

R code: Recall the advertising example. From fitting $\text{lm}()$ we get the estimated coeffs $(\hat{\beta})$ so ex.

$$\text{sales} = 2.93 + 0.046 \cdot \text{TV} + 0.139 \cdot \text{Radio} + (-0.001) \cdot \text{Newspaper}$$

Recall • residuals $:= e_i = y_i - \hat{y}_i$; $\text{len}(y_i, \hat{y}_i) = n$.

• $\hat{\beta} = (X^T X)^{-1} X^T Y$ with $\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$
 (p x p)

(need that variables are uncorrelated and have equal variance)

• $\sigma^2 = \frac{RSS}{n-p} = \frac{\sum_{i=1}^n e_i^2}{n-p}$
 degrees of freedom $:= (n-p)$ makes σ^2 unbiased.

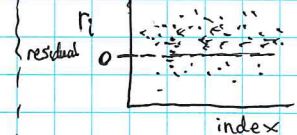
(p := rank of X)

s.t. $\text{Var}(\hat{\beta}_j) = [\sigma^2 (X^T X)^{-1}]_{jj}$ ← ie jth element in the diagonal of $\text{Var}(\hat{\beta})$

$$\text{Var}(\hat{\beta}) = \begin{bmatrix} \text{Var}(\hat{\beta}_1) & \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) & \dots \\ \vdots & \text{Var}(\hat{\beta}_2) & \\ \vdots & & \ddots \\ \vdots & & & \text{Var}(\hat{\beta}_p) \end{bmatrix}$$

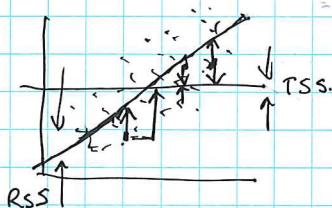
* SEE Script 1.4

• Residual standard error (RSS) $\equiv \hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-p}}$
 helps measures the goodness of fit.



• Issue: residuals may not be ^{normally} ~~uniformly~~ (randomly) distributed
 ↳ they don't fluctuate around ^{randomly} zero, i.e. skewed.

• RSS: residuals wrt regression line (\hat{y}_i) TSS: residuals wrt mean level/value (\bar{y}_i) , ie
 $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ $TSS = \sum_{i=1}^n (y_i - \bar{y}_i)^2$



Note $TSS \geq RSS$.

• Sample variance of y-values $:= \frac{1}{n-1} \cdot TSS$
 • Remaining variance after fitting the model $:= \frac{1}{n-1} \cdot RSS$

• Proportion of variance that remains after fitting the model $[0,1]$ $:= \frac{RSS}{TSS} = \frac{RSS/(n-1)}{TSS/(n-1)}$ (smaller value better)

• $R^2 = 1 - \frac{RSS}{TSS}$ $:=$ proportion of variance that is explained by the model

$R^2 \in [0,1]$ → larger values are better

- Not affected by scale of y
- R^2 is the same as (sample correlation)², but R^2 is more general. (in simple linear regression)

• Adding more variables to the model can only increase R^2 (as RSS can only decrease)
 → means we can't fairly compare models of which one is a submodel of the other.
 (Account for this with R_{adj} !)

- R^2 (cont'd) To account for comparing ^{subsets of} ~~stat~~ models, we adjust R^2 according to the model's degrees of freedom $(n-p)$:
 i.e. $R^2_{adj} = 1 - \frac{RSS/(n-p)}{TSS/n}$ \Rightarrow adding unnecessary variables to the model will decrease R^2_{adj} (as $RSS/(n-p)$ can only increase)
 ie $R^2_{adj} < R^2$ σ^2 is unknown

- Standard error of $\hat{\beta}_j$. We generally can't calculate $\text{Var}(\hat{\beta}_j) = [\sigma^2(X^T X)^{-1}]_{jj}$ but we can calculate $\widehat{\text{Var}}(\hat{\beta}_j) = [\hat{\sigma}^2(X^T X)^{-1}]_{jj}$
 \rightarrow Standard of error of $\hat{\beta}_j = \sqrt{\widehat{\text{Var}}(\hat{\beta}_j)} = \sqrt{[\hat{\sigma}^2(X^T X)^{-1}]_{jj}}$

- Distribution of least squares estimates assuming Gaussian errors.

Since we assume a linear model and require $\varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \sim \mathcal{N}(0, \sigma^2)$, we can show that

$$(i) \left[\begin{aligned} \hat{\beta} &\sim \mathcal{N}_p(\beta, \sigma^2(X^T X)^{-1}) \equiv \mathcal{N}_p(\beta, \text{Var}(\hat{\beta})) \\ \hat{\beta}_j &\sim \mathcal{N}(\beta_j, [\sigma^2(X^T X)^{-1}]_{jj}) \equiv \mathcal{N}(\beta_j, \text{Var}(\hat{\beta}_j)) \end{aligned} \right]$$

- p-values: prob (observing a test statistic that is at least as extreme than the one we ~~see~~ ^{see} when the null hypothesis is true)

ie for

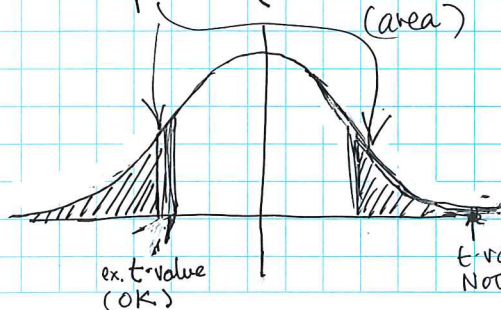
- From the normal distribution of $\hat{\beta}_j$ we derive
 Under the null-hypothesis, $\hat{\beta}_j \sim \mathcal{N}(0, [\sigma^2(X^T X)^{-1}]_{jj} \equiv \text{Var}(\hat{\beta}_j))$
 \rightarrow we can get $\frac{\hat{\beta}_j}{\sqrt{\text{Var}(\hat{\beta}_j)}} \sim \mathcal{N}(0, 1)$ (standard normal)
 with $\begin{cases} \beta_j = 0 \text{ if } H_0 \text{ (null hypothesis)} \\ \beta_j \neq 0 \text{ if } H_a \text{ (alternative hypothesis)} \end{cases}$

Substituting the unknown quantity σ^2 with the estimate $\hat{\sigma}^2$ we get the t-test statistic

ie $T_j = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2(X^T X)^{-1}_{jj}}} \sim t_{n-p} \text{ under } H_0$

the p-value in the t-test refers to/is defined as:

(under H_0) $p = \text{Prob}(|T\text{-statistic}| > |t\text{-value}| \mid H_0 \text{ is true})$



• An individual t-test for H_0 should be interpreted as quantifying the effect of the j th predictor ^{variable} after having subtracted the linear effects of all other predictor variables on Y .

(slightly different distr from $\mathcal{N}(0, 1)$; thicker tails)
 As $n-p \rightarrow \infty$, $t_{n-p} \approx \mathcal{N}(0, 1)$

• p values can vary substantially depending on which/how many variables are included in the model \Rightarrow need to consider which ones to include.

Lecture 3: Least Squares Estimates (Properties); tests, confidence regions, cont'd 01 Mar 2018

p Values, cont'd

ex. If we regress $\text{Sales} \sim \text{TV} + \text{Radio} + \text{Newspaper}$, then the p-value for Newspaper:

$$H_0: \text{sales} \sim \text{TV} + \text{Radio}$$

$$H_a: \text{sales} \sim \text{TV} + \text{Radio} + \text{Newspaper}$$

- An insignificant p-value says that there is no / little evidence in the data that shows H_a is better than H_0 . ie. so if you're already using radio + TV, it doesn't help to add Newspaper

ex. $\text{Sales} \sim \text{Newspaper}$: $H_0: \text{sales} = \beta_1 + \epsilon$

F-test. Question: For a given variable after adding all the other variables, does it help to add this one?

- Partial F-test (anova) can compare subset and the whole model.

ex. Model 1: $\text{sales} \sim \text{TV} + \text{radio}$
Model 2: $\text{sales} \sim \text{TV} + \text{radio} + \text{newspaper}$

If we want to look at only newspaper, then

Model 1: $\text{sales} \sim 1$, Model 2: $\text{sales} \sim \text{newspaper}$

F-statistic example: checking full vs. empty model (football data)

It is important for checking / looking at multiple testing problems.

- If F-stat is bad it means that the model is not good (start over).
- If good / significant, continue analysis, looking at specific variables.

R internally creates dummy variables.

ex. $\text{balance}_i = \beta_1 + \beta_2 \cdot \text{income}_i + \beta_3 x_i + \epsilon_i$, $x_i = \begin{cases} 1 & \text{if } i \text{ is student} \\ 0 & \text{otherwise} \end{cases}$

$$= \begin{cases} \beta_1 + \beta_2 \cdot \text{income}_i + \beta_3 + \epsilon_i & \text{if student} \\ \beta_1 + \beta_2 \cdot \text{income}_i + \epsilon_i & \text{otherwise} \end{cases}$$

↪ slope is the same (β_2)
y-int differs: student: $\beta_1 + \beta_3$
other: β_1

