

RECAP

①

Multivariate regression model $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{i,p-1} + \varepsilon_i, i=1, \dots, n$
 $E(\varepsilon_i) = 0, \text{Var}(\varepsilon_i) = \sigma^2$

Or vector form: $y = \bar{X}\beta + \varepsilon$
 $n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1$

Least Squares: $\min_{\beta} Q(\beta) = (y - \bar{X}\beta)'(y - \bar{X}\beta)$

$\Rightarrow \frac{\partial Q}{\partial \beta} = 0 \Rightarrow$ Normal Equations $\boxed{(\bar{X}'\bar{X})\beta = \bar{X}'y}$

$$\approx \text{Cov}(\bar{X})\beta = \text{Cov}(\bar{X}, y)$$

$$\text{Cov}(\bar{X}) = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \dots \\ \vdots & \vdots & \ddots \\ \vdots & \vdots & \vdots & \text{Var}(x_{p-1}) \end{pmatrix} \quad \begin{pmatrix} \text{Cov}(x_1, y) \\ \text{Cov}(x_2, y) \\ \vdots \\ \text{Cov}(x_{p-1}, y) \end{pmatrix}$$

How X's are related

How each X is related to y

② $\text{Cov}(\bar{X})$ is diagonal, i.e. all X's are uncorrelated

$$\Rightarrow \hat{\beta} = (\bar{X}'\bar{X})^{-1} \bar{X}'y \approx (\underbrace{\text{Cov}(\bar{X})}_{\text{diagonal}})^{-1} \text{Cov}(\bar{X}, y)$$

$$\approx \begin{pmatrix} \text{Cov}(x_1, y) \\ \text{Cov}(x_2, y) \\ \vdots \\ \text{Cov}(x_{p-1}, y) \end{pmatrix} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_{p-1} \end{pmatrix}$$

Meaning - if x 's are uncorrelated, $\hat{\beta}_j$ measures how x_j and y are related. ②

① $\text{Cov}(X)$ not diagonal $\rightarrow \hat{\beta} = (X'X)^{-1} X'y$

$$\approx (\text{Cov}(X))^{-1} \text{Cov}(X, y)$$

$$= \begin{pmatrix} \text{corr}(x_1, y) \text{ \& other } x-y \text{ correlations} \\ \text{corr}(x_2, y) \text{ --- } \dots \\ \vdots \\ \text{corr}(x_p, y) \text{ --- } \dots \end{pmatrix}$$

Meaning $\hat{\beta}_j$ measures not only how x_j and y are related, but also how other x 's relate to y .

Properties $E(\hat{\beta}) = \beta$ unbiased

$$V(\hat{\beta}) = \sigma^2 (X'X)^{-1} \approx \frac{\sigma^2}{n} (\text{Cov}(X))^{-1}$$

$$\sigma^2 \uparrow \Rightarrow V(\hat{\beta}) \uparrow$$

do more noise means more uncertainty

$$n \uparrow \Rightarrow V(\hat{\beta}) \downarrow$$

do more data means less uncertainty

more spread in x 's (large values on diagonal of $\text{Cov}(X)$) $\Rightarrow V(\hat{\beta}) \downarrow$

but if x 's are highly correlated (large values off-diagonal in $\text{Cov}(X)$) $\Rightarrow V(\hat{\beta}) \uparrow$ in general

So - correlated x 's in data makes estimation more difficult (more uncertain) (see Sect 5 of Lab 2) ③

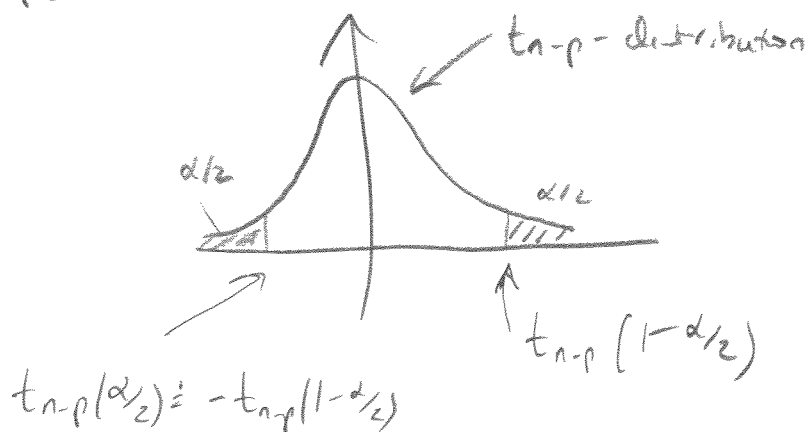
Testing

$$t_j = \frac{\hat{\beta}_j - 0}{SE(\hat{\beta}_j)}, \text{ where } SE(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}_{jj}}$$

$\underbrace{\hspace{10em}}$
jth diagonal element
of $\mathbf{X}'\mathbf{X}$.

If true $\beta_j = 0$, t_j follows the t_{n-p} -distribution.

So if observed t_j is extreme for t_{n-p} , we reject the null $\beta_j = 0$

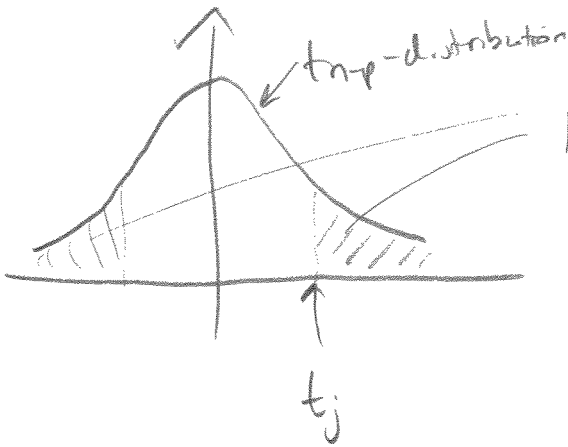


If $|t_j| > t_{n-p}(1-\alpha/2)$, reject $\beta_j = 0$ at level α .

Note, this means we reject the null, WHEN IT IS TRUE, with probability α , so keep α small (e.g. 1%)

OR compute p-values

(4)



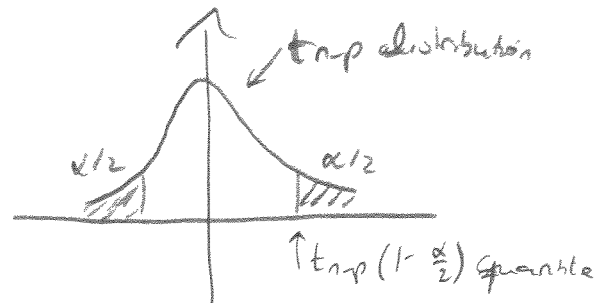
probability mass in t_{n-p} distribution with values more extreme than t_j = p-value

observed t-value

∪ p-value is small, states that t_j is unusual to get if null is true - so we start to question the null

OR compute confidence intervals

Know $\frac{\hat{\beta}_j - \beta_j}{SE(\hat{\beta}_j)} \sim t_{n-p}$ ← true β_j



Means $Prob\left(\left|\frac{\hat{\beta}_j - \beta_j}{SE(\hat{\beta}_j)}\right| \leq t_{n-p}(1 - \alpha/2)\right) = 1 - \alpha$

∴ $Prob\left(-t_{n-p}(1 - \alpha/2) \leq \frac{\hat{\beta}_j - \beta_j}{SE(\hat{\beta}_j)} \leq t_{n-p}(1 - \alpha/2)\right) = 1 - \alpha$ (Random)

∴ $Prob\left(\underbrace{\hat{\beta}_j - t_{n-p}(1 - \alpha/2) SE(\hat{\beta}_j)}_{\text{Random interval}} \leq \beta_j \leq \underbrace{\hat{\beta}_j + t_{n-p}(1 - \alpha/2) SE(\hat{\beta}_j)}_{\text{Random interval}}\right) = 1 - \alpha$

So, random interval

⑤

$$[\hat{\beta}_j \pm t_{n-p}(1-\alpha/2) SE(\hat{\beta}_j)]$$

Covers the β_j with probability $1-\alpha$.

So, if interval does not cover 0, reject

the null $\beta_j=0$ at level α .

Model selection via testing

• Keep only β_j 's that are significant (null rejected in t-test)

• Complication

① Multiple Testing

② Dependent tests

⇓
X's are correlated

⇒ $\hat{\beta}$'s are correlated

⇒ tests are correlated

⇓
Individual testing

of β_j can be
very misleading

(lect 5, Lab 2)

Each test (for each β_j) has probability α of leading to a false rejection, i.e.

null $\beta_j=0$ is rejected even when it is true!
That means, we keep β_j when we shouldn't.

⇓
If we do P tests; Prob(at least one false rejection) ⇒ α

In fact $P=10 \Rightarrow$ Prob(at least one false rejection) $\approx 40\%$
 $\alpha=5\%$

$P=25 \Rightarrow$ - - - - - $\approx 72\%$

Comparing models using the F-test

⑥

Model 1 : k variables in model, # parameters $p_1 = k + 1$

Model 2 : m variables (a subset of the k in model 1)
parameters $p_2 = m + 1$

Fit both models to the data and compute their residual sum of squares, $RSS(\text{Model 1})$, $RSS(\text{Model 2})$

Null hypothesis : Model 2 is sufficient (formally testing β 's for variables in model 1 but not in model 2 are all 0)

$$F_{\text{obs}} = \frac{\left[\frac{RSS(2) - RSS(1)}{p_1 - p_2} \right]}{\left[\frac{RSS(1)}{n - p_1} \right]}$$

} Drop in RSS / extra parameter it cost

} Best estimate of σ^2

If null is true (model 2 sufficient) $F_{\text{obs}} \sim F_{p_1 - p_2, n - p_1}$

So if $F_{\text{obs}} > \bar{F}_{p_1 - p_2, n - p_1}(1 - \alpha)$

, reject model 2.
(in favor of model 1)

$1 - \alpha$ quantile of
F distribution

Question: Which models to compare?