

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Online Review Course of Undergraduate Probability and Statistics

Review Lecture 17

Linear Regression, part 2

Chris A. Mack
Adjunct Associate Professor

Course Website: www.lithoguru.com/scientist/statistics/review.html

© Chris Mack, 2014 1

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Linear Regression Steps

1. Create an x-y scatterplot of the data
 - Is the assumption of a linear relationship reasonable?
 - Are there any obvious outliers?
2. What do you know about the underlying pdfs of X and Y?
 - Do you have reason to assume normal distributions, or are you doing so for convenience?
3. Perform the regression
4. Plot the residuals
 - Check the assumption of iid normal

© Chris Mack, 2014 4

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

MLE Linear Regression

- Our model: $E[Y|X] = aX + b$

$$a = \frac{cov(X,Y)}{var(X)} \quad b = E[Y] - aE[X]$$
- Our LSE estimators:

$$a = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad b = \bar{y} - a\bar{x}$$
- For Y_i normally distributed, the MLE is the same as the LSE

© Chris Mack, 2014 2

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Linear Regression Graphs

- Create an x-y scatterplot of the data

Model: $y_i = ax_i + b + \epsilon_i$

Predicted Value: $\hat{y}_i = ax_i + b$

Residual: $\epsilon_i = y_i - \hat{y}_i$

© Chris Mack, 2014 5

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

MLE Linear Regression Assumptions

- Y is linearly related to X
- X_i is known with certainty; only Y_i has uncertainty, some of which is explained by X_i
- All Y_i are iid normally distributed
- There are no outliers: our regression is not robust, and even one bad data point will mess up our results

© Chris Mack, 2014 3

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Linear Regression Graphs

- Create an x-y scatterplot of the residuals

Plot ϵ_i vs. \hat{y}_i or x_i

Plot a histogram of ϵ_i

Note that $\bar{\epsilon} = 0$

Look for violations of the iid normal assumption – do you see any trend in the residuals?

© Chris Mack, 2014 6

THE UNIVERSITY OF TEXAS
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

Goodness of Fit

- Goodness of fit metric: $R^2 = r^2$

$$r = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} \quad R^2 = \frac{\text{cov}^2(X, Y)}{\text{var}(X)\text{var}(Y)}$$

- Also, we can show that

$$R^2 = 1 - \frac{\text{var}(\epsilon)}{\text{var}(Y)}$$

R^2 is the fraction of the variance of Y that is explained by the linear fit

© Chris Mack, 2014 7

THE UNIVERSITY OF TEXAS
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

Review #17: What have we learned?

- What are the assumptions for a standard least-squares linear regression?
- Why is it a good idea to plot our residuals?
- What is the goodness of fit and how do we interpret its value?
- How does one construct confidence intervals for the linear fit parameters and for predicted values from our fit equation?

© Chris Mack, 2014 10

THE UNIVERSITY OF TEXAS
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

Uncertainty of Fit Parameters

- The regression fit is based on a sample of data

Population Model: $y_i = \alpha x_i + \beta + \epsilon_i$ where $\epsilon_i \sim N(0, \sigma_\epsilon^2)$

- To create confidence intervals for a , b , \hat{y}_i , and r we need to know their sampling distributions
 - They are unbiased and t-distributed (DF = $n - 2$)

$$\text{var}(a) = \frac{\text{var}(\epsilon)}{n \text{var}(X)} \quad \text{var}(b) = \frac{\text{var}(\epsilon)}{n} \left(1 + \frac{\bar{x}^2}{\text{var}(X)} \right)$$

© Chris Mack, 2014 8

THE UNIVERSITY OF TEXAS
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

Uncertainty of Predictions

- Uncertainty in \hat{y}_i comes from the spread of the residuals and from uncertainty in the best fit parameters a and b
- Again, the sampling distribution will be Student's t with DF = $n - 2$

$$\text{var}(\hat{y}_i) = \frac{\text{var}(\epsilon)}{n} \left(1 + \frac{(x_i - \bar{x})^2}{\text{var}(X)} \right)$$

© Chris Mack, 2014 9