

Contents

Day 27	1
Multiple Linear Regression	2
ANOVA for Multiple Linear Regression	3
ANOVA Table	3
t-Test for Slope in Multiple Linear Regression	4
Model Selection	5

Day 27

Multiple Linear Regression

Model

$$\mu_{y|x} = \beta_o + \beta_1 x_1 + \dots + \beta_p x_p$$

$$y = \mu_{y|x} + \epsilon_1, \epsilon \sim N(0, \sigma)$$

Least Squares Line:

$$\hat{y} = b_o + b_1 x_1 + \dots + b_p x_p$$

$$y_i = \hat{y}_i + e_i$$

Assumptions

- Linear relationship
 - Between y and each x_j in the model*
- Normally distributed residuals
 - Normal q-q plot
- Residuals have mean of 0 and standard deviation $\sigma = 0$
 - Independent of \hat{y}
- Independent residuals
 - Residual plot: e_i vs \hat{y}_i
- All the variables in the model are independent (usually settle for uncorrelated)*

Assumptions 2-4 are for inference

* check with scatter plot matrix

Simplified multiple linear regression model:

- $\mu_{y|x} = \beta_o + \beta_1 x_1 + \beta_2 x_2$
- $\hat{y} = b_o + b_1 x_1 + b_2 x_2$

Interpreting b_o and β_o :

- Average value of y when all x -variables are 0

Interpreting b_1 and β_1 :

- Average change in y for a 1 unit increase in x_1
 - Both are the same, pick one
 - Holding (the effect of) x_2 constant. Think of a partial derivative (∂)
 - After accounting for the other variables in the model

ANOVA for Multiple Linear Regression

Hypothesis:

- Population mean is the same and is unaffected by anything.
- ANOVA for Linear Regression: population mean is estimated well by null model, μ_{x_1, \dots, x_p}
- Equivalent: $\beta_1 = 0 \dots \beta_p = 0$

ANOVA Table

try to fill in

Interpretation:

$$F_{\text{observed}} \sim F(P, n - P - 1)$$

- If p-value \leq significance level \implies reject H_0
 - \therefore Our model is significantly better than the null model at explaining changes in y

IMPORTANT!!!!

- This means one or more x-variables in the model are required for the better model. It DOES NOT tell us which one(s), and it CERTAINLY DOES NOT mean they are all important!
- If p-value $>$ significance level \implies fail to reject H_0
 - \therefore our model is not significantly better than the null model at explaining changes in y . We prefer the null model.

Important: some x-variables may still be important predictors. However, we may not see their effect if “more important predictors” are left out of the model.

t-Test for Slope in Multiple Linear Regression

Model:

$$\mu_{y|x_1 \dots x_p} = \beta_o + \beta_1 x_1 + \dots + \beta_p x_p$$

- H_o: In this model, $\beta_j = 0$
 - β_j is the slope multiplying x_j , $1 \leq j \leq p$
- H_a: In this model: $\beta_o \neq 0$

$$t = \frac{\text{stat} - \text{parameter}}{\text{SE}} = \frac{b_j - \beta_j}{SE_{b_j}}$$

$$t_{\text{observed}} \sim t(n - p - 1)$$

Interpretation:

- If p-value \leq significance level \implies reject H_o
 - $\therefore x_j$ is a significant predictor of y , even after accounting for the effect of the other variables in the model
- If p-value $>$ significance level \implies fail to reject H_o
 - $\therefore x_j$ is not a significant predictor of y in this model. We CANNOT distinguish between two competing explanations.
 1. x_j does not have a linear relationship with y
 2. The effect of x_j on y is already accounted for by other variables in the model. It is redundant.

Model Selection

General question: which model is the best?

- Step 1: [Feature engineering](#)
 - [Common sense & explanatory analysis](#)
 - Goal: identify important variables
- Step 2: [Decide on a model selection algorithm and a selection criterion](#)
 - Our algorithm: [backward selection](#)
 - Our criterion: Stop when all explanatory variables are significant (at 5 % level)
 - [DO NOT](#) use R^2 as a selection criterion in multiple linear regression
- Step 3: [Implement the algorithm](#)

Collinearity: x_1 & x_2 are highly correlated. [This is bad!](#)

- Remove least significant predictor. This means removing the variable with the highest p-value.
ONLY REMOVE ONE VARIABLE AT A TIME.