

Lecture Note 2: Means, t-tests, and Regressions

Sample average: $\hat{\mu}_x = \bar{X} = \frac{1}{N} \sum_{i=1}^N X_i = \frac{1}{N} (X_1 + X_2 + \dots + X_N)$

Properties:

① BLUE: \bar{X} is the best linear unbiased estimator

linear: $\hat{\mu}_x = a_1 X_1 + a_2 X_2 + \dots + a_N X_N$ } $\varepsilon; a_i = 1$ } optimal $a_i = \frac{1}{N}$
unbiased: $E[\hat{\mu}_x] = \mu_x$
best: $\min V[\hat{\mu}_x]$

② Law of large numbers: $\bar{X} \xrightarrow{P} \mu_x$ (consistent)

③ CLT: as $N \rightarrow \infty$, $\bar{X} \sim N(\mu_x, \frac{\sigma_x^2}{N})$

so in "large" samples, $\frac{\sqrt{N}}{\sigma_x} (\bar{X} - \mu_x) \sim N(0, 1)$

\bar{X} is unbiased:

$$E[\bar{X}] = E\left[\frac{1}{N} \sum_{i=1}^N X_i\right] = \frac{1}{N} E[X_1 + X_2 + \dots + X_N]$$

$$= \frac{1}{N} \sum_{i=1}^N E[X_i] \quad \text{red arrow to } \boxed{\text{i.i.d.}}$$

$$= \frac{1}{N} \sum_{i=1}^N \mu_X$$

$$= \frac{1}{N} N \mu_X$$

unbiased!

Variance of \bar{X} :

$$\begin{aligned} V[\bar{X}] &= V\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} V[X_1 + X_2 + \dots + X_n] \quad \text{① i.i.d} \\ &= \frac{1}{n^2} \sum_{i=1}^n V[X_i] \quad \text{② i.i.d} \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma_X^2 \\ &= \frac{1}{n^2} n \sigma_X^2 \\ &= \frac{\sigma_X^2}{n} \end{aligned}$$

$$\begin{aligned} V[aX] &= a^2 V[X] \\ V[aX + bY] &= a^2 V[X] + b^2 V[Y] \\ &\quad + 2ab \operatorname{cov}(X, Y) \end{aligned}$$

Sample Variance

How to estimate $V(X) = E[(X - E(X))^2]$

How about $\hat{\sigma}_X^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$ ①

$$E[\hat{\sigma}_X^2] = \sigma_X^2 \underbrace{\left(1 - \frac{1}{N}\right)}_{\text{bias}}$$

For unbiased, need $\frac{1}{N-1}$ instead of $\frac{1}{N}$

② \uparrow
degrees of freedom

{ Estimator ① is biased but consistent and efficient
Estimator ② is unbiased and consistent but inefficient
tradeoff bet. bias and variance

t-test

$$t = \frac{\hat{\delta} - \theta_0}{SE(\hat{\delta})} \quad \text{null}$$

\Rightarrow CLT: as $N \rightarrow \infty$, $t \sim N(0, 1)$

One sample: use \bar{X} to test $\mu_X = 0$

$$V[\bar{X}] = \frac{\sigma_X^2}{N} \rightarrow SE[\bar{X}] = \frac{\sigma_X}{\sqrt{N}} \rightarrow t = \frac{\bar{X}}{\hat{\sigma}_X / \sqrt{N}}$$

Two sample: $\bar{X}_W, \bar{X}_B, \mu_W, \mu_B, \sigma_W^2, \sigma_B^2$ $\bar{X}_W \sim N(\mu_W, \frac{\sigma_W^2}{N_W})$

$$\bar{X}_B \sim N(\mu_B, \frac{\sigma_B^2}{N_B})$$

$$t = \frac{\bar{X}_W - \bar{X}_B}{SE[\bar{X}_W - \bar{X}_B]}$$

$$SE = \sqrt{V[\bar{X}_W - \bar{X}_B]} = \sqrt{V[\bar{X}_W] + V[\bar{X}_B] - 2\text{cov}(\bar{X}_W, \bar{X}_B)}$$

$$= \frac{\bar{X}_W - \bar{X}_B}{\sqrt{\frac{\sigma_W^2}{N_W} + \frac{\sigma_B^2}{N_B}}} = \frac{\bar{X}_W - \bar{X}_B}{\sqrt{SE_W^2 + SE_B^2}}$$

OLS Estimator

$$X, Y, \hat{Y} = b_0 + b_1 X \quad \text{"best fit"}$$

Min mean sq error \rightarrow error

$$U = Y - \hat{Y} = Y - b_0 - b_1 X$$

$$\min_{b_0, b_1} E[(Y - b_0 - b_1 X)^2]$$

population

$$\min_{\hat{b}_0, \hat{b}_1} \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{b}_0 - \hat{b}_1 X_i)^2$$

sample

In sample, $U_i = Y_i - \hat{Y}_i$ is called residual

$$\text{Optima: } \beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)}, \quad \hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

Can view $\hat{\beta}_1$ as an estimator for β_1 in the model:

$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

$$\hat{\beta}_1 \Rightarrow \beta_1$$

Gauss-Markov Th^m

- fixed (non-random) X_i
- Assumptions about $Y_i = \beta_0 + \beta_1 X_i + U_i$

① $E[U_i] = 0$ for all i

② $V[U_i] = \sigma^2$ for all i

③ $\text{cov}(U_i, U_j) = 0$ for $i \neq j$

Theorem: OLS is BLUE.