

* Announcement

① Quiz 10 : Thu (11/19) 9:00 AM

~ Fri (11/20) 9:00 AM

Ch10.3 ~ Ch11.1

Exponential Dist'n, $X \sim \text{Exp}(\lambda)$
 $\rightarrow f(x) = \lambda e^{-\lambda x}$ for $x \geq 0$. Density to find median
 $E(X) = \frac{1}{\lambda}$, $h = \frac{\ln 2}{\lambda}$ (Find h s.t. $F(h) = S(h)$)
 \uparrow half-life such that
 $= \text{median}$

Normal Dist'n.
 $X \sim N(\mu_x, \sigma_x^2)$, $Y \sim N(\mu_y, \sigma_y^2) \Rightarrow X+Y \sim ?$
 If indep, $X+Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$
 95% CI of $\mu_x - \mu_y$: $(\bar{x} - \bar{y} \pm 2 \sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}})$

STAT 88: Lecture 35

- Bias - Variance decomposition.

$$MSE = Bias^2 + Var.$$

You should be able to tell from figure if Bias/Var small or large



Contents

Section 11.2: The German Tank Problem, Revisited

Section 11.3: Least Squares Linear Regression

Warm up:

German tanks were numbered $1, 2, 3, \dots, N$, 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 ?

① $E(\bar{x}) = \frac{N+1}{2}$ ← population mean
 Unif $\{1, 2, \dots, N\}$

$$\Rightarrow 2E(\bar{x}) - 1 = N$$

$$\Rightarrow E(2\bar{x} - 1) = N$$

$$T_1 = 2\bar{x} - 1 = 2 \cdot \frac{20+31+43+78+92}{5} - 1 = 109.6$$

② $T_2 = \max\{X_1, \dots, X_n\}$



$$\text{Expected length of a gap} = \frac{N-5}{6} \rightarrow \frac{N-n}{n+1} \quad n \text{ Sample Size.}$$

$$E(\text{Gap } 6) = E(N - T_2) = \frac{N-n}{n+1}$$

$$\Rightarrow N - E(T_2) = \frac{N-n}{n+1}$$

$$\Rightarrow E(T_2) = N - \frac{N-n}{n+1} = \frac{n}{n+1} \cdot (N+1)$$

$$\Rightarrow \frac{n+1}{n} \cdot E(T_2) - 1 = N$$

$$\Rightarrow E\left(\frac{n+1}{n} \cdot T_2 - 1\right) = N$$

$$T_3 = \frac{6}{5} \cdot 92 - 1 = 109.4$$

$$\begin{aligned} Var(T_3) &= Var\left(\frac{n+1}{n} T_2 - 1\right) \\ &= \left(\frac{n+1}{n}\right)^2 Var(T_2) \\ &= \left(\frac{6}{5}\right)^2 \cdot Var(T_2) \\ &> 1 \end{aligned}$$

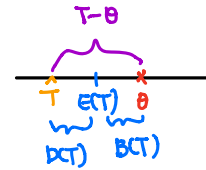
T_2 : biased estimator.

Last time

Bias and Variance

We score how good an estimator T of a parameter θ is by

$$\text{MSE}_\theta(T) = E_\theta((T - \theta)^2).$$



And we showed

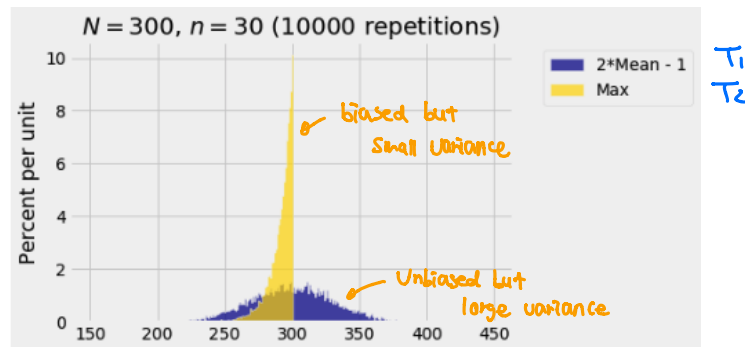
$$\text{MSE}_\theta(T) = B_\theta^2(T) + \text{Var}_\theta(T),$$

where

$$B_\theta(T) = E_\theta(T) - \theta \quad \text{and} \quad \text{Var}_\theta(T) = E_\theta((T - E_\theta(T))^2).$$

bias
Variance

The best estimator is *not* always unbiased.



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

11.2. The German Tank Problem

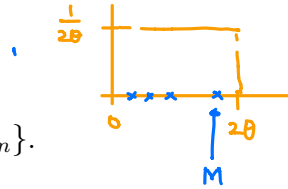
Practice for finding an unbiased estimator

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

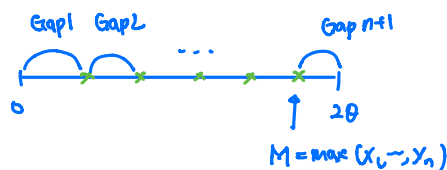
$$M = \max\{X_1, \dots, X_n\}.$$

Is M a biased estimator?

$$E(M) < 2\theta \Rightarrow M \text{ is biased.}$$



Find $E(M)$.



$$\begin{aligned} \text{Expected length of a Gap} &= \frac{2\theta - 0}{n+1} \leftarrow \text{total length} \\ &\quad \leftarrow \# \text{ Gaps} \end{aligned}$$

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

" $2\theta - E(M)$

Discrete: $\underbrace{1 \ 2 \ 3 \ \dots \ N}_{\text{Gap 1} \ \text{Gap 2} \ \text{Gap 3} \ \dots \ \text{Gap } n+1}$ $\frac{N-n}{n+1}$

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

(Discrete case)
 $E(T_2) = \frac{n}{n+1} (N+1)$

(See appendix for another way to calculate $E(M)$ using cdf of M)

Find an unbiased estimator for 2θ .

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

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

$$\Rightarrow E\left(\underbrace{\frac{n+1}{n} M}_T\right) = 2\theta$$

$$T = \frac{n+1}{n} \cdot M \quad \text{unbiased estimator.}$$

11.3. Least Squares Linear Regression

Let (X, Y) be a random pair of father and son heights from the population:

X : father height, and Y : son height.

We want to estimate Y , call this \hat{Y} , by the function

$$\hat{Y} = aX + b,$$

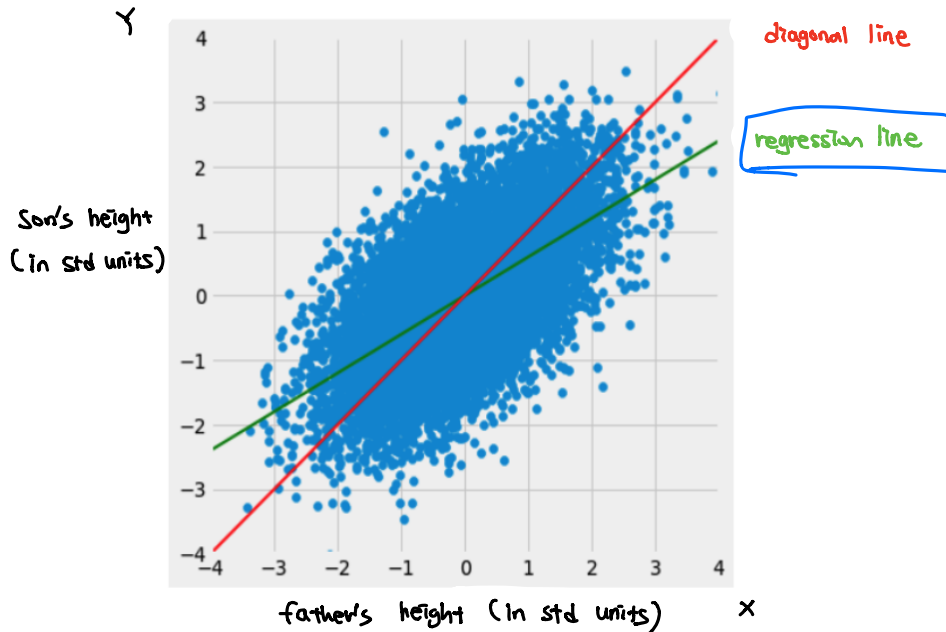
for some slope a and intercept b .

You plug in X into $\hat{Y} = aX + b$ to predict Y . To find a and b , in Data 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. The formulas are:

$$\text{slope of the regression line} = \underbrace{r}_{\text{Correlation}} \frac{\text{SD of } Y}{\text{SD of } X},$$

and

intercept of the regression line = (average of Y) – slope \times (average of X).

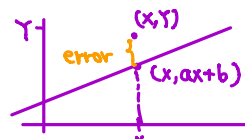


We will now derive the formulas mathematically using calculus and properties of expectation and variance.

$$T \sim \theta$$

$$MSE(T) = E((T - \theta)^2)$$

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,

$$MSE(a, b) = E(\overbrace{(Y - (aX + b))^2}^{\text{error}}).$$


Note that we average over all random (X, Y) pairs in the population. We have to find the values of a and b that minimize this function.

Notation

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

Best Intercept for a Fixed Slope Fix slope a , and solve $\frac{\partial MSE(a, b)}{\partial b} = 0$. Since

$$\begin{aligned} MSE(a, b) &= E((Y - (aX + b))^2) \\ &= E(((Y - aX) - b)^2) \\ &= E((Y - aX)^2 - 2b(Y - aX) + b^2) \quad \text{by additivity} \\ &= E((Y - aX)^2) - 2b \cdot E(Y - aX) + b^2. \end{aligned}$$

$E(b^2)$

Solve $\frac{\partial MSE(a, b)}{\partial b} = 0$ for b :

$$\begin{aligned} & -2E(Y - aX) + 2b = 0 \\ \Rightarrow & b = E(Y - aX) \\ & = E(Y) - aE(X) \\ & = \mu_Y - a\mu_X \\ & \quad \text{write it as } \hat{b}_a \end{aligned}$$

Best Slope For each fixed slope a , we first plug in the best intercept we just found.
The error becomes

$$\begin{aligned}
 \text{Error} &= Y - (aX + \hat{b}_a) = Y - (aX + \mu_Y - a\mu_X) \\
 &= Y - aX - \mu_Y + a\mu_X \\
 &= Y - \mu_Y - a(X - \mu_X) \\
 &= D_Y - aD_X.
 \end{aligned}$$

$\underbrace{\mu_Y}_{\text{Deviation of } Y}$ $\underbrace{X - \mu_X}_{\text{Deviation of } X}$

Then

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

$$E(D_Y^2) = E((Y - \mu_Y)^2) = \text{Var}(Y) = \sigma_Y^2$$

$$E(D_X^2) = \sigma_X^2$$

Solve $\frac{d\text{MSE}(a, \hat{b}_a)}{da} = 0$ for a :

$$\begin{aligned}
 0 - 2E(D_X D_Y) + 2a \cdot \sigma_X^2 &= 0 \\
 \Rightarrow a &= \frac{E(D_X D_Y)}{\sigma_X^2} = \frac{E((X - \mu_X)(Y - \mu_Y))}{\sigma_X^2}
 \end{aligned}$$

write it as \hat{a}

So the regression line is

$$\hat{Y} = \hat{a}X + \hat{b},$$

→ Best linear predictor
in the sense that it minimizes
 $\text{MSE} = E((Y - (\hat{a}X + \hat{b}))^2)$

where

$$\hat{a} = \frac{E(D_X D_Y)}{\sigma_X^2} \text{ and } \hat{b} = \mu_Y - \hat{a} \cdot \mu_X.$$

Correlation $E(D_X D_Y)$ is called the **covariance** of X and Y . If X is father's height (ft) and Y is son's height (ft), then $E(D_X D_Y)$ has unit ft^2 .

If we divide it by $\sigma_X \sigma_Y$,

$$r = \frac{E(D_X D_Y)}{\sigma_X \sigma_Y} \quad \frac{\cancel{\text{ft}^2}}{\cancel{\text{ft}^2}}$$

is unitless and called the correlation coefficient of X and Y . This tells you

$$\text{Covariance } E(D_X D_Y) = r \sigma_X \sigma_Y,$$

so

$$\hat{a} = \frac{E(D_X D_Y)}{\sigma_X^2} = \frac{r \cancel{\sigma_X} \sigma_Y}{\sigma_X^2} = \frac{r \sigma_Y}{\sigma_X}.$$

Appendix

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Unif}(0, 2\theta)$. Let $M = \max\{X_1, \dots, X_n\}$. Calculate the density of M by first calculating the CDF of M .

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

So,

$$f(m) = \frac{dF(m)}{dm} = nm^{n-1} \cdot \frac{1}{(2\theta)^n}.$$

Now we calculate

$$\begin{aligned} E(M) &= \int_0^{2\theta} mf(m)dm \\ &= \frac{n}{(2\theta)^n} \int_0^{2\theta} m^n dm \\ &= \frac{n}{(2\theta)^n} \frac{m^{n+1}}{n+1} \Big|_0^{2\theta} \\ &= (2\theta) \frac{n}{n+1}. \end{aligned}$$