

STA302/1001 - Methods of Data Analysis I

(Week 01 - lecture B)

Wei (Becky) Lin

May 15-19, 2017



From last lecture (Week 1- Lect A)

- Notes about syllabus.
- A functional relationship vs a statistical relationship.
- Simple linear regression: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$. $\rightarrow \beta_0, \beta_1, \sigma^2$: 3 parameters
- Least Square Estimation of β_0, β_1 .
 - Gauss-Markov conditions
 - LS estimators b_0, b_1

$$b_0 = \bar{y} - b_1 \bar{x}, \quad b_1 = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_1^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_{xx}}$$

- The unbiased estimator for σ^2 .

$$s^2 = \hat{\sigma}^2 = MSE = \frac{SSE}{n-2}, \quad SSE = \sum_1^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n \ell_i^2$$

- Variance of LS estimators: $V(b_0)$ and $V(b_1)$

$$V(b_0) = V\left(\sum w_i Y_i\right) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{s_{xx}}\right)$$

$$V(b_1) = V\left(\sum k_i Y_i\right) = \frac{\sigma^2}{s_{xx}}$$

- Estimator of $V(b_0)$ and $V(b_1)$: replace unknown σ^2 by MSE.

Last lecture (Week 1- Lect A)

$$Q = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - b_0 - b_1 X_i)^2$$

- Show LS estimators b_0, b_1 which minimize Q are BLUE
 - LS estimators b_0, b_1 are linear estimators

$$b_0 = \sum_1^n w_i Y_i = \sum_1^n \left(\frac{1}{n} - k_i \bar{X} \right) Y_i = \sum_1^n \left\{ \frac{1}{n} - \frac{(X_i - \bar{X}) \bar{X}}{\sum_1^n (X_i - \bar{X})^2} \right\} Y_i$$

$$b_1 = \sum_1^n k_i Y_i = \sum_1^n \frac{X_i - \bar{X}}{\sum_1^n (X_i - \bar{X})^2} Y_i$$

and

$$\sum k_i = 0, \quad \sum k_i X_i = 1, \quad \sum k_i^2 = 1/S_{xx}$$

②

$$\sum w_i = 1, \quad \sum w_i X_i = 0$$

- Unbiasedness: $E(b_0) = \beta_0, E(b_1) = \beta_1$
- Show that b_0, b_1 have the smallest variance among the class of all linear unbiased estimators by linear algebra.

Week 01- Lecture B: Learning objectives & Outcomes

- Properties of Least Squares Fitted Line. same as "estimated reg. line"
- Why squared residuals instead of absolute residuals in OLS estimation?
- Introduction to other types of regressions
- Normal regression model.
- Maximum Likelihood Estimation (MLE) of β_0, β_1 .
- Review of statistical inference.

Properties of LS fitted line (Summary)

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

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

0. $E(\hat{Y}_i) = E(Y_i)$

1. $\sum_1^n e_i = 0$.

2. $\sum_1^n e_i^2$ is a minimum.

3. $\sum_1^n Y_i = \sum_1^n \hat{Y}_i$

4. $\sum_1^n X_i e_i = 0$

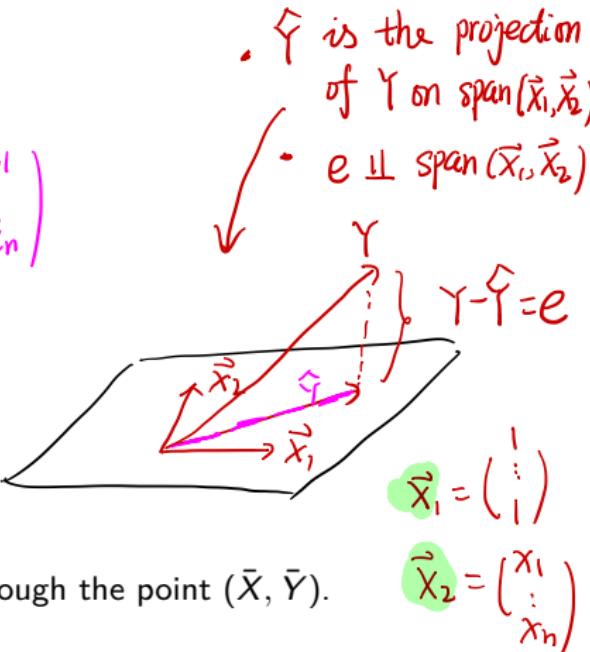
5. $\sum_1^n \hat{Y}_i e_i = 0$

6. The fitted regression line always goes through the point (\bar{X}, \bar{Y}) .

We will prove them one by one.

$$\begin{matrix} \uparrow \\ \vec{X}_1 \\ \uparrow \\ \vec{X}_2 \end{matrix}$$

$$e_i = Y_i - \hat{Y}_i$$



• two vectors \vec{u}, \vec{v} are orthogonal iff $\langle u, v \rangle = \sum_i^n u_i v_i = 0$

• e is perpendicular to the subspace spanned by \vec{X}_1, \vec{X}_2 , so

$$\langle e, \vec{X}_1 \rangle = \sum_{i=1}^n e_i \cdot 1 = 0, \quad \langle e, \vec{X}_2 \rangle = \sum_{i=1}^n X_i e_i = 0, \quad \langle e, \hat{Y} \rangle = 0$$

Properties of LS fitted line: (0) $E(\hat{Y}_i) = E(Y_i)$

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \rightarrow \begin{aligned} E(\varepsilon_i) &= 0 \\ V(\varepsilon_i) &= \sigma^2 \\ \text{cov}(\varepsilon_i, \varepsilon_j) &= 0 \end{aligned}$$

Proof:

$$E(Y_i) = \beta_0 + \beta_1 X_i + E(\varepsilon_i) = \underbrace{\beta_0 + \beta_1 X_i}_{\text{||}}$$

Since $\hat{Y}_i = b_0 + b_1 X_i$ so we have

$$E(\hat{Y}_i) = E(b_0) + E(b_1)X_i = \underbrace{\beta_0 + \beta_1 X_i}_{\text{||}}$$

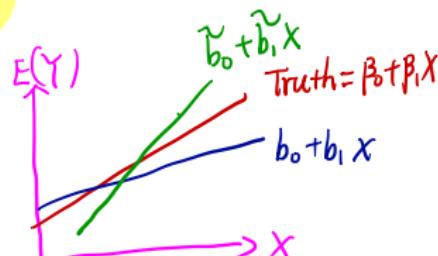
And b_0 and b_1 are unbiased estimators of β_0, β_1 respectively, it follows that

$$E(\hat{Y}_i) = \beta_0 + \beta_1 X_i = E(Y_i)$$

Q: How to understand $E(\hat{Y}_i) = E(Y_i)$

A: Given different samples of (\vec{X}, \vec{Y}) , we have different estimated lines: $\hat{b}_0 + \hat{b}_1 X$

But on avg, it leads to the truth reg. line



Properties of LS fitted line: (1) $\sum_1^n e_i = 0 \rightarrow$ Proof 2: $e_i = Y_i - \hat{Y}_i$

$$= Y_i - (\bar{Y} + b_1 X_i - \bar{X})$$

Proof:

- Recall b_0, b_1 satisfy the Normal equations:

$$\left\{ \begin{array}{l} \sum_{i=1}^n Y_i = nb_0 + b_1 \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i Y_i = b_0 \sum_{i=1}^n X_i + b_1 \sum_{i=1}^n X_i^2 \end{array} \right.$$

$$\begin{aligned} \sum_1^n e_i &= \sum_{i=1}^n (Y_i - \bar{Y} - b_1(X_i - \bar{X})) \\ &= \sum_1^n Y_i - n\bar{Y} - b_1 \sum (X_i - \bar{X}) \\ &= n\bar{Y} - n\bar{Y} - b_1 \cdot 0 \\ &= 0 \end{aligned}$$

- $\sum_1^n e_i = \sum_1^n (Y_i - \hat{Y}_i) = \sum_1^n (Y_i - b_0 - b_1 X_i)$ This gives

$$\sum_1^n e_i = \sum_{i=1}^n Y_i - nb_0 - b_1 \sum_{i=1}^n X_i$$

$$\begin{aligned} \hat{Y}_i &= b_0 + b_1 X_i \\ b_0 &= \bar{Y} - b_1 \bar{X} \\ \Rightarrow e_i &= \bar{Y} + b_1 (X_i - \bar{X}) \end{aligned}$$

- Thus we have the sum of residuals is zero, $\sum_1^n e_i = 0$, by first Normal equation.

Properties of LS fitted line: (2) $\sum_1^n e_i^2$ is a minimum



$$\hat{Y}_i = b_0 + b_1 X_i$$

where b_0, b_1 minimize $Q = \sum_{i=1}^n e_i^2 \Rightarrow Q$ reaches its minimum at b_0, b_1 .

Proof:

This was the requirement to be satisfied in deriving the LS estimators of the regression parameters since the criterion

$$Q = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - b_0 - b_1 X_i)^2$$

to be minimized when the LS estimators b_0 and b_1 are used for estimating β_0 and β_1 .

LS approach: Q is minimized at b_0, b_1 ,

$$\left\{ \begin{array}{l} \text{step 1: } \frac{\partial Q}{\partial \theta} = 0 \rightarrow \hat{\theta}, \quad \theta = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} \\ \text{step 2: check } \det(H) > 0 \text{ and } H(1,1) > 0 \end{array} \right.$$

$$\text{Properties of LS fitted line: (3)} \sum_1^n Y_i = \sum_1^n \hat{Y}_i \quad \leftarrow \text{(1)} \sum_{i=1}^n e_i = \sum (Y_i - \hat{Y}_i) = 0$$

↑
 observed ↑
 Fitted value

Proof:

Want to show that the sum of the **observed values** Y_i equals the sum of the **fitted values** \hat{Y}_i . From (1) we have $\sum_1^n e_i = \sum_1^n (Y_i - \hat{Y}_i) = 0$, so (3) follows.

$$\begin{aligned}
 \sum_1^n \hat{Y}_i &= \sum_1^n (b_0 + b_1 X_i) &> \boxed{\hat{Y}_i = b_0 + b_1 X_i = \bar{Y} + b_1 (X_i - \bar{X})} \\
 &= \sum_1^n (\bar{Y} - b_1 \bar{X} + b_1 X_i), \quad b_0 = \bar{Y} - b_1 \bar{X} &\Downarrow \\
 &= n\bar{Y} - b_1 n\bar{X} + b_1 \sum_1^n X_i &= \sum (\bar{Y} + b_1 (X_i - \bar{X})) \\
 &= n\bar{Y} - b_1 n\bar{X} + b_1 n\bar{X} &= n\bar{Y} + b_1 \sum_1^n (X_i - \bar{X}) \\
 &= n\bar{Y} &= n\bar{Y} \\
 &= \boxed{\sum_{i=1}^n Y_i} &= \sum Y_i \quad QED.
 \end{aligned}$$

Properties of LS fitted line: (4) $\sum_1^n X_i e_i = 0$

Proof:

Want to prove that the sum of the weighted residuals is zero when the i^{th} residual is weighted by the i^{th} predictor variable value. Again, the **Normal equations for b_0 and b_1** are

$$\sum_{i=1}^n Y_i = nb_0 + b_1 \sum_{i=1}^n X_i, \quad \left\{ \sum_{i=1}^n X_i Y_i = b_0 \sum_{i=1}^n X_i + b_1 \sum_{i=1}^n X_i^2 \right.$$

And we have

$$\begin{aligned} \sum_{i=1}^n X_i e_i &= \sum_{i=1}^n X_i (Y_i - \hat{Y}_i) = \sum_{i=1}^n X_i (Y_i - b_0 - b_1 X_i) \\ &= \left\{ \sum_{i=1}^n X_i Y_i - b_0 \sum_{i=1}^n X_i - b_1 \sum_{i=1}^n X_i^2 \right\} \end{aligned}$$

$$= 0$$

Thus, $\sum_1^n X_i e_i = 0$ by the second Normal Equation.

Properties of LS fitted line: (5) $\sum_1^n \hat{Y}_i e_i = 0$

Proof:

Want to prove that the sum of weighted residuals is zero when the weight of i^{th} residual is \hat{Y}_i .

$$\begin{aligned}\sum_1^n \hat{Y}_i e_i &= \sum_1^n (b_0 + b_1 X_i) e_i, \quad \hat{Y}_i = b_0 + b_1 X_i \\ &= b_0 \sum_1^n e_i + b_1 \sum_1^n X_i e_i \\ &= 0, \text{ by properties (1) and (4)}\end{aligned}$$

Properties of LS fitted line: (6) always goes thru (\bar{X}, \bar{Y})

Proof:

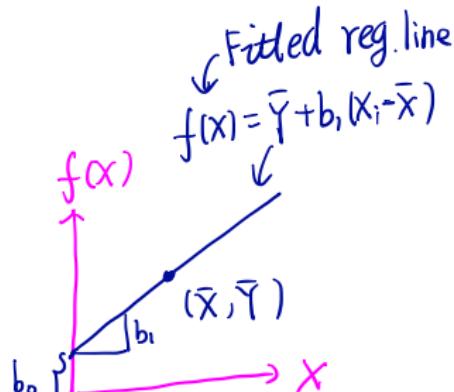
Recall the fitted regression line is

$$\hat{Y}_i = b_0 + b_1 X_i.$$

Replacing b_0 by $\bar{Y} - b_1 \bar{X}$ gives

$$\hat{Y}_i = \bar{Y} + b_1(X_i - \bar{X}),$$

When $X = \bar{X}$, we have $\hat{Y}_i = \bar{Y}$. That is, the fitted regression line always goes through the point (\bar{X}, \bar{Y}) .

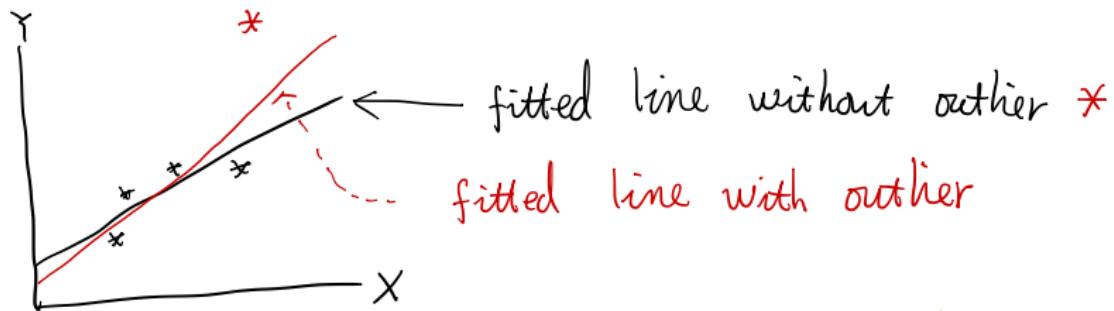


Why minimizes $\underbrace{Q = \sum e_i^2}_{\text{LSE}}$ instead of $\underbrace{Q_1 = \sum |e_i|}_{\text{MAD/LAD}}$

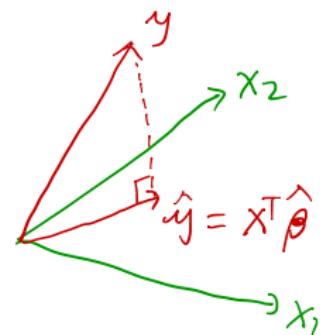
The reasons why we prefer Least Squares estimates instead of the minimum absolute deviation (MAD) estimates (or the Least Absolute Deviation (LAD) Estimator):

1. Squaring always gives a positive value and emphasizes large differences.
2. There is a strong link between OLS estimation and linear algebra. \hat{Y} is a linear function of Y . In fact, it is a projection onto a subspace defined by the independent variables.
3. When ϵ is normally distributed, the maximum likelihood estimate is the OLS estimates (coming soon). **OLS vs MLE**
4. OLS is BLUE.
5. Least squares method gives a single unique answer, with least absolute value you can get infinite which makes it harder to interpret the results.
6. A lot of nice properties happen with OLS. Such as in the normal error model, it is more efficient. (Reference: Huber, Robust Statistics, p.10)
7. For a sum of independent variables, variances can be added but not for the standard deviation.

(1) Also implies that Q is sensitive to outliers while Q_1 is less sensitive to outliers.



(2) Predicted outcome \hat{y} are the
orthogonal projection of y }
on subspace spanned by \vec{x} }



Different criterion of regression

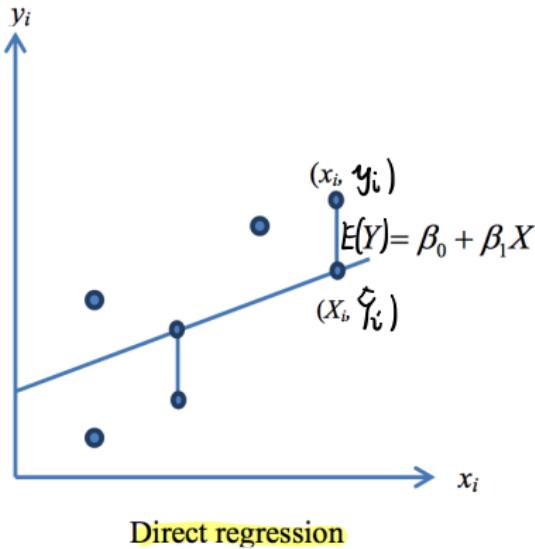
Direct regression

LSE: minimizes

$$Q = \sum_{i=1}^n e_i^2$$

LAD: minimizes

$$Q = \sum_{i=1}^n |e_i|$$

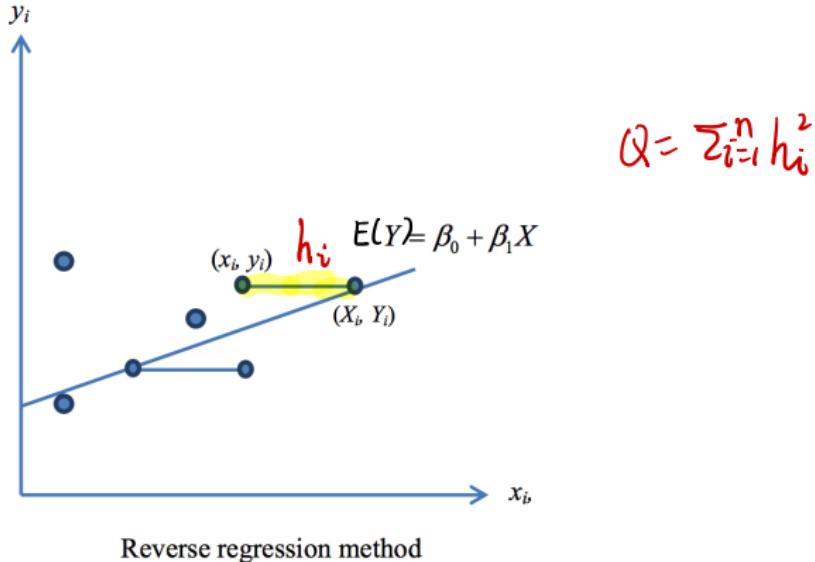


Direct regression

- This method is often referred to the Ordinary Least Squares (OLS) estimation which minimizes $Q = \sum (Y_i - b_0 - b_1 X_i)^2$
- Least Absolute Deviation (LAD) Regression Method: the estimates of β_0, β_1 are chosen such that the sum of absolute deviations is minimum.

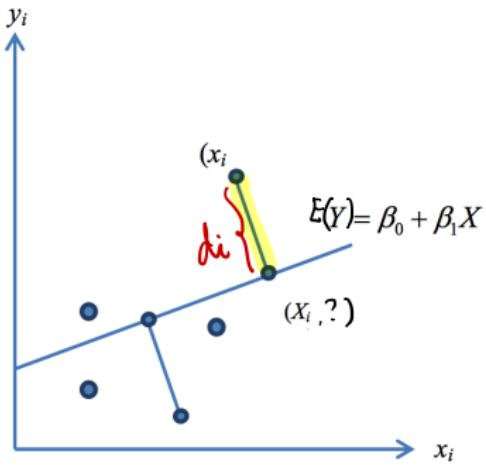
$$Q = \sum |Y_i - b_0 - b_1 X_i|$$

Reverse regression method



This reverse (or inverse) regression approach minimizes the sum of squares of horizontal distance between the observed data points and the fitted line to estimate β_0, β_1

Orthogonal regression method

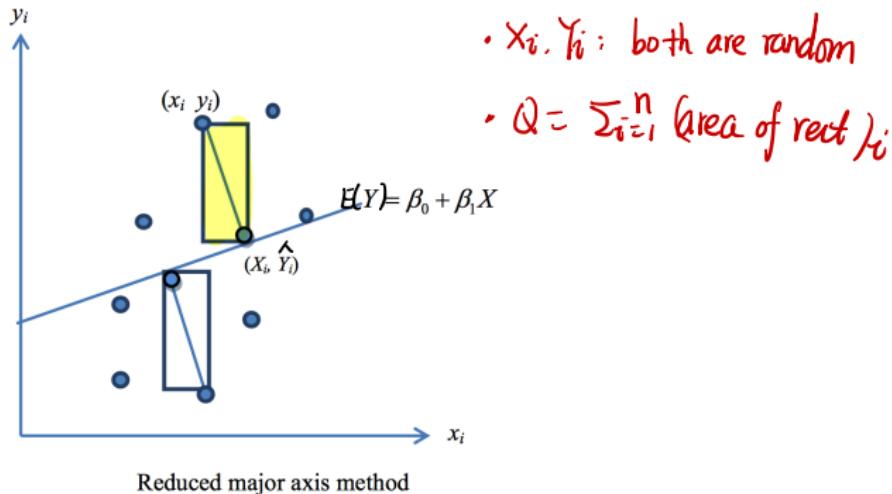


Major axis regression method

- both x_i and y_i are random
- $Q = \sum_{i=1}^n d_i^2$
 \uparrow
 d_i = Perpendicular distance between y_i and the fitted line.

Generally when uncertainties are involved in dependent and independent variables both, then orthogonal regression is more appropriate. The least squares principle in orthogonal regression minimizes the squared perpendicular distance between the observed data points and the fitted line to obtain the estimates of regression coefficients. This is also known as major axis regression method.

Reduced Major Axis regression



The area of rectangles defined between corresponding observed data points and nearest point on the fitted line can also be minimized. Such an approach is more appropriate when the uncertainties are present in study and explanatory variables both and is called as reduced major axis regression.

More discussion of different criterion of regression: Chapter 2.2, 2.3, 2.11-2.13, Linear Models and Generalizations, 3rd Edition, by C.R. Rao, H. Toutenburg, Shalabh, C. Heumann.

Normal Error Regression Model

Normal Error regression model

$$\left\{ \begin{array}{l} \text{needs only GM assumptions} \\ \text{(doesn't specify the distrn of } \epsilon_i) \end{array} \right. \quad \left\{ \begin{array}{l} \text{① } Y = \beta_0 + \beta_1 X + \epsilon \\ \text{② } E(\epsilon_i) = 0, \forall i \\ \text{③ } V(\epsilon_i) = \sigma^2, \forall i \\ \text{④ } \text{cov}(\epsilon_i, \epsilon_j) = 0, j \neq i \end{array} \right.$$

- Motivation: Least Squared estimators b_0, b_1 are BLUE no matter what distribution of the error terms ϵ_i . However, to set up interval estimates and make tests, we need to make an assumption about the form of the distribution of the ϵ_i .

- Normal error model description

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, \dots, n.$$

where $\epsilon \sim_{\text{iid}} N(0, \sigma^2)$. → stronger assumption than G-M.

$$\left\{ \begin{array}{l} \text{① } E(\epsilon_i) = 0 \\ \text{② } V(\epsilon_i) = 0 \\ \text{③ } \epsilon_i \perp \epsilon_j \\ \text{④ Normal} \end{array} \right.$$

- It follows that Y_i are independent normals (since mean varies with X_i).

$$Y_i \sim_{\text{iid}} \overset{\text{indep}}{N}(\beta_0 + \beta_1 X_i, \sigma^2) \quad \begin{matrix} E(Y_i) \\ \downarrow \\ V(Y_i) \end{matrix}$$

- $\epsilon_i \sim N(0, \sigma^2), Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$
- The first 2 moments determine the Normal distrn.

Review: Likelihood function and MLE

- Likelihood of the parameters given the data is proportional to the probability of the data given the parameters. i.e. $L(\theta|y) = P(Y|\theta)$
- If $Y_i \sim \text{IID } f(\theta), i = 1, \dots, n$, then the likelihood function is

not a pdf

$$L(\theta) = c f(y_1; \theta) f(y_2; \theta) \dots f(y_n; \theta) = c \prod_{i=1}^n f(y_i; \theta).$$

It's information in relative sense.

where c is an arbitrary non-negative constant.

- The Maximum Likelihood Estimate for a parameter is the value that maximizes the likelihood function.

$$\hat{\theta} = \max_{\theta \in \Omega} L(\theta)$$

- To maximize $L(\theta)$ is equivalent to maximize the log likelihood function $\ell(\theta)$

$$\ell(\theta) = \log L(\theta), \quad \hat{\theta} = \max_{\theta \in \Omega} \ell(\theta)$$

Symbol log in this course stands for the natural logarithm.

Maximum Likelihood Estimation: example

Consider a random sample $y_i, i = 1, \dots, n$ from $N(\mu, \sigma^2 = 1)$

- Step 1 • Likelihood function of μ

$$L(\mu) = \prod_{i=1}^n f(y_i; \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(y_i - \mu)^2}{2}\right\} = c \exp\left\{-\frac{\sum_1^n (y_i - \mu)^2}{2}\right\}$$

- Step 2 • Log likelihood function of μ

$$\ell(\mu) = \log L(\mu) = \log c - \frac{\sum_1^n (y_i - \mu)^2}{2}$$

- Step 3 • Find the MLE of μ , set the derivative to μ to zero and solve it for MLE

$$\frac{\partial \ell(\theta)}{\partial \theta} = \sum_1^n (y_i - \mu) = 0$$

- $\hat{\mu}_{MLE} = \bar{y}$.
- $$\frac{\partial^2 \ell(\theta)}{\partial \theta^2} = -n < 0 \Rightarrow L(\mu) at \bar{Y} reaches its maximum.$$

Assume p.d.f $f(y; \theta)$ is smooth enough
standard procedure to find MLE:

$$\textcircled{1} \text{ Likelihood } L(\theta; y_1, \dots, y_n) = C \prod_{i=1}^n f(y_i; \theta) \cdots f(y_n; \theta)$$

\Downarrow monotonic transf.

$$\textcircled{2} \text{ log-Likelihood } l(\theta) = \log L(\theta) = \sum_{i=1}^n \log f(y_i; \theta)$$

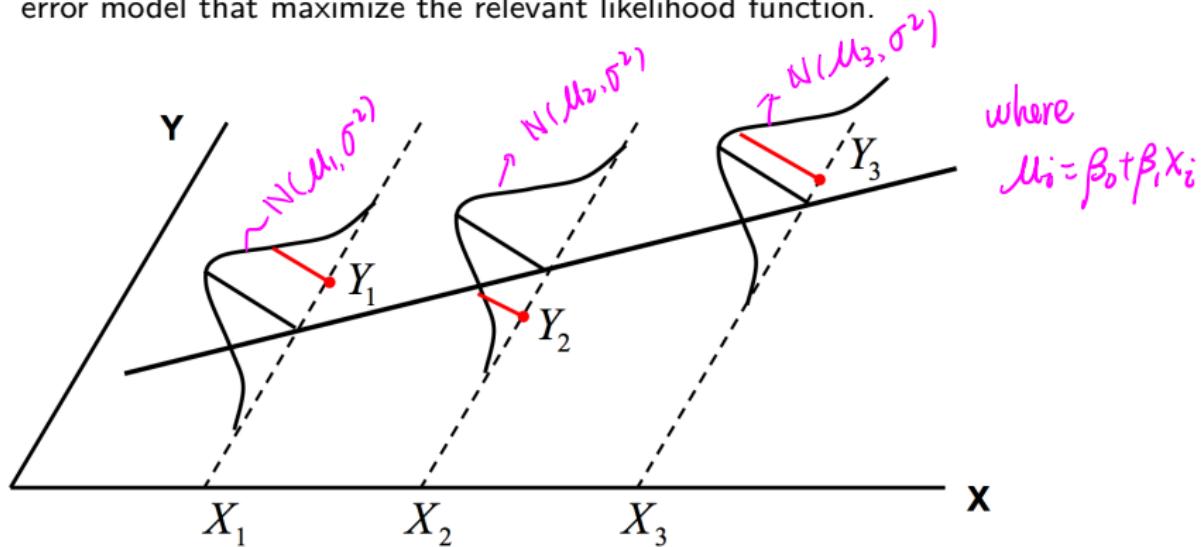
$$\textcircled{3} \quad \frac{\partial l(\theta)}{\partial \theta} = 0, \text{ if } \theta = (\begin{matrix} \theta_1 \\ \theta_2 \end{matrix}), \quad l_\theta(\theta) = \left(\begin{matrix} \frac{\partial l(\theta)}{\partial \theta_1} \\ \frac{\partial l(\theta)}{\partial \theta_2} \end{matrix} \right) = 0$$

$\longrightarrow \hat{\theta} \longleftarrow$

$$\textcircled{4} \text{ check } \frac{\partial^2 l(\theta)}{\partial \theta^2} \Big|_{\theta=\hat{\theta}} < 0$$

Maximum Likelihood Estimation for Regression

- Similarly, we find the values of parameters $\theta = (\beta_0, \beta_1, \sigma^2)$ in the normal error model that maximize the relevant likelihood function.



- $Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$, the density function at y_i is

$$f(y_i; \beta_0, \beta_1, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\}$$

MLE for Regression

① • Likelihood function

$$L(\beta_0, \beta_1, \sigma^2) = \prod_{i=1}^n f(y_i; \theta) = \left\{ \frac{1}{2\pi\sigma^2} \right\}^{n/2} \exp\left\{-\frac{\sum_1^n (y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}\right\}$$

② • Log likelihood function

$$\ell(\beta_0, \beta_1, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{\sum_1^n (y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}$$

③ Partial differentiation of $\ell(\beta_0, \beta_1, \sigma^2)$ yields

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad \left\{ \begin{array}{l} \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_0} = \frac{1}{\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i) \\ \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_1} = \frac{1}{\sigma^2} \sum x_i (y_i - \beta_0 - \beta_1 x_i) \end{array} \right. \quad (1)$$

$$\hat{\beta}_1 = \frac{\sum xy}{\sum x^2} \Leftarrow \left\{ \begin{array}{l} \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_1} = \frac{1}{\sigma^2} \sum x_i (y_i - \beta_0 - \beta_1 x_i) \\ \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (y_i - \beta_0 - \beta_1 x_i)^2 \end{array} \right. \quad (2)$$

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n} \quad \left\{ \begin{array}{l} \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (y_i - \beta_0 - \beta_1 x_i)^2 \end{array} \right. \quad (3)$$

MLE vs LSE for Regression

Set partial derivatives equal to zero and solve them to obtain the MLE of $\beta_0, \beta_1, \sigma^2$.

Parameters	MLE	Same as LSE?
β_0	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	Same. $\hat{\beta}_0 = b_0$
β_1	$\hat{\beta}_1 = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_1^n (x_i - \bar{x})^2}$	Same. $\hat{\beta}_1 = b_1$
σ^2	$\hat{\sigma}^2 = \frac{\sum_1^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{n}$	Different. $\hat{\sigma}^2 = \frac{n-2}{n} MSE$

Note that the MLE of σ^2 is a biased estimator of the error term variance.

Upcoming: Inference for SLR

- We knew how to **point estimate** of β_0, β_1 in SLR from data, Using Ordinary Least Squares (OLS) or MLE.
- Now we take up the inferences concerning the regression parameters β_0, β_1
 - How accurate are their estimates?
 - How to obtain an interval estimate for each of them?
 - How to test a specific parameter value of interest?
- Use Statistical Inference
 - Confidence intervals
 - Hypothesis Tests

Before we give answers, we will have a short review of statistical inference.

Review on Statistical Inference

Review: Central Limit Theorem (sample mean distribution)

CLT: X_1, \dots, X_n ~ one distn with $\mu, \sigma^2 < 0$, then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{\sum X_i - n\mu}{\sigma\sqrt{n}} \xrightarrow{D} N(0, 1)$$

1. Consider a random sample X_1, \dots, X_n from a distribution with mean μ and variance $\sigma^2 < \infty$.
2. Point estimate for the unknown mean is the sample mean, $\bar{X} = \sum X_i/n$.
3. **Central Limit Theorem (CLT):** as n approaches infinity, the random variables $\sqrt{n}(\bar{X} - \mu)$ converge in distribution to a normal $N(0, \sigma^2)$.

$$\sqrt{n}(\bar{X} - \mu) = \frac{1}{\sqrt{n}}(\sum X_i - n\mu) \xrightarrow{D} N(0, \sigma^2)$$

or put it in a brief way, $\bar{X} \xrightarrow{D} N(\mu, \sigma^2/n)$. *Normal is the limiting distn*

4. Note that if $X_i \sim_{iid} N(\mu, \sigma^2)$, then \bar{X} is exactly Normally distributed.

$$\sqrt{n}(\bar{X} - \mu) \sim N(0, \sigma^2)$$

Make sure you understand the difference between 3 and 4.

exact approximation.

四 $\frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}} \xrightarrow{D} N(0, 1)$, where $\theta = E(\hat{\theta})$

四 CLT: $\theta = \mu$, $\hat{\theta} = \bar{x}$, $E(\hat{\theta}) = E(\bar{x}) = \mu$

$$\Rightarrow \text{Var}(\hat{\theta}) = \text{Var}(\bar{x}) = \text{Var}\left(\frac{1}{n} \sum_i^n x_i\right) = \frac{1}{n^2} (\sigma^2 + \dots + \sigma^2) = \frac{\sigma^2}{n}$$

$$\begin{aligned} \Rightarrow \frac{\hat{\theta} - \mu}{\sqrt{\sigma^2/n}} &= \frac{\sqrt{n}(\hat{\theta} - \mu)}{\sigma} * \frac{n}{n} \\ &= \frac{\sqrt{n}(\sum x_i - n\mu)}{n\sigma} \end{aligned}$$

$$= \frac{\sum x_i - n\mu}{\sqrt{n}\sigma} \xrightarrow{D} N(0, 1)$$



all follow
the same formula.

$$E(\sum x_i) = n E(x_i) = n\mu$$

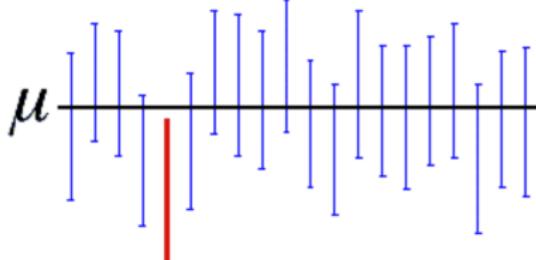
$$\text{Var}(\sum x_i) = n\sigma^2 \rightarrow \sqrt{n\sigma^2} = \sqrt{n}\sigma$$

Review: Confidence Interval for μ given σ^2

Two-side 95% CI for μ (assume σ^2 is known)

$$\bar{X} \pm Z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} = (\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}})$$

- A **confidence interval (CI)** is a type of interval estimate of a population parameter. It is an **observed interval** which is calculated from the observations, in principle different from sample to sample.
- How frequently the observed interval contains the parameter is determined by the **confidence level** or confidence coefficient.
- More specifically, **confidence level** suggests that if CI are constructed across many separate data analyses of replicated experiments, the proportion of such intervals that contain the true value of the parameter will match the given confidence level.



Replicate the same experiments 100 times, $(-2) \times 100$ times, the CI contains the true μ .

Example 2.1: Confidence Interval for μ given σ^2

- Observed a random sample (5,6,2,5,3) of size $n=5$, further assume $\sigma^2 = 1$.
- 95% CI (two-sided) for μ

$$\begin{aligned} 0.95 &= P\left(\frac{|\bar{X} - \mu|}{1/\sqrt{n}} < 1.96\right) & Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \rightarrow N(0,1) \\ &= P\left(|\bar{X} - \mu| < 1.96 \frac{1}{\sqrt{n}}\right) & P\left(\frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} \leq 1.96\right) = 0.95 \\ &= P\left(\bar{X} - 1.96 \frac{1}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{1}{\sqrt{n}}\right) & \Downarrow \\ & & P\left(\bar{X} - 1.96 \frac{1}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{1}{\sqrt{n}}\right) = 0.95 \end{aligned}$$

- For the given data, we find $\bar{X} = 4.2$ and the 95% CI for μ

$$(4.2 - 1.96/\sqrt{5}, 4.2 + 1.96/\sqrt{5}) = (3.323461, 5.076539)$$

Review: Hypothesis Testing

$$X_1, \dots, X_n \sim N(\mu, 1)$$

$$\begin{cases} H_0: \mu = 0 \\ H_a: \mu > 0 \end{cases}$$

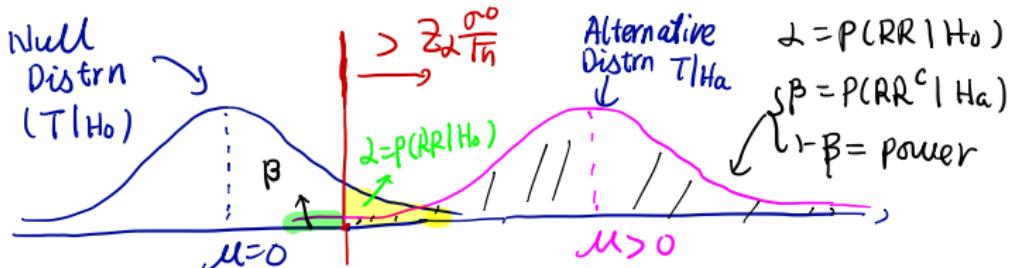
test statistic

$$T = \frac{\bar{X} - \mu}{\sigma_0 / \sqrt{n}}$$

distrn under H_0

$$\text{③ } T|_{H_0} = \frac{\bar{X} - \mu}{\sigma_0 / \sqrt{n}} \sim N(0, 1)$$

- A **hypothesis test** is a **statistical test** that is used to determine whether there is enough evidence in a sample of data to infer the **null hypothesis** is true.
- P-value for a testing indicates where the observed data sits on the **null distribution**.
- In some textbook, they view P-value as a way to determine "likely" or "unlikely" of observing a more extreme statistic in the direction of alternative hypothesis than the one observed assuming null hypothesis were true.



Reject H_0 when $\frac{\bar{X} - \mu}{\sigma_0 / \sqrt{n}} > z_\alpha \Leftrightarrow \bar{X} \geq z_\alpha \frac{\sigma_0}{\sqrt{n}}$

Review: Hypothesis Testing

R.R. = Rejection Region

Reality/ Truth

		Null is True $H_0: T$	Null is False $H_0: F$
		Reject Null	Correct Decision
Decision	Reject Null	Type I error "False Positive"	Correct Decision
	Fail to Reject Null	Correct Decision	Type II error "False Negative"

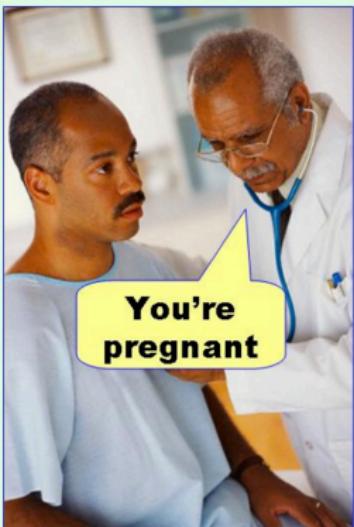
- Type I error = $P(\text{reject } H_0 | H_0 \text{ is true}) = P(\text{R.R.} | H_0) = \alpha$
- Type II error = $P(\text{fail to reject } H_0 | H_0 \text{ is false}) = P(\text{fail to reject } H_0 | H_a \text{ is true}) = P(\text{R.R.}^c | H_a) = \beta$
- Significance level α : maximum allowable probability of making type I error (usually 5%)
- Test statistics: a statistics (function of data variable) whose distribution under H_0 is known.

Null distribution

Review: Hypothesis Testing

Type I error

(false positive)



Type II error

(false negative)



<https://effectsizefaq.com/index/faqs/page/3/>

Review: Hypothesis Testing

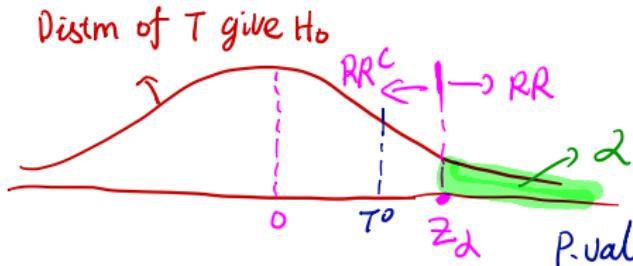
$$H_0: \theta \in \mathcal{R}_0$$

$$H_a: \theta \notin \mathcal{R}_0$$

- If we reject H_0 (i.e. H_0 is false), then we accept H_a .
- However, if we fail to reject H_0 , it doesn't imply that H_0 is true.

- P-value: the probability (under H_0) that the test statistics is equally or more extreme than the one actually observed.
 - If $P\text{-value} < \alpha$, reject H_0 . $\Rightarrow H_a$ is True.
 - If $P\text{-value} > \alpha$, do not reject H_0 . This doesn't imply that H_0 is true.
- Smaller P-value \rightarrow less likely H_0 is to be true
 - $p\text{-value} < 0.001 \rightarrow$ strong evidence against H_0
 - $p\text{-value} < 0.05 \rightarrow$ moderate evidence against H_0
 - $p\text{-value} < 0.1 \rightarrow$ weak evidence against H_0
 - $p\text{-value} > 0.1 \rightarrow$ no evidence against H_0

$$\begin{aligned} H_0: \mu = 0 \\ H_a: \mu > 0 \end{aligned}$$



$$\begin{aligned} T^0 &= \frac{\bar{x}^0 - 0}{\sigma^0 / \sqrt{n}} \\ &= \text{observed } T \\ &\quad (\text{plugging data}) \end{aligned}$$

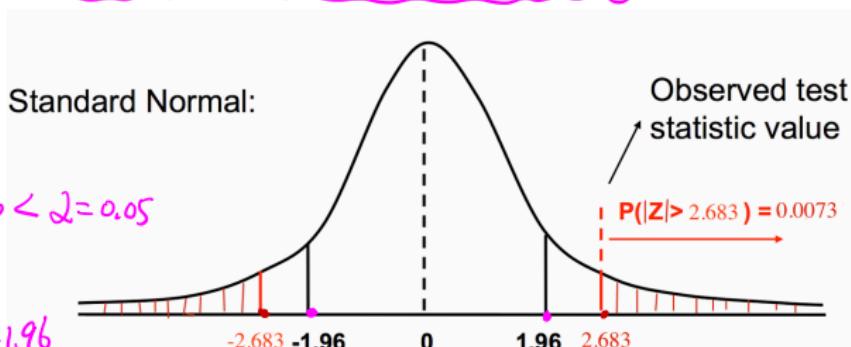
Example 2.1: Hypothesis Testing

- Observed a random sample $(5, 6, 2, 5, 3)$ of size $n=5$, assume $\sigma^2 = 1$.
- 95% CI (two-sided) for μ : $(3.323461, 5.076539) \leftarrow \bar{X} \pm 1.96 \frac{1}{\sqrt{n}}$
- Now interested in testing $H_0 : \mu = 3$, $H_a : \mu \neq 3$
 - test statistics:

$$T = \frac{\bar{X} - \mu}{1/\sqrt{5}} \sim_{H_0} Z = N(0, 1)$$

- given data, under H_0 , $T_{obs} = (4.2 - 3)/\sqrt{5} = 2.683$
- p-value = $P(|Z| > T_{obs}) = P(|Z| > 2.683) = 0.0073$
- make decision: we have strong evidence to reject the H_0 since p-value is less than 0.01, or because the $T_{obs} > 1.96$ for $\alpha = 5\%$.
- Note that $\mu = 3 \notin 95\% \text{ CI: } (3.323461, 5.076539)$. ③

Methods
①, ②, ③;
consistent.



① $P(\text{value}) = 0.0073 < \alpha = 0.05$

② $T_{obs} = 2.683$

$|T_{obs}| > z_{\alpha/2} = 1.96$

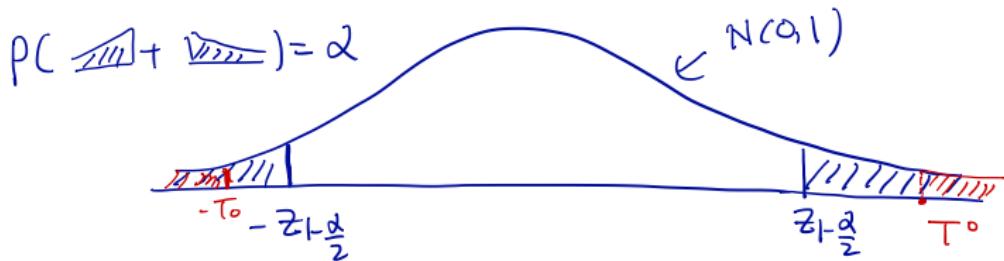
① ② give consistent result: reject H_0 .

Review: P-value and CI agree about statistical significance

- You can use either P values or confidence intervals to determine whether your results are statistically significant. (*I Reject H₀*)
- The confidence level is equivalent to $1 - \alpha$ level. So, if your significance level is 0.05, the corresponding confidence level is 95%.
 - If the P value is less than your significance (α) level, the hypothesis test is statistically significant.
 - If the CI does not contain the null hypothesis value (H_0 value), the results are statistically significant.
 - If $P\text{-value} < \alpha$, the CI will not contain the null hypothesis value.
- To understand why the results always agree, let's recall how both the significance level and confidence level work.
 - The significance level defines the distance the sample mean must be from H_0 to be considered statistically significant.
 - The confidence level defines the distance for how close the confidence limits are to sample mean.
 - Both the significance level and the confidence level define a distance from a limit to a mean. Guess what? The distances in both cases are exactly the same!

Assume two-sided test: $H_0: \theta = \theta_0$ $H_a: \theta \neq \theta_0$

$T|_{H_0} \sim$ A known distm, say $N(0,1)$



④ 3 ways to make decision: Accept or Reject H_0

$$T^* = T(\text{plugging in observed data})|_{H_0}$$

① P-value = $P(|Z| > T^*)$, P-value $< \alpha$: Reject

② $|T^*| > z_{1-\frac{\alpha}{2}}$: Reject

③ $\theta_0 \notin (1-\alpha)$ confidence interval.

Upcoming topic and reading

- Review of distribution theory.
- Inference concerning β_0, β_1 .
- Interval estimation of mean response.
- Prediction interval.
- Difference between prediction interval and confidence interval.

Reading for next week: CH1: 1.8, CH2: 2.1,2.2, 2.4,2.5,2.6, 2.7

Practice problems after Week 1- lecture B

1. Keep trying the practice problem

- 1.3, 1.5, 1.6, 1.7, 1.8, 1.11, 1.16, 1.18, 1.20(*), 1.21(*), 1.24(*), 1.29, 1.30, 1.33, 1.36, 1.39a, 1.40, 1.41a.
- For questions marked (*), the SAS code & output is posted with the solutions.
- You only need to interpret that output.

2. Show that

- $\text{Cov}(\bar{Y}, b_1) = 0$ (Hint: $\text{Cov}(\bar{Y}, b_1) = \text{Cov}\{\sum \frac{1}{n} Y_i, \sum k_i Y_i\}$)
- Using above prove that $V(b_0) = \sigma^2(1/n + \bar{X}^2/S_{xx})$
- $\text{Cov}(b_0, b_1) = -\sigma^2 \bar{X}/S_{xx}$
- $\text{Cov}(\bar{Y}, b_0) = \sigma^2/n$

3. A no-intercept model is: $Y_i = \beta_1 X_i + \epsilon_i, \epsilon \sim_{iid} N(0, \sigma^2)$

- find the LS estimate of β_1 , denoted it as b_1
- find the variance of b_1 , $V(b_1)$
- find an unbiased estimate of $V(b_1)$