

Warmup 2:00 - 2:10

Stat 88 lec 35

German tanks were numbered 1, 2, 3, ..., N, with N unknown, during World War 2 and the Allies needed to estimate N. They captured 5 tanks numbered 20, 31, 43, 78 and 92. Can you find an unbiased estimate of N?

Using \bar{X}

$$E(\bar{X}) = \frac{N+1}{2} \Rightarrow 2E(\bar{X}) - 1 = N$$

$$T_1 = 2\bar{X} - 1 \rightarrow \text{unbiased}$$

$$\text{since } E(T_1) = 2E(\bar{X}) - 1 = N \quad \checkmark$$

$$T_1 = 2(52.8) - 1 = \boxed{104.6}$$

Using $M = \max(x_1, \dots, x_n)$

g.g.p. $\overbrace{\dots 20 \dots 31 \dots 43 \dots 78 \dots 92 \dots}^m \dots N$

expected length of a g.g.p. $= \frac{N-n}{n+1}$, $n=5$ sample size

$$E(M) = \frac{N-n}{n+1} \cdot n + n = n \left[\frac{N-n}{n+1} + 1 \right]$$

$$= n \left[\frac{N-n+nt}{n+1} \right] = \frac{n}{n+1} (N+1)$$

$$\Rightarrow \frac{n+1}{n} E(M) - 1 = N$$

$$T_2 = \frac{n+1}{n} M - 1 \rightarrow \text{unbiased estimator}$$

$$T_2 = \frac{6}{5} (92) - 1 = \boxed{109.4}$$

We need to calculate $\text{Var}(T_1)$ and $\text{Var}(T_2)$

to see which is the better estimator.

Last time

Sec 9.1, 9.2 Bias and variance

We score how good an estimator T of a parameter Θ is by $MSE_{\Theta}(T) = E((T-\Theta)^2)$.

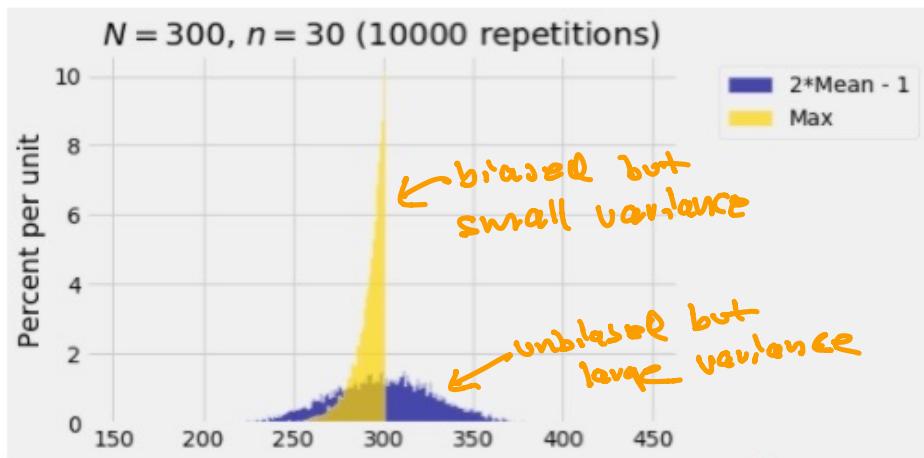
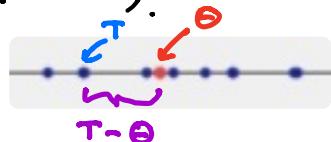
We showed,

$$MSE_{\Theta}(T) = (B_{\Theta}(T))^2 + \text{Var}_{\Theta}(T)$$

$$\text{where } B_{\Theta}(T) = E(T) - \Theta \quad (\text{bias})$$

$$\text{Var}_{\Theta}(T) = E((T-E(T))^2) \quad (\text{variance}),$$

The best estimator isn't always unbiased.

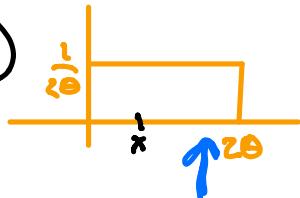


To find an unbiased estimator start with a statistic whose expectation is a linear function of the parameter.

- Today
- ① Sec 11.2 Practice finding an unbiased estimator
 - ② Sec 11.3 least squares regression.

① Sec 11.2 Practice finding an unbiased estimator

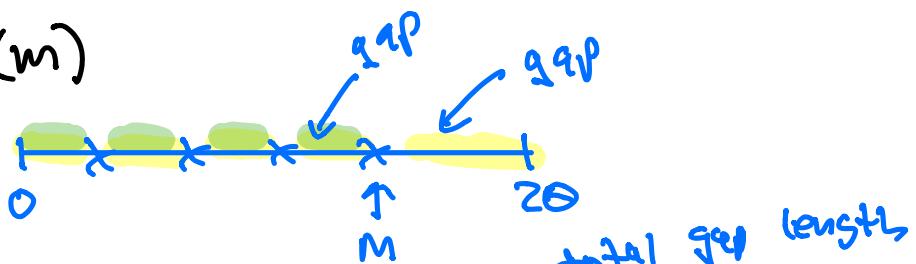
Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(0, 2\theta)$



let $M = \max(X_1, \dots, X_n)$

Is M a biased estimator of 2θ ?

Find $E(M)$



$$\text{average gap length} = \frac{2\theta}{n+1}$$

↑ # gaps

$$E(M) = \frac{2\theta}{n+1} \cdot n$$

discrete case.

$$E(M) = \frac{n-n}{n+1} \cdot n + n \\ = (n/n+1)(n+1)$$

Appendix: See appendix for another way to calculate $E(M)$ using the CDF of M .

Find an unbiased estimator for 2θ .

Solve for 2θ :

$$\frac{n+1}{n} E(M) = 2\theta$$

$$E(M) = \frac{2\theta}{n+1} \cdot n$$

$$\Rightarrow T = \frac{n+1}{n} M \quad \text{is unbiased}$$

② sec 11.3 least squares regression

Let (x, y) be a random pair of father, son heights from the population

x = father height

y = son height

We want to estimate y , call this \hat{y} , by the function $\hat{y} = ax + b$ for some a, b ,

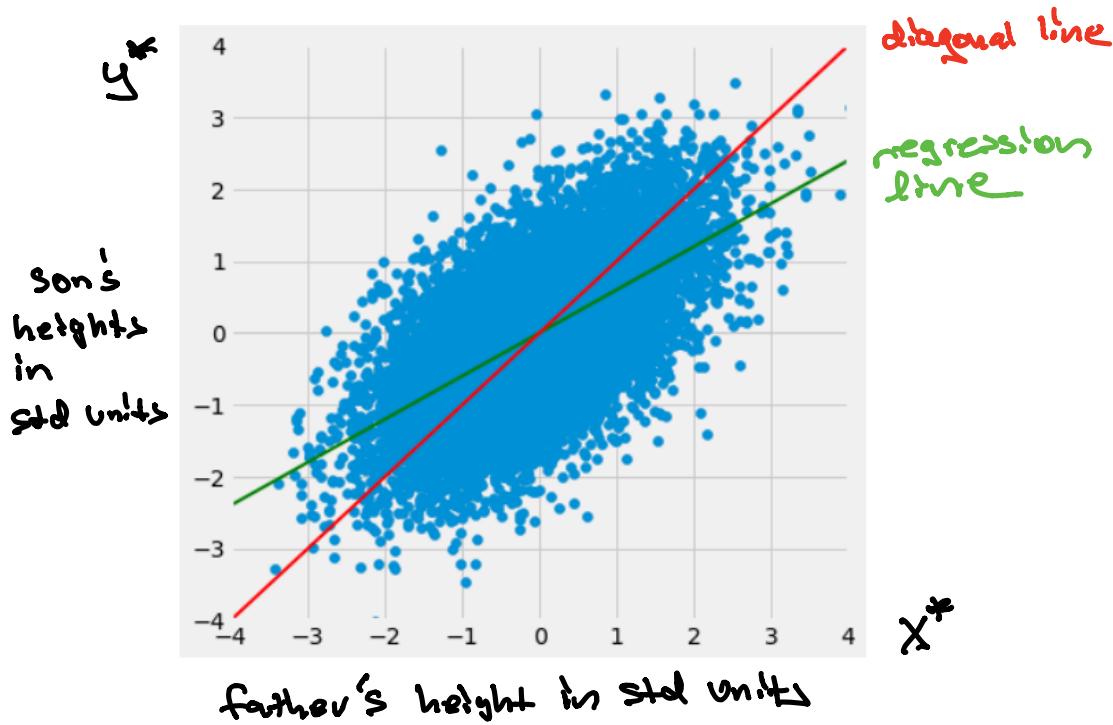
you plug in x into $\hat{y} = ax + b$ to predict y .

To find a and b ,
in date 8 you collected n pairs
 $(x_1, y_1), \dots, (x_n, y_n)$ and made a
scatter plot.

The regression line is the "best"
fitting line $\hat{y} = ax + b$
through your scatter plot.

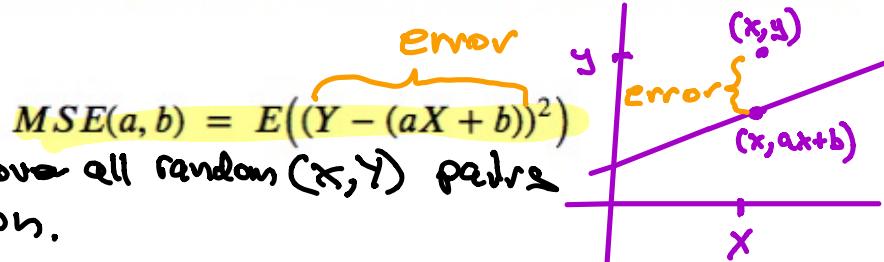
$$a = \text{slope of the regression line} = r \frac{\text{SD of } y}{\text{SD of } x}$$

$$\begin{aligned} b &= \text{intercept of the regression line} \\ &= (\text{average of } y) - \text{slope} \times (\text{average of } x) \end{aligned}$$



Mean Squared Error

For the random point (X, Y) , the mean squared error of a linear predictor of Y based on X depends on the slope a and intercept b of the line used. So let us define $MSE(a, b)$ to be the mean squared error when we use the line $aX + b$ to predict Y . That is,



We average over all random (X, Y) pairs in the population.

Notation

- $E(X) = \mu_X, SD(X) = \sigma_X$
- $E(Y) = \mu_Y, SD(Y) = \sigma_Y$

We need to take the partial derivative one variable at a time.

Step 1 Best intercept (b) for a fixed slope (a).

Fix slope, a , and solve $\frac{d}{db} \text{MSE}_a(b) = 0$ for b .

$$\begin{aligned}\text{MSE}_a(b) &= E\left((y - ax - b)^2\right) \\ &= E\left(\left((y - ax) - b\right)^2\right) = 0\end{aligned}$$

random variables
constant

FOIL

$$\begin{aligned}&= E\left((y - ax)^2 - 2(y - ax)b + b^2\right) \\ &= E(y^2) - 2E(y)E(x) + E(x^2)\end{aligned}$$

$E(b^2) = b^2$
since b not random.

Solve $\frac{d}{db} \text{MSE}_a(b) = 0$ for b

$$-2(E(y) - aE(x)) + 2b = 0$$

Solve for b

$$\hat{b}_a = E(y) - aE(x)$$

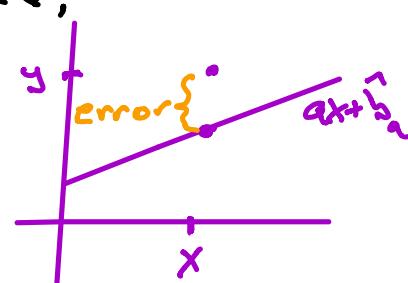
we will write $\hat{b}_a = \bar{y} - a\bar{x}$

Step 2 Find the best slope (a),

$$\text{error} = Y - \left(aX + \hat{b}_a \right)^{M_y - a\mu_x}$$

$$= Y - aX - M_y + a\mu_x$$

$$= \underbrace{Y - M_y}_{D_y} - a \underbrace{(X - \mu_x)}_{D_x}$$



remember
 $E(D_y^2) = \sigma_y^2$

$$\begin{aligned} MSE(a) &= E((D_Y - aD_X)^2) \\ &= E(D_Y^2) - 2aE(D_X D_Y) + a^2 E(D_X^2) \\ &= \sigma_Y^2 - 2aE(D_X D_Y) + a^2 \sigma_X^2 \end{aligned}$$

$$\frac{d}{da} MSE(a) = 0$$

$$\Rightarrow -2E(D_X D_Y) + 2a \sigma_X^2 = 0$$

$$\Rightarrow \hat{a} = \frac{E(D_X D_Y)}{\sigma_X^2}$$

$$= \frac{E((X - \mu_x)(Y - \mu_y))}{\sigma_X^2} \quad \begin{array}{l} \text{we will show equal} \\ \text{to } r \frac{\sigma_Y}{\sigma_X} \end{array}$$

So regression line is
 $\hat{y} = \hat{a}x + \hat{b}$ where $\hat{a} = \frac{E(D_x D_y)}{\sigma_x^2}$
 $\hat{b} = \mu_y - \hat{a}\mu_x$

Correlation

What is $E(D_x^2)$? $D_x = x - \mu_x$

$\longrightarrow \text{Var}(x)$

$E(D_x D_y)$ is the covariance of x and y ,

If x = father's ht. (ft)

y = son's ht. (ft)

$E(D_x D_y)$ has units ft^2

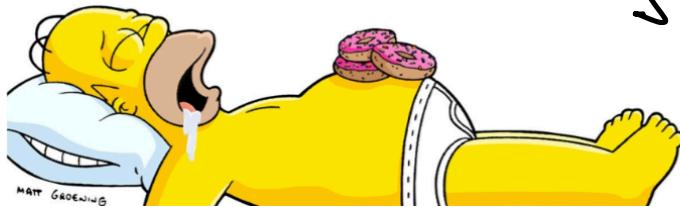
Dividing it by $\sigma_x \sigma_y$,

$r = \frac{E(D_x D_y)}{\sigma_x \sigma_y}$ is unitless and called the correlation coefficient of x, y ,

Covariance $E(D_x D_y) = r \sigma_x \sigma_y$

$$\text{so } \hat{a} = \frac{E(D_x D_y)}{\sigma_x^2} = \frac{r \sigma_x \sigma_y}{\sigma_x^2} = \boxed{\frac{r \sigma_y}{\sigma_x}}$$

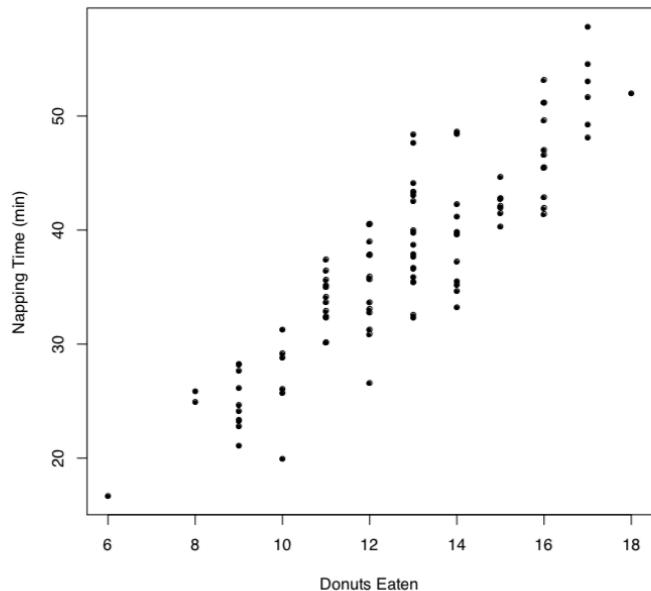
11F



$$\hat{y} = \hat{a}x + \hat{b}$$

$$\hat{a} = r \frac{s_y}{s_x}$$

$$\hat{b} = \bar{y} - \hat{a}\bar{x}$$



Suppose

$$\bar{x} = 11.5 \quad SD(x) = 2.5$$

$$\bar{y} = 30 \quad SD(y) = 8$$

$$r = .5$$

Suppose I tell you that Homer ate 14 donuts today. What is your best prediction for the time he spent napping?

$$\hat{y} = \hat{a}x + \hat{b}$$

$$\hat{a} = .5 \left(\frac{8}{2.5} \right) = 1.6$$

$$\hat{b} = 30 - 1.6(11.5) = 11.6$$

$$\text{so } \hat{y} = 1.6(14) + 11.6 = 34 \text{ min}$$

Appendix

Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unit}(0, 2\theta)$

Let $M = \max(X_1, \dots, X_n)$

Calculate the density of M by first calculating the cdf of M .

$$\begin{aligned} F(m) &= P(M \leq m) = P(X_1 \leq m, \dots, X_n \leq m) \\ &= P(X_1 \leq m)^n \\ &= \left(m \left(\frac{1}{2\theta}\right)\right)^n \end{aligned}$$

$$f(m) = \frac{d}{dm} F(m) = n \left(m \left(\frac{1}{2\theta}\right)\right)^{n-1} \cdot \frac{1}{2\theta}$$

$$= \boxed{n m^{n-1} \cdot \frac{1}{(2\theta)^n}}$$

Calculate $E(M)$

$$\begin{aligned} E(m) &= \int_0^{2\theta} m f(m) dm = \frac{n}{(2\theta)^n} \int_0^{2\theta} m^n dm \\ &= \frac{n}{(2\theta)^n} \frac{m^{n+1}}{n+1} \Big|_0^{2\theta} \\ &= \boxed{\left(2\theta\right) \frac{n}{n+1}} \end{aligned}$$