpha*se_beta_vec, coefs + t_alpha*se_beta_vec)
```

→ confint(fit, level = 0.95)

$\text{fit} \sim \text{lmlRes} \sim \text{pred1} + \text{pred2} + \text{pred3}, \text{data} = \dots$



$$\begin{aligned} \text{RSS} &= \vec{\mathbf{Y}}^T (I - H) \vec{\mathbf{Y}} \\ \hat{\sigma}^2 &= \frac{\text{RSS}}{n - p - 1} \\ \text{se}(\hat{\beta}_k) &= \sqrt{\hat{\sigma}^2 [(\mathbf{X}^T \mathbf{X})^{-1}]_{kk}} \end{aligned}$$

	2.5 %	97.5 %
(Intercept)	-16.6190567	239.3262733
Brain	0.9153051	3.2054285
Height	-5.2304287	-0.2334296
Weight	-0.3999266	0.4010465

## SE of the regression line

$$E(\hat{y}) = E(\vec{x}^T \vec{\beta}_{\text{hat}}) = \vec{x}^T \vec{\beta}$$

$$\text{var}(\hat{y}) = \vec{x}^T \text{var}(\vec{\beta}_{\text{hat}}) \vec{x} = b^2 \vec{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \vec{x}$$

**Theorem E.** Under the MLR assumptions, we can use  $\vec{\beta}_{\text{hat}}$  to predict at any  $\vec{X} = \vec{x}$ :

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p = \vec{x}^T \vec{\beta}_{\text{hat}}.$$

Then  $E(\hat{y}) = \vec{x}^T \vec{\beta}$ ,  $\text{var}(\hat{y}) = \sigma^2 \vec{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \vec{x}$  and  $\hat{y}$  is independent of RSS.



95% exact confidence interval for  $\beta_0 + \cdots + \beta_p x_p$ :

$$\hat{y} \pm t_{n-p-1}(\alpha/2) \cdot \text{se}(\hat{y})$$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n - p - 1}$$

$$\text{se}(\hat{y}) = \sqrt{\hat{\sigma}^2 \vec{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \vec{x}}$$

For prediction intervals of the response:

$$E(y - \hat{y}) = 0, \text{var}(y - \hat{y}) = \sigma^2 \left[ 1 + \vec{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \vec{x} \right]$$

↑                    ↑

and  $y$  is independent of  $\hat{y}$  and RSS.

# CLs of population mean

**Example 3.** Interested in answering whether person's brain size and body size predictive of his or her intelligence, some researchers (Willerman, et al, 1991) collected the following data on a sample of  $n = 38$  college students:

- **Response ( $y$ ):** Performance IQ scores (PIQ) from the revised Wechsler Adult Intelligence Scale. Potential
- $x_1$ : Brain size based on the count obtained from MRI scans (given as count/10,000).
- $x_2$ : Height in inches.
- $x_3$ : Weight in pounds

Obtain a 95% CI for the mean IQ score when brain size = 96, height = 72, weight = 120.

```
> x_new <- c(1, 96, 72, 120) ←  $\hat{y} = \vec{x}^T \hat{\beta}_{\text{hat}}$ 
> fitted <- x_new %*% coefs ←
> se_fitted <- sqrt(s2_hat*t(x_new)%*%solve(t(X)%*%X)%*%x_new)
> c(fitted - t_alpha*se_fitted, fitted + t_alpha*se_fitted)
```

$\downarrow$        $\downarrow$        $\downarrow$

Lm()

```
> new <- data.frame(Intercept=1, Brain=96,
+                     Height=72, Weight=120)
> predict(fit, newdata = new, interval = 'confidence', level=0.95)
```

$\downarrow$        $\downarrow$        $\downarrow$

$\hat{y}$        $\hat{y} - t_{n-p}(\alpha_2) se(\hat{y})$        $\hat{y} + t_{n-p}(\alpha_2) se(\hat{y})$   
 $\downarrow$        $\downarrow$        $\downarrow$   
 fit      lwr      upr  
 112.5171 91.83446 133.1998

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n - p - 1}$$

$$se(\hat{y}) = \sqrt{\hat{\sigma}^2 \vec{x}^T (\vec{X}^T \vec{X})^{-1} \vec{x}}$$

$\uparrow$        $\uparrow$   
 $\downarrow$        $\downarrow$

```
predict(fit, newdata = new, interval = 'prediction', level=0.95)
fit      lwr      upr
112.5171 67.28453 157.7497
```

# Testing $H_0 : \beta_k = 0$

*Making sense of R output*

08/05/2021

$$\begin{aligned}\vec{y} &= \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{np} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \\ &= \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{np} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{k-1} \\ \underbrace{\beta_k}_{\text{n } x \text{ p}} \\ \beta_{k+1} \\ \vdots \\ \beta_p \end{pmatrix}\end{aligned}$$

take this  
↓  
out

# Hypothesis testing

$$\hat{\sigma}_{MLE}^2 = \frac{RSS}{n} \quad \hat{\sigma} = \frac{RSS}{n-p-1}$$

**Theorem G.** Under the MLR assumptions, find the LRT for  $H_0 : \beta_1 = 0$  versus  $H_1 : \beta_1 \neq 0$ .

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}_{n \times (p+1)}$$

Solution:  $\Rightarrow \mathbb{H}_0 = \{ \beta_0 \in \mathbb{R}, \beta_2 \in \mathbb{R}, \cdots, \beta_p \in \mathbb{R}, \beta^2 > 0 \}$

$$\Rightarrow \mathbb{H} = \{ \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}, \beta_2 \in \mathbb{R}^2, \cdots, \beta_p \in \mathbb{R}^p, \beta^2 > 0 \}$$

$$\sup_{\mathbb{H}_0} L(\hat{\beta}, \hat{\sigma}^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} e^{-\frac{(Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_p X_{ip})^2}{2\hat{\sigma}^2}}$$

$$\sup_{\mathbb{H}_0} L(\beta, \sigma^2) = \sup_{\mathbb{H}_0} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_p X_{ip})^2}{2\sigma^2}} \left( \frac{1}{\sqrt{2\pi\hat{\sigma}^2}} \right)^n e^{-\frac{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}}{2 \frac{1}{n} \vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}}}$$

$$X_0 = \begin{pmatrix} 1 & X_{12} & \cdots & X_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n2} & \cdots & X_{np} \end{pmatrix}_{n \times p}$$

$$RSS_0 = \vec{Y}^T (\mathbf{I} - \mathbf{H}_0) \vec{Y}$$

$$H_0 = X_0 (X_0^T X_0)^{-1} X_0^T$$

$$\lambda(\vec{Y}_n) = \frac{\sup_{\mathbb{H}_0} L}{\sup_{\mathbb{H}} L} = \left( \frac{RSS}{RSS_0} \right)^{\frac{n}{2}} = \left( \frac{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}}{\vec{Y}^T (\mathbf{I} - \mathbf{H}_0) \vec{Y}} \right)^{\frac{n}{2}} = \left( \frac{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}}{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y} + \vec{Y}^T (\mathbf{H} - \mathbf{H}_0) \vec{Y}} \right)^{\frac{n}{2}}$$

$$SS_E \downarrow$$

$$SS_E + SS_B \downarrow$$

$$R = \left\{ \lambda(\vec{Y}_h) \leq c \right\} = \left\{ \left[ \frac{\frac{1}{1 + \frac{\vec{Y}^T (H - H_0) \vec{Y}}{\vec{Y}^T (I - H) \vec{Y}}}}{\frac{\vec{Y}^T (H - H_0) \vec{Y}}{\vec{Y}^T (I - H) \vec{Y}}} \right]^{\frac{n}{2}} \leq c \right\} = \left\{ \frac{\vec{Y}^T (H - H_0) \vec{Y}}{\vec{Y}^T (I - H) \vec{Y}} \stackrel{\text{SS}_B}{\geq} c \right\}$$

$$\textcircled{1} \quad \vec{Y}^T (I - H) \vec{Y} \sim \chi_{n-p-1}^2.$$

$$\textcircled{2} \quad \vec{Y}^T (H - H_0) \vec{Y} \sim \chi_{\dim(\mathbb{B}) - \dim(\mathbb{B}_0)}^2 \text{ under } H_0$$

$$\begin{aligned} \text{rank}(H - H_0) &= \text{trace}(H - H_0) = \cancel{\text{trace } H} - \text{trace } H_0 = \text{rank } X - \text{rank } X_0 \\ &\stackrel{\text{II}}{=} 1. \end{aligned}$$

$F_{1, n-p-1}$

\textcircled{3} Independence :

It suffices to show  $(I - H) \vec{Y} \perp \!\!\! \perp (H - H_0) \vec{Y}$

$$\begin{aligned} \text{cov}((I - H) \vec{Y}, (H - H_0) \vec{Y}) &= (I - H)(H - H_0) \\ &= H - H_0 - H + \cancel{H H_0} \\ &= H - H_0 - H + H_0 \\ &= 0. \end{aligned}$$

Under  $H_0$ ,  $\frac{\vec{Y}^T (H - H_0) \vec{Y}}{\vec{Y}^T (I - H) \vec{Y}} / \frac{\text{SS}_E}{(\dim(\mathbb{B}) - \dim(\mathbb{B}_0))}$

$\sim F_{\dim(\mathbb{B}) - \dim(\mathbb{B}_0), n-p-1}$ .

$$t_{n-p-1}^2$$

$$H H_0 = H_0.$$

$$\cancel{H X_0} (X_0^T X_0)^{-1} X_0^T = X_0$$

projecting column vectors of  $X_0$  onto  $\text{span}(X)$ .  $\Rightarrow$  projecting onto the span that contains all column vectors of  $X_0$

# Hypothesis testing

**Theorem G'.** Now find the LRT for  $H_0 : \beta_1 = 0$  and  $\beta_2 = 0$  versus  $H_1 : H_0$  is not true.

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \ddots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}_{n \times (p+1)}$$

$$\mathbf{X}_0 = \begin{pmatrix} 1 & x_{13} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n3} & \cdots & x_{np} \end{pmatrix}_{n \times (p-1)}$$

$$\sup_{\{\boldsymbol{\theta}\}} L(\boldsymbol{\mu}, \boldsymbol{\theta}^2) = \left( \frac{1}{2\pi \cdot RSS} \right)^{\frac{n}{2}} e^{-\frac{n}{2}}$$

$$\sup_{\{\boldsymbol{\theta}_0\}} L(\boldsymbol{\mu}, \boldsymbol{\theta}_0^2) = \left( \frac{1}{2\pi \cdot RSS_0} \right)^{\frac{n}{2}} e^{-\frac{n}{2}}$$

$$RSS_0 = \vec{Y}^T (\mathbf{I} - \underbrace{\mathbf{H}_0}_{\downarrow}) \vec{Y}$$

$$R = \left\{ \frac{\vec{Y}^T (\mathbf{H} - \mathbf{H}_0) \vec{Y}}{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}} \geq c \right\}$$

$$\text{Under } H_0, \quad \frac{\vec{Y}^T (\mathbf{H} - \mathbf{H}_0) \vec{Y}}{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y}} / \frac{(d\text{im } \mathbf{H}) - (d\text{im } \mathbf{H}_0)}{(n-p-1)} \sim F_{2, n-p-1} \quad \text{11}$$

# Hypothesis testing

**Theorem G'.** Now find the LRT for  $H_0$  : a certain set of  $\beta'$ s are 0 versus  $H_1$  :  $H_0$  is not true.

$$\mathbf{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{pmatrix}$$

$$F \text{ statistic} = \frac{\vec{Y}^T (\mathbf{H} - \mathbf{H}_0) \vec{Y} / m}{\vec{Y}^T (\mathbf{I} - \mathbf{H}) \vec{Y} / (n - p - 1)} \stackrel{\substack{= \dim(\mathbb{H}) - \dim(\mathbb{H}_0) \\ = \# \text{ of zero } \beta's \text{ in } H_0}}{\sim} F_{m, n-p-1} \text{ under } H_0.$$

$$\mathbf{x}_0 = \begin{pmatrix} & \cdot & \\ & & \end{pmatrix}$$

$$H_0 = \mathbf{x}_0 (\mathbf{x}_0^T \mathbf{x}_0)^{-1} \mathbf{x}_0^T$$

# Output from R

```
> summary(fit)

Call:
lm(formula = PIQ ~ Brain + Height + Weight, data = dat) ←

Residuals:
    Min      1Q Median      3Q     Max 
-32.74 -12.09 -3.84 14.17 51.69 ←

Coefficients: ↓ ↓
              Estimate Std. Error t value Pr(>|t|) 
(Intercept) 1.114e+02 6.297e+01 1.768 0.085979 .
Brain        2.060e+00 5.634e-01 3.657 0.000856 *** 
Height       -2.732e+00 1.229e+00 -2.222 0.033034 * 
Weight        5.599e-04 1.971e-01  0.003 0.997750 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.79 on 34 degrees of freedom ←
Multiple R-squared:  0.2949   Adjusted R-squared:  0.2327 
F-statistic: 4.741 on 3 and 34 DF, p-value: 0.007215
```

$$X = \begin{pmatrix} 1 & X_{11} & X_{12} & X_{13} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & X_{n3} \\ \text{Brain} & \text{Height} & \text{Weight} \end{pmatrix}$$

$$\begin{aligned} H_0: \underline{\beta_{\text{brain}}} &= 0 \\ F_{\text{stat}} &\sim \underline{F_{1, n - \text{ncol}(X)}} \end{aligned}$$

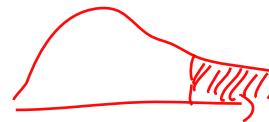
```
X <- cbind(1, dat$Brain, dat$Height, dat$Weight) ←
y <- dat$PIQ
H <- X %*% solve(t(X) %*% X) %*% t(X) ←

X0 <- X[, -2]
H0 <- X0 %*% solve(t(X0) %*% X0) %*% t(X0)
I <- diag(n)
numerator <- t(y) %*% (H - H0) %*% y / I
denominator <- t(y) %*% (I - H) %*% y / (n - ncol(X))
F_stat <- numerator / denominator ←

F_stat
[1] 13.37159
pf(F_stat, 1, n - ncol(X), lower.tail = FALSE)
[1] 0.0008556322
sqrt(F_stat)
[1] 3.656719
```

$$\hat{\sigma}^2 = \frac{RSS}{n-p-1} \sim \underline{n-p-1}$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$



## Output from R

```
> summary(fit)

Call:
lm(formula = PIQ ~ Brain + Height + Weight, data = dat)

Residuals:
    Min      1Q  Median      3Q     Max 
-32.74 -12.09 - 3.84 14.17 51.69 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 1.114e+02 6.297e+01   1.768 0.085979 .  
Brain        2.060e+00 5.634e-01   3.657 0.000856 *** 
Height       -2.732e+00 1.229e+00  -2.222 0.033034 *  
Weight        5.599e-04 1.971e-01   0.003 0.997750    
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.79 on 34 degrees of freedom
Multiple R-squared:  0.2949, Adjusted R-squared:  0.2327 
F-statistic: 4.741 on 3 and 34 DF,  p-value: 0.007215
```

$$X = \begin{pmatrix} 1 & | & | & | \\ | & | & | & | \\ 1 & | & | & | \end{pmatrix}$$

$$H_0: \beta_{\text{brain}} = \beta_{\text{height}} = \beta_{\text{weight}} = 0.$$

$$F_{\text{stat}} \sim F_{3, n - \text{ncol}(X)}$$

```
X0 <- X[, -c(2,3,4)] H0 <- X0 %*% solve(t(X0) %*% X0) %*% t(X0) I <- diag(n)
numerator <- t(y) %*% (H-H0) %*% y / 3
denominator <- t(y) %*% (I-H) %*% y / (n - ncol(X))
F_stat <- numerator / denominator

F_stat
[1] 4.740931 pf(F_stat, 3, n - ncol(X), lower.tail = FALSE)
[1] 0.00721527 C_{0.05}
```

Fail to reject  $H_0$

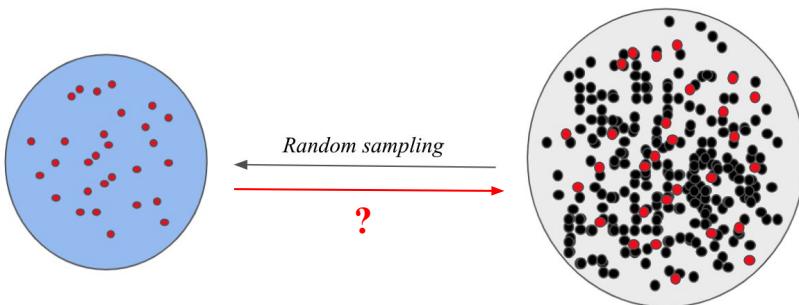
# Bayesian statistics

*8.6 of Rice*

08/05/2021

# Population vs. samples

How to learn from samples about  $f(x)$ ?



- **Point estimate:** a particular value  $\hat{\theta}$  that best approximates the parameter of interest.
  - Maximum likelihood estimators; Method of moment estimators;
  - Sufficiency; Estimation error quantification.
- **Interval estimate:** an interval  $[\theta - a, \theta + b]$  that would contain the true parameter  $\theta$  with a certain degree of confidence.
- **Hypothesis testing:** whether to reject a hypothesis.
  - $\theta > a?$
  - $\theta < a?$
- **Regression:** A special parametric model
$$Y | \beta, \sigma^2 = \underbrace{\mathbf{z}\beta}_{\text{Explanatory variables}} + N(0, \sigma^2)$$

Formulate the research problem

*Population  $f(x)$*

Make propositions about  $f(x)$

*Parametric  $f(x | \theta)$*

Survey design & Data collection

$X_1, \dots, X_n$

Parameter estimation of  $\theta$

Conclusion

# Prior information

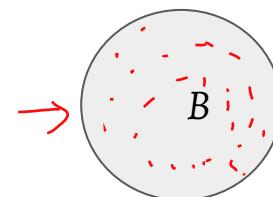
**Example 4.** What is a patient's probability of having liver disease if they are an alcoholic?

- Overall, 10% of patients have liver disease.  $P(A) = 0.10$ .
- Overall, 5% of patients are alcoholic.  $P(B)=0.05$ .
- Among those patients diagnosed with liver disease, 7% are alcoholics.  $P(B|A)=0.07$ .

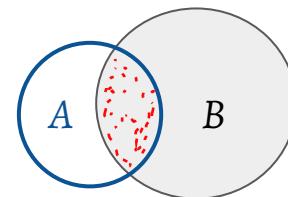
*from your experience*  
*experiment*

$$\text{Bayes rule: } P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

*Data information*      *Prior information*  
↑                          ↓  
*Posterior information*    *Marginalization*



Without  
prior



With prior

$$P(A|B) = \frac{0.07 \times 0.1}{0.05} = 0.14 > \underbrace{P(A) = 0.1}$$

# Prior information

**Example 5.** We want to estimate the average height of UCB students,  $\mu$ . Sample heights were collected from 100 students.

- **Model:**

$$x_j | \mu, \sigma^2 \sim N(\mu, \sigma^2), \quad j = 1, \dots, n.$$

- **Expert prior:** According to CDC.gov, average heights for US men and women ( $> 20$ ) are 69.0 inches and 63.5 inches.

Figure: 'Expert' prior  $\pi(\mu)$

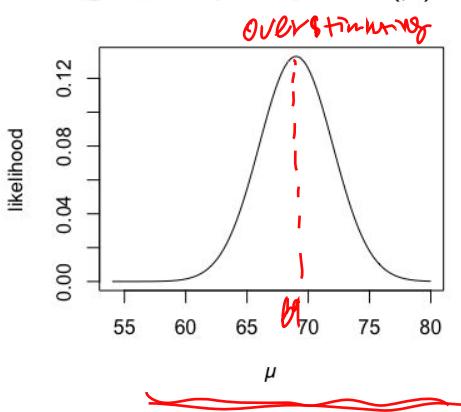
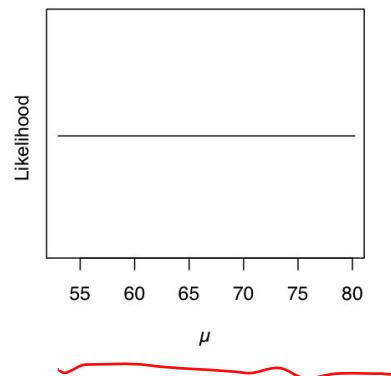


Figure: Naïve flat prior  $\pi(\mu)$



$$\pi(\mu, \sigma^2 | \mathbf{X}_n) = \frac{f(\mathbf{X}_n | \mu, \sigma^2) \pi(\mu) \pi(\sigma^2)}{f(\mathbf{X}_n)}.$$

- Is the frequentist prior a good idea?
- Is the expert prior correct for the UCB student population?

# Posterior distribution

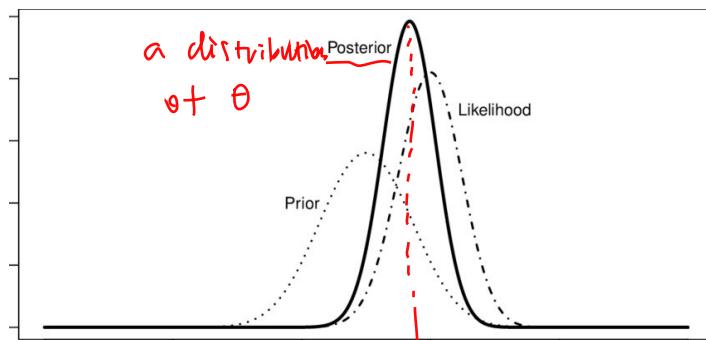
**Definition.** In the parametric model  $f(x | \theta)$  that we assume, we regard  $\theta$  as a random variable. Since  $\theta$  is random, we assign a pdf/pmf which captures our belief about  $\theta$ :

$\pi(\theta) \rightarrow$  prior distribution on  $\theta$ ,

Then my modified belief about  $\theta$  can be summarized in the posterior distribution

$$\pi(\theta | \mathbf{X}_n) = \frac{f(\mathbf{X}_n | \theta)\pi(\theta)}{\int_{\Theta} f(\mathbf{X}_n | \theta)\pi(\theta)d\theta} \propto f(\mathbf{X}_n | \theta)\pi(\theta).$$

⇒  $\hat{\theta}_B = E(\theta | \mathbf{X}_n)$  or median( $\theta | \mathbf{X}_n$ ) or mode( $\theta | \mathbf{X}_n$ ).



- When  $n$  is small, prior has greater impact;
- When  $n \rightarrow \infty$ , data information **dominates** the prior information.

$$X \sim \text{Beta}(\alpha, \beta), \quad E[X] = \frac{\alpha}{\alpha + \beta}, \quad \text{model } X = \frac{\alpha - 1}{\alpha + \beta - 2}.$$

## Posterior distribution

**Example 5.** You flip a coin 20 times and get 13 heads. Do you have a fair coin?

**Frequentist:**  $x_1, \dots, x_{20}$  iid Bernoulli( $p$ )  
 $f(\bar{x}_n | p) = p^{\sum_i x_i} (1-p)^{n - \sum_i x_i}$

$$\left. \begin{array}{l} \text{MLE: } \hat{p}_{\text{MLE}} = \bar{x}_n = \frac{13}{20} = 0.65 \\ \text{MM: } \end{array} \right\}$$

Estimators are based on  $\bar{x}_n$

$\hat{p}_{\text{MLE}} \rightarrow$  Sampling distribution

$$\hat{p}_{\text{Bayes}} = \text{model}(p | \bar{x}_n) = \frac{\sum_i x_i}{n} = \bar{x}_n$$

$$\text{or } \hat{p}_{\text{Bayes}} = E(p | \bar{x}_n) = \frac{\sum_i x_i + 1}{n + 2} = \frac{14}{22} = 0.63$$

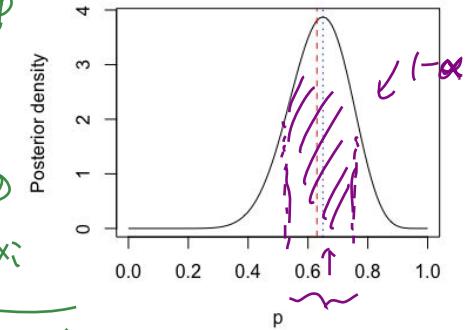
**Bayesian:**

$$\pi(p | \bar{x}_n) = \frac{f(\bar{x}_n | p) \pi(p)}{\int_0^1 f(\bar{x}_n | p) \pi(p) dp}$$

$$= \frac{p^{\sum_i x_i} (1-p)^{n - \sum_i x_i}}{\int_0^1 p^{\sum_i x_i} (1-p)^{n - \sum_i x_i} dp}$$

$$= \frac{p^{\sum_i x_i} (1-p)^{n - \sum_i x_i}}{\text{Beta}(\sum_i x_i + 1, n - \sum_i x_i + 1)}$$

$$\sim \text{Beta}\left(\sum_i x_i + 1, n - \sum_i x_i + 1\right)$$



We can easily quantify the uncertainty of the Bayesian estimators.

# Posterior distribution with informative prior

**Example 5.** You flip the same coin 20 times again and get 10 heads. What is your modified belief on  $p$ ?

Informative prior:  $\pi(p) = \text{Beta}(14, 8)$ .

$$\pi(p | \bar{x}_n) = \frac{f(\bar{x}_n | p) \pi(p)}{\int_0^1 f(\bar{x}_n | p) \pi(p) dp} = \frac{p^{\sum_i x_i} (1-p)^{n - \sum_i x_i} p^{13} (1-p)^7}{\int_0^1 p^{\sum_i x_i + 13} (1-p)^{n - \sum_i x_i + 7} dp}$$

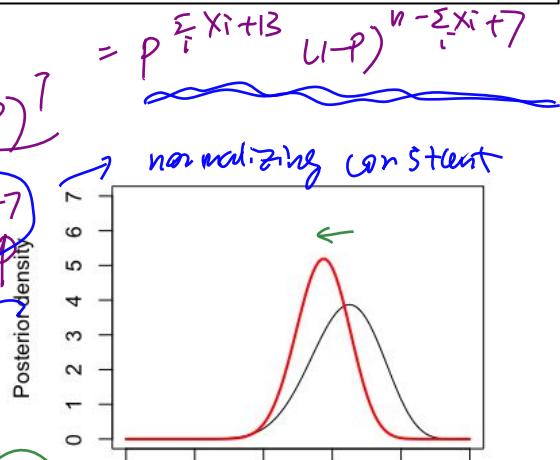
$$\sim \text{Beta}\left(\sum_i x_i + 14, n - \sum_i x_i + 8\right)$$

$$\hat{p}_{\text{Bayes}} = \frac{\sum_i x_i + 14}{n + 22} = \frac{n}{n+22} \cdot \frac{\sum_i x_i}{n} + \frac{22}{n+22} \cdot \frac{14}{22}$$

$n \rightarrow \infty$ ,  $\frac{n}{n+22} \rightarrow 1$ .

$$\hat{p}_{\text{Bayes}} = \bar{x}_n, \quad \frac{22}{n+22} \rightarrow 0$$

Frequentist estimator based on data information on the prior.



Our information on  $p$  is updated and more accurate.

$\frac{22}{n+22} \uparrow$

## Posterior distribution with informative prior

$$P(\theta \in [L(\bar{X}_n), U(\bar{X}_n)]) \leftarrow \frac{\bar{X}_n - Z_{\alpha/2} \frac{s}{\sqrt{n}}}{\bar{X}_n + Z_{\alpha/2} \frac{s}{\sqrt{n}}}$$

**Example 5.** You flip the same coin 20 times again and get 10 heads. What is your modified belief on  $p$ ?

Bayesian CI:

$$P(\theta \in I | \mathbf{X}_n) = 1 - \alpha.$$

Credible Interval

Bayesian HT:

Accept  $H_0$  if  $\underline{P(\theta \in \Theta_0 | \mathbf{X}_n)} \geq P(\theta \in \Theta_1 | \mathbf{X}_n)$ , else reject.

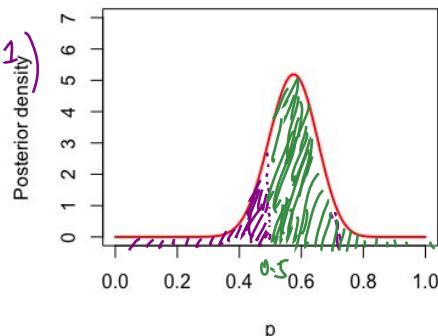
$$H_0: p \leq 0.5 \text{ vs. } H_1: p > 0.5$$

$$(0, 0.5]$$

$$(0.5, 1)$$

$$P(p \leq 0.5 | \mathbf{X}_n) =$$

$$P(p > 0.5 | \mathbf{X}_n) =$$

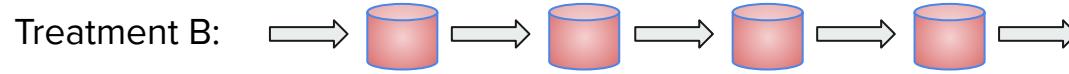
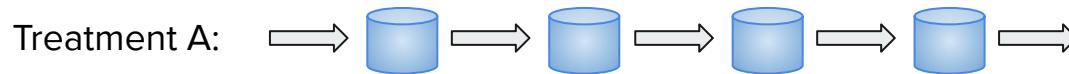


Reject  $H_0$  and conclude  $p > 0.5$ .

## Application: Bayesian sequential update

**Example 6.** In large-scale A/B testing, it is costly and time-consuming to collect large number of samples. What if you want to save money or want to know the test results as the experimentation goes on.

- Collect data chunks

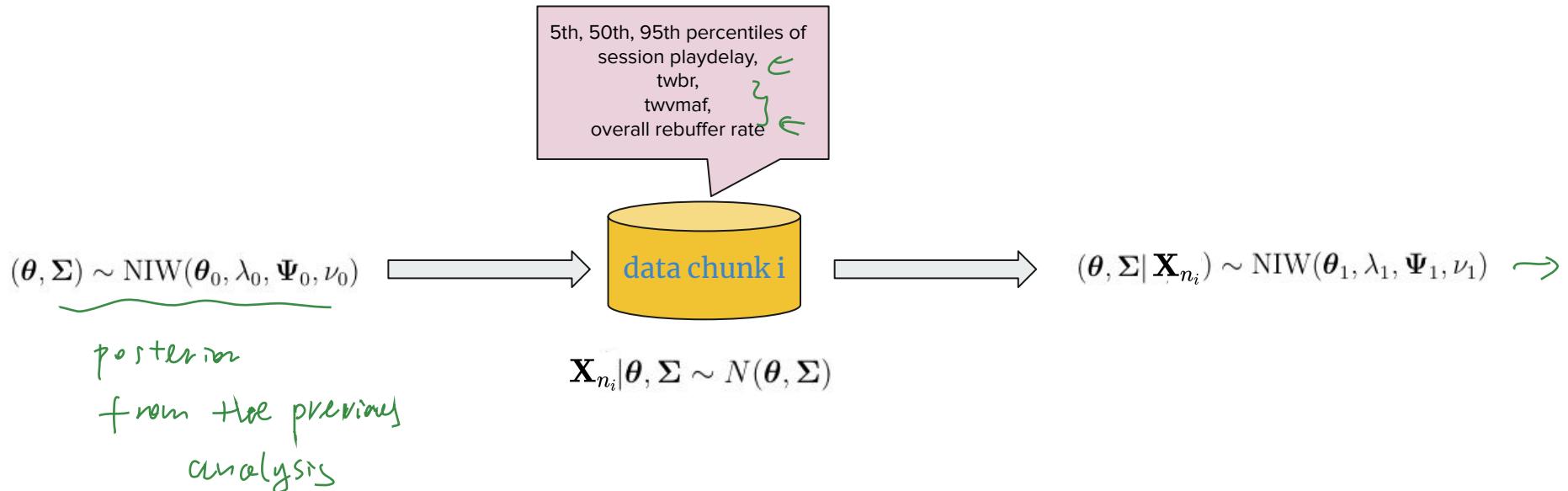


- Allow to fully adaptive experimentation



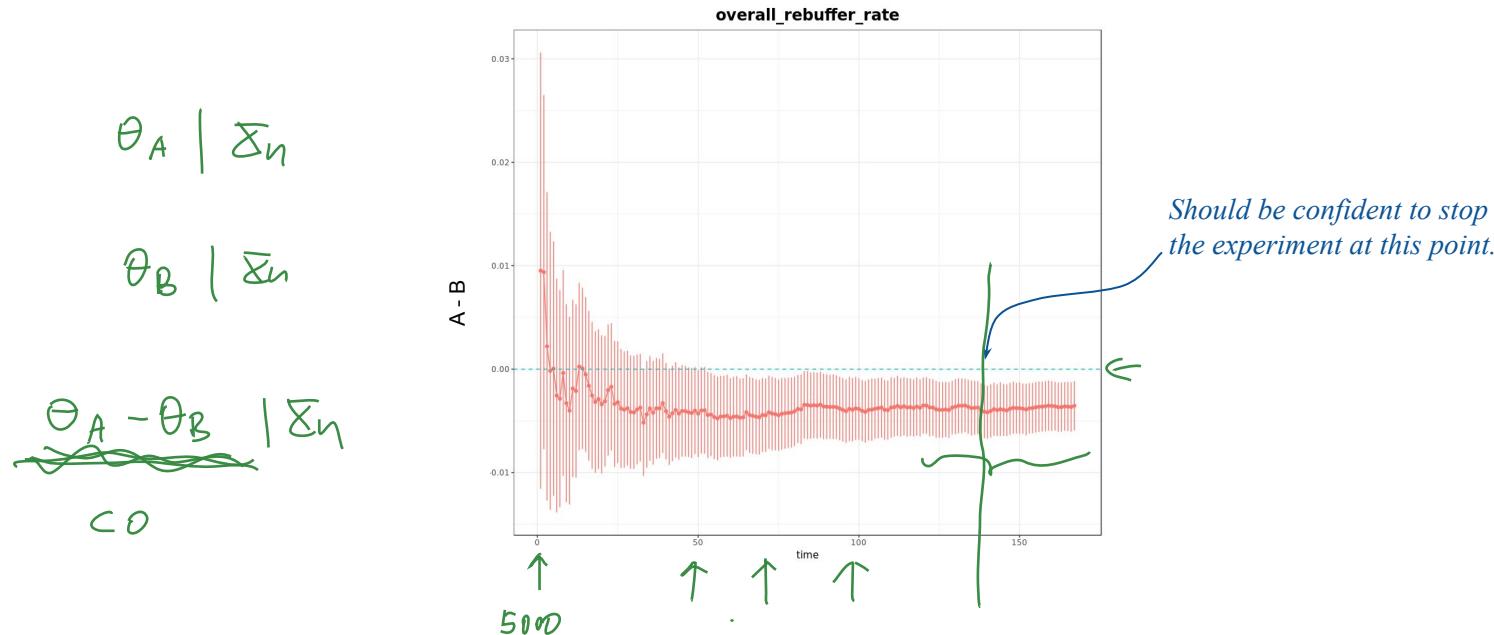
## Application: Bayesian sequential update

**Example 6.** In large-scale A/B testing, it is costly and time-consuming to collect large number of samples. What if you want to save money or want to know the test results as the experimentation goes on.



## Application: Bayesian sequential update

**Example 6.** In large-scale A/B testing, it is costly and time-consuming to collect large number of samples. What if you want to save money or want to know the test results as the experimentation goes on.



## Next week ...

- Review session & practice final;
- Final exam.