Section 2.1: Intro to Simple Linear Regression & Least Squares

Jared S. Murray The University of Texas at Austin McCombs School of Business Suggested reading: OpenIntro Statistics, Chapter 7.1, 7.2

Regression: General Introduction

- Regression analysis is the most widely used statistical tool for understanding relationships among variables
- It provides a conceptually simple method for investigating functional relationships between one or more factors and an outcome of interest
- The relationship is expressed in the form of an equation or a model connecting the response or dependent variable and one or more explanatory or predictor variables

Why?

Straight prediction questions:

- For how much will my house sell?
- How many runs per game will the Red Sox score this year?
- Will this person like that movie? (e.g., Netflix)

Explanation and understanding:

- What is the impact of getting an MBA on lifetime income?
- How do the returns of a mutual fund relate to the market?
- Does Walmart discriminate against women when setting salaries?

Example: Predicting House Prices

Problem:

Predict market price based on observed characteristics

Solution:

- Look at property sales data where we know the price and some observed characteristics.
- Build a decision rule that predicts price as a function of the observed characteristics.

What characteristics do we use?

We have to define the variables of interest and develop a specific quantitative measure of these variables

- Many factors or variables affect the price of a house
 - size
 - number of baths
 - garage
 - neighborhood
 - ► ...

To keep things super simple, let's focus only on size.

The value that we seek to predict is called the dependent (or output) variable, and we denote this:

► *Y*, e.g. the price of the house (thousands of dollars)

The variable that we use to aid in prediction is the independent, explanatory, or input variable, and this is labelled

► X, e.g. the size of house (thousands of square feet)

What does this data look like?

Size		Price	
	0.80		70
	0.90		83
	1.00		74
	1.10		93
	1.40		89
	1.40		58
	1.50		85
	1.60		114
	1.80		95
	2.00		100
	2.40		138
	2.50		111
	2.70		124
	3.20		161
	3.50		172

It is much more useful to look at a scatterplot:

plot(Price ~ Size, data = housing)



In other words, view the data as points in the $X \times Y$ plane.

Appears to be a linear relationship between price and size:

As size goes up, price goes up.



The line shown was fit by the "eyeball" method.

Recall that the equation of a line is:

$$Y = b_0 + b_1 X$$

Where b_0 is the intercept and b_1 is the slope.

The intercept value is in units of Y (\$1,000). The slope is in units of Y per units of X (\$1,000/1,000 sq ft).



Our "eyeball" line has $b_0 = 35$, $b_1 = 40$.

Can we do better than the eyeball method?

We desire a strategy for estimating the slope and intercept parameters in the model $\hat{Y} = b_0 + b_1 X$

A reasonable way to fit a line is to minimize the amount by which the fitted value differs from the actual value.

This amount is called the residual.

What is the "fitted value"?



The dots are the observed values and the line represents our fitted values given by $\hat{Y} = b_0 + b_1 X$.

What is the "residual"' for the *i*th observation'?



We can write $Y_i = \hat{Y}_i + (Y_i - \hat{Y}_i) = \hat{Y}_i + e_i$.

Least Squares

Ideally we want to minimize the size of all residuals:

- If they were all zero we would have a perfect line.
- Trade-off between moving closer to some points and at the same time moving away from other points.

The line fitting process:

- Take each residual e_i and assign it a weight e_i². Bigger residuals = bigger "mistakes" = higher weights
- Minimize the total of these weights to get best possible fit.

Least Squares chooses b_0 and b_1 to minimize $\sum_{i=1}^{N} e_i^2$

$$\sum_{i=1}^{N} e_i^2 = e_1^2 + e_2^2 + \dots + e_N^2 = (Y_1 - \hat{Y}_1)^2 + (Y_2 - \hat{Y}_2)^2 + \dots + (Y_N - \hat{Y}_N)^2$$

Least Squares

LS chooses a different line from ours:

- $b_0 = 38.88$ and $b_1 = 35.39$
- ▶ What do *b*₀ and *b*₁ mean again?



Least Squares in R

The 1m command fits linear (regression) models

```
fit = lm(Price ~ Size, data = housing)
print(fit)
##
## Call:
## lm(formula = Price ~ Size, data = housing)
##
## Coefficients:
## (Intercept)
                      Size
        38.88
                     35.39
##
```

```
fit = lm(Price ~ Size, data = housing)
summary(fit)
##
## Call:
## lm(formula = Price ~ Size, data = housing)
##
## Residuals:
      Min 10 Median 30 Max
##
## -30.425 -8.618 0.575 10.766 18.498
##
## Coefficients:
##
      Estimate Std. Error t value Pr(>|t|)
## (Intercept) 38.885 9.094 4.276 0.000903 ***
## Size 35.386 4.494 7.874 2.66e-06 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.14 on 13 degrees of freedom
## Multiple R-squared: 0.8267, Adjusted R-squared: 0.8133
## F-statistic: 62 on 1 and 13 DF, p-value: 2.66e-06
```

2nd Example: Offensive Performance in Baseball

1. Problems:

- Evaluate/compare traditional measures of offensive performance
- Help evaluate the worth of a player
- 2. Solutions:
 - Compare prediction rules that forecast runs as a function of either AVG (batting average), SLG (slugging percentage) or OBP (on base percentage)

2nd Example: Offensive Performance in Baseball



Baseball Data - Using AVG

Each observation corresponds to a team in MLB. Each quantity is

the average over a season.



• Y = runs per game; X = AVG (batting average)

LS fit: Runs/Game = -3.93 + 33.57 AVG

Baseball Data – Using SLG



- ► Y = runs per game
- X = SLG (slugging percentage)

LS fit: Runs/Game = -2.52 + 17.54 SLG

Baseball Data – Using OBP



- ► Y = runs per game
- ► X = OBP (on base percentage)

LS fit: Runs/Game = -7.78 + 37.46 OBP

Baseball Data

- What is the best prediction rule?
- Let's compare the predictive ability of each model using the average squared error

$$\frac{1}{n}\sum_{i=1}^{n}e_{i}^{2}=\frac{\sum_{i=1}^{n}\left(\widehat{Y}_{i}-Y_{i}\right)^{2}}{n}$$

Place your Money on OBP!

Average Squared Error			
AVG	0.083		
SLG	0.055		
OBP	0.026		



 $\hat{Y}_{n+1} = b_0 + b_1 x_{n+1}$

- b₀ is the intercept and b₁ is the slope
- ▶ We find *b*⁰ and *b*¹ using *Least Squares*
- ► For a new value of the independent variable OBP (say x_{n+1}) we can predict the response Y_{n+1} using the fitted line

From now on, terms "fitted values" (\hat{Y}_i) and "residuals" (e_i) refer to those obtained from the least squares line.

The fitted values and residuals have some special properties...

The Fitted Values and X

plot(predict(fit) ~ Size, data = housing, ylab = "fitted va")



cor(predict(fit), housing\$Size)

[1] 1

The Residuals and X plot(resid(fit) ~ Size, data = housing, ylab = "residuals")



mean(resid(fit)); cor(resid(fit), housing\$Size)

```
## [1] -9.633498e-17
```

```
## [1] 2.120636e-17
```

```
(i.e., zero). What's going on here?
```

A Deeper Look at Least Squares Estimates

Least squares estimates have some special properties:

- The fitted values \hat{Y} and x were **very** dependent
- The residuals $Y \hat{Y}$ and x had no apparent relationship
- The residuals $Y \hat{Y}$ had a sample mean of zero

What's going on? And what exactly are the least squares estimates?

We need to review sample covariance and correlation

Covariance

Measure the *direction* and *strength* of the linear relationship between Y and X



Correlation is the standardized covariance:

$$\operatorname{corr}(X,Y) = \frac{\operatorname{cov}(X,Y)}{\sqrt{s_x^2 s_y^2}} = \frac{\operatorname{cov}(X,Y)}{s_x s_y}$$

The correlation is scale invariant and the units of measurement don't matter: It is always true that $-1 \leq \operatorname{corr}(X, Y) \leq 1$.

This gives the direction (- or +) and strength (0 \rightarrow 1 in absolute value)

of the linear relationship between X and Y.





34

Only measures linear relationships:

corr(X, Y) = 0 does not mean the variables are not related!



Also be careful with influential observations...

The Least Squares Estimates

The values for b_0 and b_1 that minimize the least squares criterion

are:

$$b_1 = r_{xy} imes rac{s_y}{s_x}$$
 $b_0 = ar{Y} - b_1 ar{X}$

where,

- \bar{X} and \bar{Y} are the sample mean of X and Y
- $corr(x, y) = r_{xy}$ is the sample correlation
- s_x and s_y are the sample standard deviation of X and Y

These are the **least squares estimates** of β_0 and β_1 .

The Least Squares Estimates

The values for b_0 and b_1 that minimize the least squares criterion

are:

$$b_1 = r_{xy} imes rac{s_y}{s_x}$$
 $b_0 = ar{Y} - b_1 ar{X}$

How do we interpret these?

- b_0 ensures the line goes through (\bar{x}, \bar{y})
- ▶ b₁ scales the correlation to appropriate units by multiplying with s_y/s_x (what are the units of b₁?)

```
# Computing least squares estimates "by hand"
y = housing$Price; x = housing$Size
rxy = cor(y, x)
sx = sd(x)
sy = sd(y)
ybar = mean(y)
xbar = mean(x)
b1 = rxy*sy/sx
b0 = ybar - b1*xbar
print(b0); print(b1)
## [1] 38.88468
## [1] 35.38596
```

```
# We get the same result as lm()
fit = lm(Price~Size, data=housing)
print(fit)
##
## Call:
## lm(formula = Price ~ Size, data = housing)
##
## Coefficients:
## (Intercept)
                      Size
        38.88
               35.39
##
```

Properties of Least Squares Estimates

Remember from the housing data, we had:

- $corr(\hat{Y}, x) = 1$ (a perfect linear relationship)
- corr(e, x) = 0 (no linear relationship)
- mean(e) = 0 (sample average of residuals is zero)

Why?

What is the intuition for the relationship between \hat{Y} and e and X? Lets consider some "crazy" alternative line:



Fitted Values and Residuals

This is a bad fit! We are underestimating the value of small houses and overestimating the value of big houses.



Clearly, we have left some predictive ability on the table!

Summary: LS is the best we can do!!

As long as the correlation between e and X is non-zero, we could always adjust our prediction rule to do better.

We need to exploit all of the predictive power in the X values and put this into \hat{Y} , leaving no "Xness" in the residuals.

In Summary: $Y = \hat{Y} + e$ where:

- \hat{Y} is "made from X" using a linear equation; $\operatorname{corr}(X, \hat{Y}) = \pm 1$.
- e has no linear relationship with X; $\operatorname{corr}(X, e) = 0$.
- On average (over the sample), our prediction error is zero: $\bar{e} = \sum_{i=1}^{n} e_i = 0.$

Recall: Variability/variance is a measure of risk, or unpredictability.

We'd like to use information in one variable (X) to help predict another (Y)

We can quantify how much the total error or variance goes down after using information in X...

Using a "good" X reduces variability in Y...



Using a "bad" X doesn't

When X has low predictive power, the story is different:

House price (Y) vs. the number of stop signs within a two block radius of a house (X).

See that in this case, the marginal and the Conditionals are not that different!



Remember that $Y = \hat{Y} + e$

Since \hat{Y} and e are uncorrelated, i.e. $\operatorname{corr}(\hat{Y}, e) = 0$,

$$\operatorname{var}(Y) = \operatorname{var}(\hat{Y} + e) = \operatorname{var}(\hat{Y}) + \operatorname{var}(e)$$
$$\frac{\sum_{i=1}^{n} (Y_i - \bar{Y})^2}{n-1} = \frac{\sum_{i=1}^{n} (\hat{Y}_i - \bar{\hat{Y}})^2}{n-1} + \frac{\sum_{i=1}^{n} (e_i - \bar{e})^2}{n-1}$$

Given that $\bar{e} = 0$, and the sample mean of the fitted values $\hat{Y} = \bar{Y}$ (why?) we get to write:

$$\sum_{i=1}^{n} (Y_i - \bar{Y})^2 = \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^{n} e_i^2$$



- SST is measuring total variation in Y/total error in Y using the simplest prediction \overline{Y} – i.e., no info about X
- SSR is measuring predictable (via our regression model) variation in Y – how much our predictions change after accounting for linear effects of X
- SSE is measuring left over, unpredictable variation in Y



Things to note:

- SST is fixed, so as SSR increases, SSE (the total error in our predictions) goes down.
- SSR describes variation that's predictable by a linear equation of X. We could get better SSR (and lower SSE) with nonlinear functions of X, but we have to be careful – more soon.

$$egin{array}{rcl} (Y_i - ar{Y}) &=& \hat{Y}_i + e_i - ar{Y} \ &=& (\hat{Y}_i - ar{Y}) + e_i \end{array}$$



Y

The Coefficient of Determination R^2

The coefficient of determination, denoted by R^2 , measures how well the fitted values \hat{Y} follow Y:

$$R^2 = rac{\mathrm{SSR}}{\mathrm{SST}} = 1 - rac{\mathrm{SSE}}{\mathrm{SST}}$$

- R² is often called the proportion of variance in Y that is "explained" by the regression line (in the mathematical – not scientific – sense!): R² = 1 - Var(e)/Var(Y)
- $0 < R^2 < 1$
- ► For simple linear regression, R² = r²_{xy}. Similar caveats to sample correlation apply!

Explanations and predictions

A better way to think about R^2 is as the proportion of variability – i.e. unpredictablility – in Y that becomes predictable when using X in a linear regression model.

 R^2 does not tell you:

- Whether there is/is not any causal relationship between X and Y (Question: What is the R² from regressing X on Y)?
- Whether your regression model is a reasonable approximation of reality
- Whether your model generalizes well outside your sample

R^2 for the Housing Data

```
summary(fit)
##
## Call:
## lm(formula = Price ~ Size, data = housing)
##
## Residuals:
      Min 10 Median 30 Max
##
## -30.425 -8.618 0.575 10.766 18.498
##
## Coefficients:
             Estimate Std. Error t value Pr(>|t|)
##
## (Intercept) 38.885 9.094 4.276 0.000903 ***
## Size 35.386 4.494 7.874 2.66e-06 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.14 on 13 degrees of freedom
## Multiple R-squared: 0.8267.Adjusted R-squared: 0.8133
```

R^2 for the Housing Data

summary(fit) ## ## Call: ## lm(formula = Price ~ Size, data = housing) ## ## Residuals: ## Min 1Q Median 3Q Max ## -30.425 -8.618 0.575 10.766 18.498 ## ## Coefficients: ## Estimate Std. Error t value Pr(>|t|) ## (Intercept) 38.885 9.094 4.276 0.000903 *** ## Size 35.386 4.494 7.874 2.66e-06 *** ## ---## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 ## ## Residual standard error: 14.14 on 13 degrees of freedom ## Multiple R-squared: 0.8267, Adjusted R-squared: 0.8133 ## F-statistic: 62 on 1 and 13 DF. p-value: 2.66e-06

R^2 for the Housing Data

```
anova(fit)
## Analysis of Variance Table
##
## Response: Price
##
            Df Sum Sq Mean Sq F value Pr(>F)
## Size 1 12393.1 12393.1 61.998 2.66e-06 ***
## Residuals 13 2598.6 199.9
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

$$R^2 = \frac{SSR}{SST} = \frac{12393.1}{2598.6 + 12393.1} = 0.8267$$

Three very similar, related ways to look at a simple linear regression... with only one X variable, life is easy!

	R^2	corr	SSE
OBP	0.88	0.94	0.79
SLG	0.76	0.87	1.64
AVG	0.63	0.79	2.49