

$$H_0: p \leq .5$$

$$H_A: p > .5$$

$$\sum_{i=1}^{100} X_i \quad (= \# \text{ of heads in 100 tosses})$$

$$\sim \text{Bin}(100, p)$$

under  $H_0$

$$\begin{aligned} \text{Var}(\text{Bin}(n, p)) \\ = np(1-p) \end{aligned}$$

$$\sim \text{Bin}(100, p = .5)$$

$$\sum_{i=1}^{100} X_i \sim N(50, \sqrt{100(.5)(.5)})$$

"5

experiment: critical region

$$\sum X_i \geq 65$$

$\Rightarrow$  3 standard deviation cent.

level of the test  $P(\text{false positive is if } H_0 \text{ is true})$

type I error

$$\begin{aligned}
 P(\text{false positive}) &= P(Z \geq 3) \\
 &= 1 - \text{norm.cdf}(3) \\
 &\approx .0044
 \end{aligned}$$

let's specify

level:  $\alpha = .01$

$$\begin{aligned}
 \text{norm.ppf}(.99, \text{loc} = 50, \\
 \text{scale} = 5)
 \end{aligned}$$

$$\approx 61.6$$

→ critical region is

62 heads or more

. . .  $\alpha = .01$  test in  
use

56 heads → p-value is the probability  
of having an at least as  
extreme as this under  $H_0$   
what is the smallest level  
under which we would have  
rejected  $H_0$

$$P\left(\sum_{i=1}^n X_i \geq 56 \mid H_0\right) = 1 - \text{norm.cdf}(56, \text{loc} = 50, \text{scale} = 5)$$

2.11

Type II error  $\rightarrow$  when you fail to reject  $H_0$   
when some  $H_A$  is actually true

power of a test (given a particular  $H_A$  that we are testing)

$\hookrightarrow P(\text{test statistic being in the critical region} \mid H_A \text{ is true})$

$$p = 0.7$$

reject if  $\sum X_i \geq 62$

under this  $H_A$

$$\sum X_i \sim N(70, \sqrt{100 \cdot 0.7 \cdot 0.3})$$

$$P(\sum_{H_A} X_i \geq 62)$$

$$= 1 - \text{norm.cdf}(62, \text{loc} = 70, \text{scale} = np.sqrt(100 \cdot 0.7 \cdot 0.3))$$

2.96

$X_1, \dots, X_{100}$

$m$  some large #

generate  $m$  bootstrap samples

↳ take sum of each sample

if  $\alpha = .01$

if the 1st percentile of  
bootstrapped sum is bigger  
than  $S_0$ , then reject  $H_0$

## Regression (Linear Regression)

Let's suppose my data is paired

$(X_1, Y_1), \dots, (X_n, Y_n)$

$X_i$ 's are not random - fixed

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma)$$

$\varepsilon_i$  mutually  
independent

$$Y_i \sim N(\alpha + \beta X_i, \sigma)$$

we are going to estimate  $\alpha, \beta, \sigma$   
using MLE

$$L(\alpha, \beta, \sigma | (X_1, Y_1), \dots, (X_n, Y_n))$$

$$= \frac{1}{(2\pi)^{n/2}} \frac{1}{\sigma^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2\right)$$

$$\ell(\alpha, \beta, \sigma) = \underline{\underline{-\frac{n}{2} \ln(2\pi) - n \ln(\sigma)}}$$

$$- \frac{1}{2\sigma^2} \sum_{i=1}^n \underline{\underline{(y_i - \alpha - \beta x_i)^2}}$$

if I maximize

$$- \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

minimizing

$$\sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

↳ finding the least square line

$$\frac{\partial \ell}{\partial \alpha} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i) = 0$$

$$\frac{\partial}{\partial \alpha} = \frac{\partial}{\partial \alpha} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2$$

$$\sum (y_i - \alpha - \beta x_i) = 0$$

$$n \alpha = \sum y_i - \beta \sum x_i$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\frac{\partial \ell}{\partial \beta} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \beta x_i - \alpha) x_i = 0$$

$$(\text{multiply by } \sigma^2) = \sum_{i=1}^n y_i x_i - \alpha \sum x_i - \beta \sum x_i^2$$

$$= \sum_{i=1}^n y_i x_i - (\bar{y} - \beta \bar{x}) \sum x_i - \beta \sum x_i^2$$

$$= \cancel{n} \cdot \bar{y} \bar{x} - (\bar{y} - \beta \bar{x}) \cancel{n} \bar{x} - \beta \cancel{n} \bar{x}^2 = 0$$

$$\beta (\bar{x}^2 - \bar{x}^2) = \bar{y} \bar{x} - \bar{y} \bar{x}$$

$$\hat{\beta} = \frac{\bar{y} \bar{x} - \bar{y} \bar{x}}{\bar{x}^2 - \bar{x}^2}$$

$$= \frac{\sum x_i y_i - \bar{x} \sum y_i}{\sum x_i^2 - n \bar{x}^2} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} = \hat{\beta}$$

$\hat{\beta}$  is normal dist'n

$$E[\hat{\beta}] = \frac{\sum (x_i - \bar{x}) E[y_i]}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x}) (\alpha + \beta x_i)}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\alpha (\sum x_i - \bar{x}) + \beta \sum x_i (x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \quad (1)$$

$$\sum x_i (x_i - \bar{x}) = \sum x_i^2 - n \bar{x}^2$$

$$(1) \quad \beta$$

$$\text{Var}(\beta) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

Likewise, we can show that the dist'n of  $\hat{\alpha}$   
 $\sim N(\alpha, \dots)$

$$SS_x = \sum (x_i - \bar{x})^2$$

SSR = sum of squared residuals

$$= \sum_{i=1}^n \underbrace{(y_i - \hat{\alpha} - \hat{\beta} x_i)}_{\hat{\epsilon}_i}^2$$

$$\sigma^2 SSR \sim \chi^2_{n-2}$$

$$\hat{\beta} \pm \sqrt{\frac{SSR}{(n-2)SS_x}} t_{\alpha/2, n-2}$$

we can also bootstrap this

$(1, 3), (2, 5), (1, 4)$

$(x_1, y_1), \dots, (x_n, y_n)$

resample the pairs  $\rightarrow$  generate bootstrap samples

calculate regression coefficients  
for each sample

find appropriate sample



percentiles

## Prediction Intervals ...

$$X_1, \dots, X_n \sim N(\mu, \sigma)$$

$X_{n+1}$  comes in the future

let's say I know  $\mu, \sigma$

can show that for any prediction

$$E[(X_{n+1} - p)^2] \geq E[(X_{n+1} - \mu)^2]$$

$\mu$  is the squared error optimal predictor

$\Rightarrow$  prediction with  $\bar{X}$

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$Var(X_{n+1} - \bar{X}) = E[(X_{n+1} - \bar{X}_n)^2]$$

$$X_{n+1} \sim N(\mu, \sigma)$$

$$\bar{X}_n \sim N(\mu, \frac{\sigma}{\sqrt{n}})$$

$$Var(aX + bY)$$

$$= a^2 Var(X)$$

$$+ b^2 Var(Y)$$

$$+ 2ab Cov(X, Y)$$

$$Var(X_{n+1} - \bar{X}_n)$$

$$= \text{Var}(X_{n+1}) + \text{Var}(\bar{X}_n)$$

$$= \sigma^2 + \frac{\sigma^2}{n} \leftarrow$$

$$= \sigma \left( 1 + \frac{1}{n} \right)$$

$\Rightarrow$  prediction error: part of error  
 come from estimate  $E(X_{n+1})$   
 rest comes from  $\text{Var}(X_{n+1})$

$\Rightarrow$  prediction errors in regression  
 models

$(X_{n+1}, Y_{n+1})$

$\nearrow Y_{n+1}$  I don't  
 choose  $X_{n+1}$

$$Y_{n+1} \sim N(\underline{\alpha + \beta x_i}, \underline{\sigma})$$

$\Rightarrow$  calculate a prediction interval  
 for  $Y_{n+1}$ ?

use the bootstrap to generate prediction  
 intervals



