



Division of
Statistics + Scientific Computation

THE UNIVERSITY OF TEXAS AT AUSTIN

Advanced Regression
Summer Statistics Institute

Day 1: Introduction to Linear Regression

Course Overview

Day 1: Intro and Simple Regression Model

Day 2: Multiple Regression, Dummy Variables and Interactions

Day 3: Transformations, Non-linear models and Model Selection

Day 4: Time Series, Logistic Regression and more...

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 variable

Regression in Business

- ▶ Optimal portfolio choice:
 - **Predict** future joint distribution of asset returns
 - **Construct** optimal portfolio (choose weights)
- ▶ Determining price and marketing strategy:
 - **Estimate** the effect of price and advertisement on sales
 - **Decide** what is optimal price and ad campaign
- ▶ Credit scoring model:
 - **Predict** future probability of default using known characteristics of borrower
 - **Decide** whether or not to lend (and if so, how much)

Why?

Straight prediction questions:

- ▶ For how much will my house sell?
- ▶ How many runs per game will the Red Sox score in 2011?
- ▶ Will this person like that movie? (Netflix prize)

Explanation and understanding:

- ▶ What is the impact of MBA on income?
- ▶ How do the returns of a mutual fund relate to the market?
- ▶ Does Walmart discriminate against women regarding salaries?

1st 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.

Predicting House Prices

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, air conditioning, etc
 - ▶ neighborhood
- ▶ Easy to quantify price and size but what about other variables such as aesthetics, workmanship, etc?

Predicting House Prices

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:

- ▶ $Y =$ price of house (e.g. thousands of dollars)

The variable that we use to guide prediction is the **explanatory (or input)** variable, and this is labelled

- ▶ $X =$ size of house (e.g. thousands of square feet)

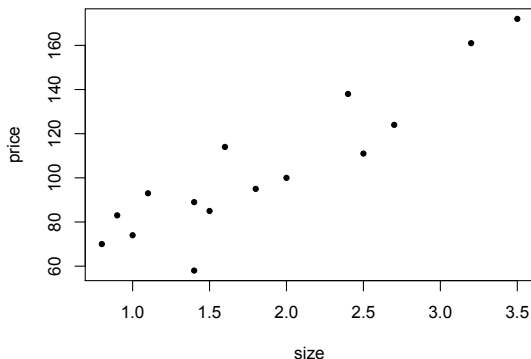
Predicting House Prices

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

Predicting House Prices

It is much more useful to look at a scatterplot



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

Regression Model

Y = response or outcome variable

$X_1, X_2, X_3, \dots, X_p$ = explanatory or input variables

The general relationship approximated by:

$$Y = f(X_1, X_2, \dots, X_p) + e$$

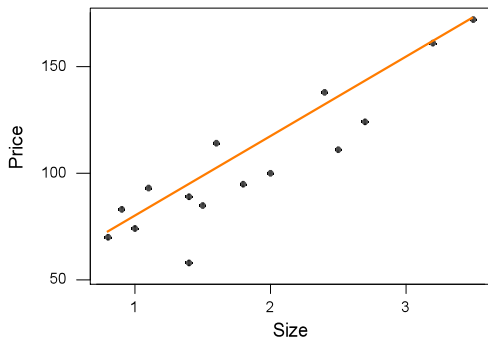
And a linear relationship is written

$$Y = b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p + e$$

Linear Prediction

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.

Linear Prediction

Recall that the equation of a line is:

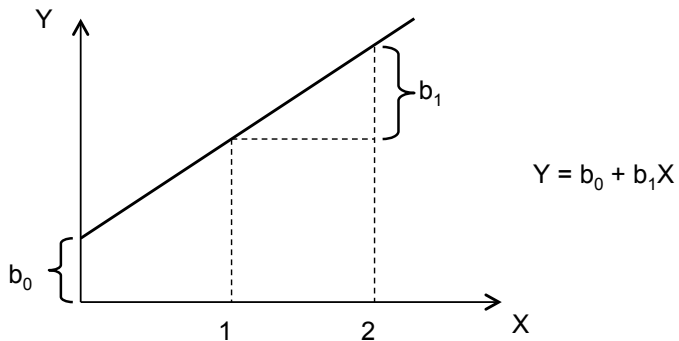
$$Y = b_0 + b_1X$$

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).

Linear Prediction



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

Linear Prediction

We can now predict the price of a house when we know only the size; just read the value off the line that we've drawn.

For example, given a house with of size $X = 2.2$.

Predicted price $\hat{Y} = 35 + 40(2.2) = 123$.

Note: Conversion from 1,000 sq ft to \$1,000 is done for us by the slope coefficient (b_1)

Linear Prediction

Can we do better than the eyeball method?

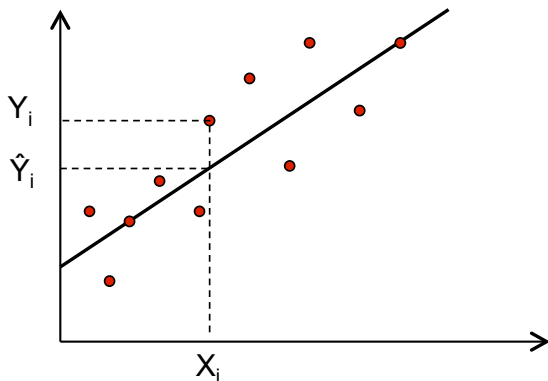
We desire a strategy for estimating the slope and intercept parameters in the model $\hat{Y} = b_0 + b_1X$

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**.

Linear Prediction

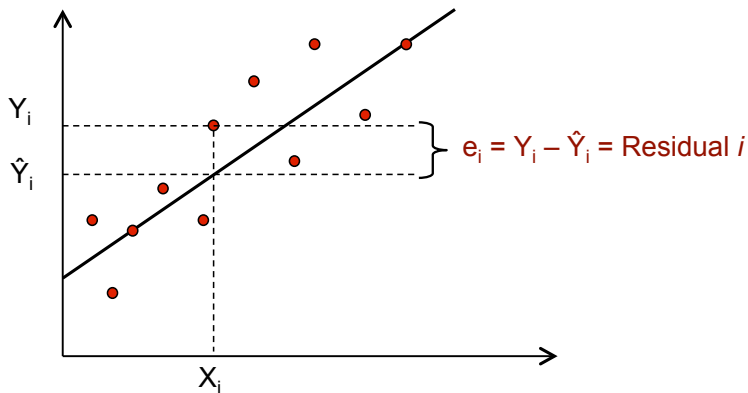
What is the “fitted value”?



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

Linear Prediction

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.

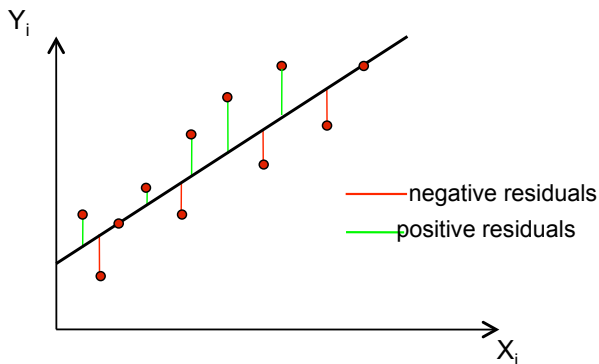
The line fitting process:

- ▶ Give weights to all of the residuals.
- ▶ Minimize the “total” of residuals to get best 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



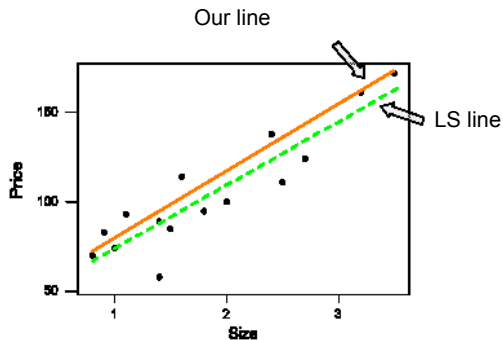
Choose the line to minimize the sum of the squares of the residuals,

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - [b_0 + b_1 X_i])^2$$

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?



Eyeball vs. LS Residuals

- ▶ eyeball: $b_0 = 35$, $b_1 = 40$
- ▶ LS: $b_0 = 38.88$, $b_1 = 35.39$

Size	Price	yhat-eyeball	yhat-LS	e-eyeball	e-LS	e2-eyeball	e2-LS
0.80	70	67	67.19	3.00	2.81	9.00	7.88
0.90	83	71	70.73	12.00	12.27	144.00	150.51
1.00	74	75	74.27	-1.00	-0.27	1.00	0.07
1.10	93	79	77.81	14.00	15.19	196.00	230.76
1.40	89	91	88.42	-2.00	0.58	4.00	0.33
1.40	58	91	88.42	-33.00	-30.42	1089.00	925.67
1.50	85	95	91.96	-10.00	-6.96	100.00	48.49
1.60	114	99	95.50	15.00	18.50	225.00	342.17
1.80	95	107	102.58	-12.00	-7.58	144.00	57.44
2.00	100	115	109.66	-15.00	-9.66	225.00	93.25
2.40	138	131	123.81	7.00	14.19	49.00	201.33
2.50	111	135	127.35	-24.00	-16.35	576.00	267.30
2.70	124	143	134.43	-19.00	-10.43	361.00	108.71
3.20	161	163	152.12	-2.00	8.88	4.00	78.86
3.50	172	175	162.74	-3.00	9.26	9.00	85.84
sum				-70.00	0.00	3136.00	2598.63

Least Squares – Excel Output

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.909209967
R Square	0.826662764
Adjusted R Square	0.81332913
Standard Error	14.13839732
Observations	15

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	12393.10771	12393.10771	61.99831126	2.65987E-06
Residual	13	2598.625623	199.8942787		
Total	14	14991.73333			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	38.88468274	9.09390389	4.275906499	0.000902712	19.23849785	58.53086763
Size	35.38596255	4.494082942	7.873900638	2.65987E-06	25.67708664	45.09483846

Excel Break...

- ▶ Scatterplots, linear function
- ▶ Regression

2nd Example: Offensive Performance in Baseball

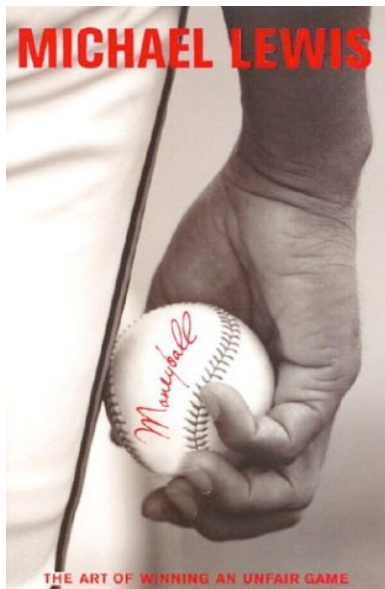
1. Problems:

- ▶ Evaluate/compare traditional measures of offensive performance
- ▶ Help evaluate the worth of a player

2. Solutions:

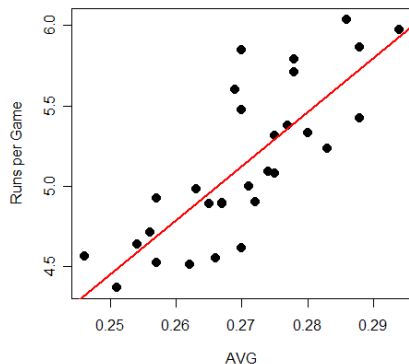
- ▶ Compare *prediction rules* that forecast runs as a function of either AVG (batting average), SLG (slugging percentage) or OBP (on base percentage)

2nd Example: Offensive Performance in Baseball



Baseball Data – Using AVG

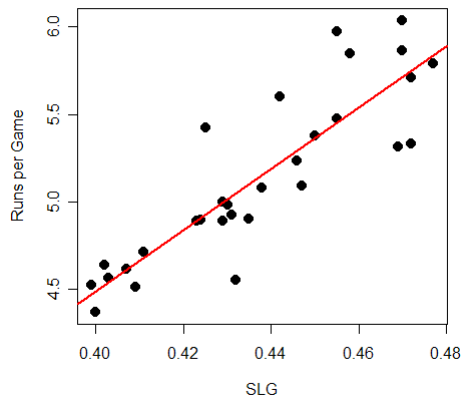
Each observation corresponds to a team in MLB. Each quantity is the average over a season.



► $Y = \text{runs per game}; X = \text{AVG (average)}$

LS fit: $\text{Runs/Game} = -3.93 + 33.57 \text{ AVG}$

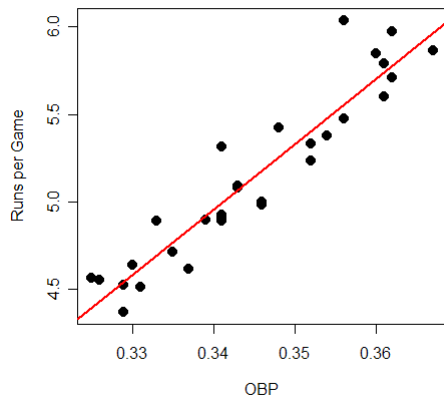
Baseball Data – Using SLG



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

LS fit: $\text{Runs/Game} = -2.52 + 17.54 \text{ SLG}$

Baseball Data – Using OBP



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

LS fit: $\text{Runs/Game} = -7.78 + 37.46 \text{ OBP}$

Baseball Data

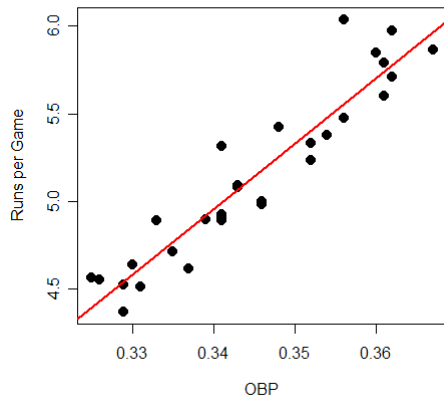
- ▶ What is the best prediction rule?
- ▶ Let's compare the predictive ability of each model using the average squared error

$$\frac{1}{N} \sum_{i=1}^N e_i^2 = \frac{\sum_{i=1}^N \left(\widehat{Runs}_i - Runs_i \right)^2}{N}$$

Place your Money on OBP!!!

Average Squared Error	
AVG	0.083
SLG	0.055
OBP	0.026

Linear Prediction



$$\hat{Y}_i = b_0 + b_1 X_i$$

- ▶ b_0 is the intercept and b_1 is the slope
- ▶ We find b_0 and b_1 using *Least Squares*

The Least Squares Criterion

The formulas for b_0 and b_1 that minimize the least squares criterion are:

$$b_1 = r_{xy} \times \frac{s_y}{s_x} \quad b_0 = \bar{Y} - b_1 \bar{X}$$

where,

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

Sample Mean and Sample Variance

- ▶ **Sample Mean:** measure of centrality

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

- ▶ **Sample Variance:** measure of spread

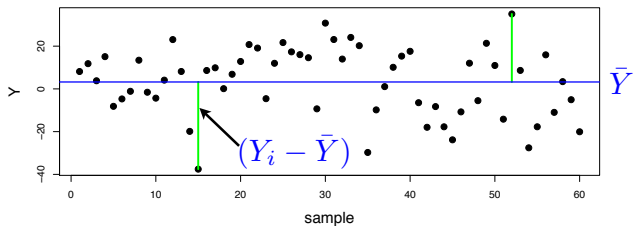
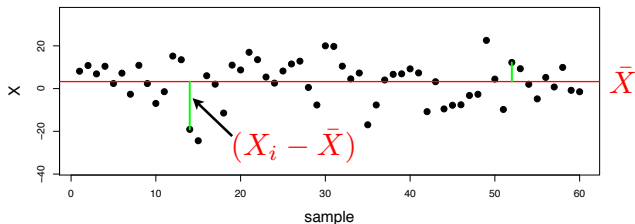
$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- ▶ **Sample Standard Deviation:**

$$s_y = \sqrt{s_y^2}$$

Example

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

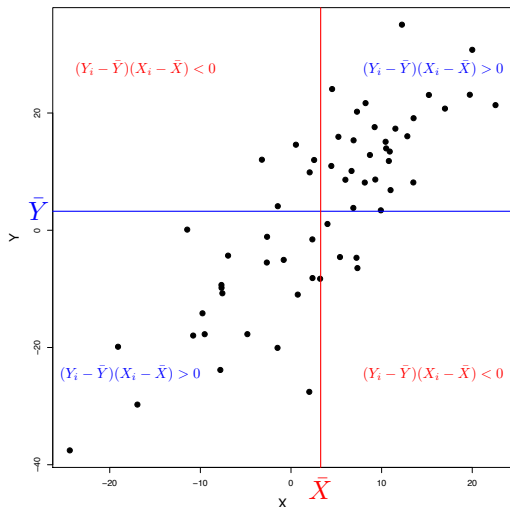


$$s_x = 9.7 \quad s_y = 16.0$$

Covariance

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

$$\text{Cov}(Y, X) = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{n - 1}$$



▶ $s_y = 15.98, s_x = 9.7$

▶ $\text{Cov}(X, Y) = 125.9$

How do we interpret that?

Correlation

Correlation is the standardized covariance:

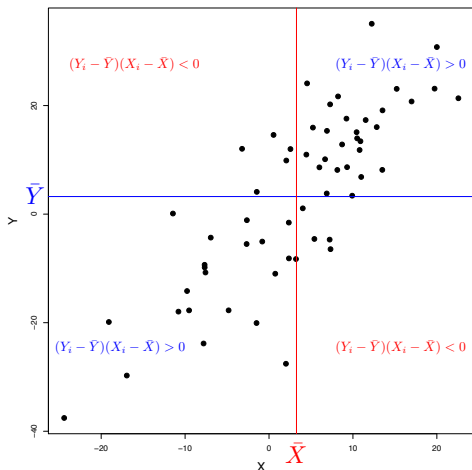
$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{s_x^2 s_y^2}} = \frac{\text{cov}(X, Y)}{s_x s_y}$$

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

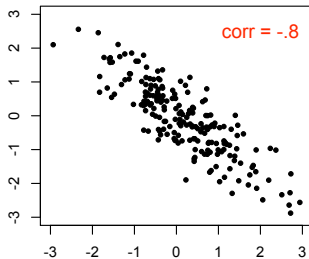
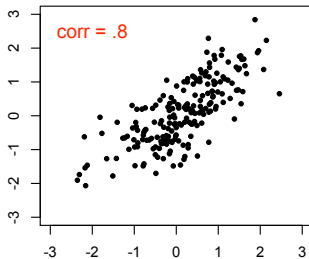
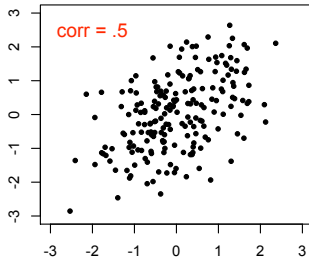
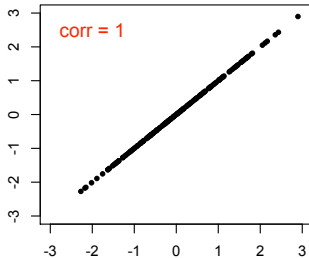
This gives the direction (- or +) and strength ($0 \rightarrow 1$) of the **linear** relationship between X and Y .

Correlation

$$\text{corr}(Y, X) = \frac{\text{cov}(X, Y)}{\sqrt{s_x^2 s_y^2}} = \frac{\text{cov}(X, Y)}{s_x s_y} = \frac{125.9}{15.98 \times 9.7} = 0.812$$



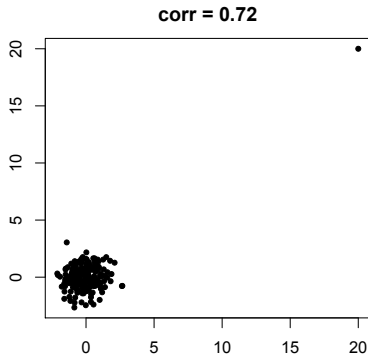
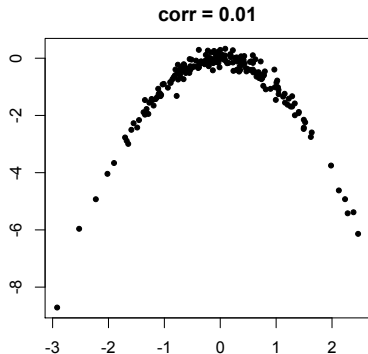
Correlation



Correlation

Only measures **linear** relationships:

$\text{corr}(X, Y) = 0$ does not mean the variables are not related!



Also be careful with influential observations. **Excel Break:** correl, stdev,...

Back to Least Squares

1. Intercept:

$$b_0 = \bar{Y} - b_1\bar{X} \Rightarrow \bar{Y} = b_0 + b_1\bar{X}$$

- ▶ The point (\bar{X}, \bar{Y}) is on the regression line!
- ▶ Least squares finds the point of means and rotate the line through that point until getting the “right” slope

2. Slope:

$$\begin{aligned} b_1 &= \text{corr}(X, Y) \times \frac{s_Y}{s_X} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\text{Cov}(X, Y)}{\text{var}(X)} \end{aligned}$$

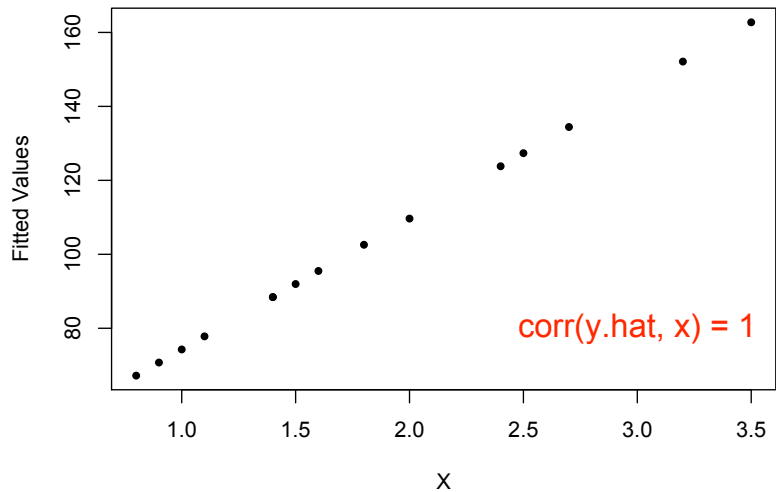
- ▶ So, the right slope is the *correlation coefficient* times a *scaling factor* that ensures the proper units for b_1

More on Least Squares

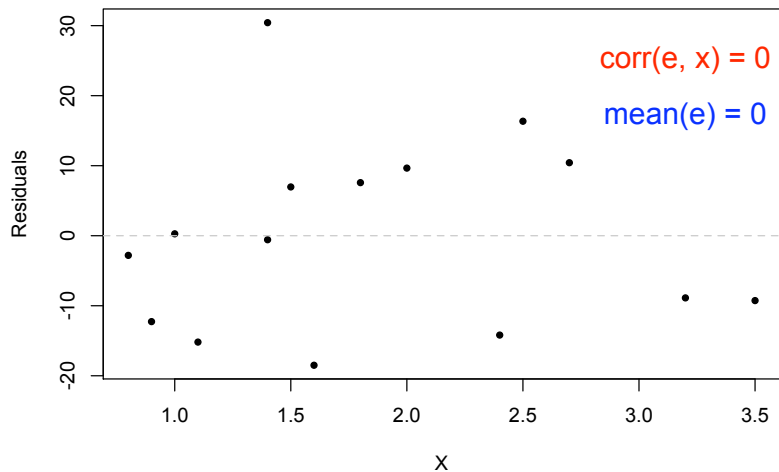
From now on, terms “fitted values” (\hat{Y}_i) and “residuals” (e_i) refer to those obtained from the least squares line.

The fitted values and residuals have some special properties. Lets look at the housing data analysis to figure out what these properties are...

The Fitted Values and X



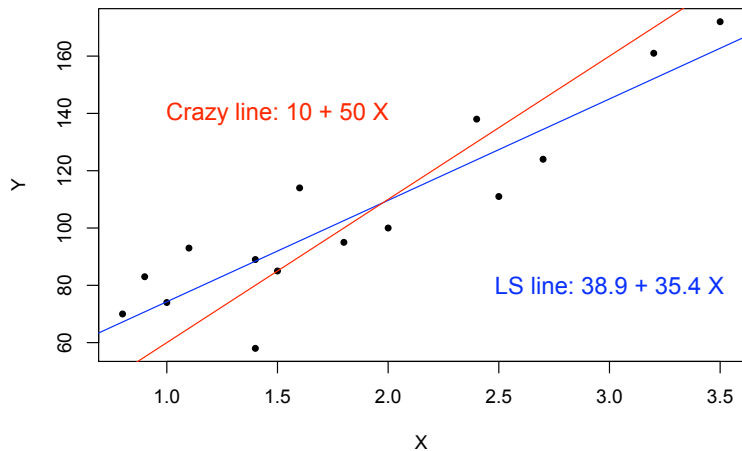
The Residuals and X



Why?

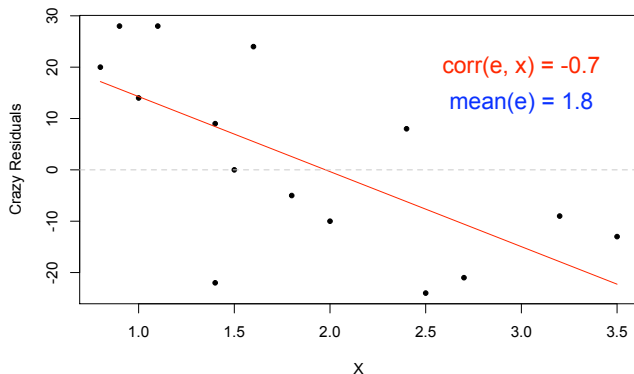
What is the intuition for the relationship between \hat{Y} and e and X ?

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!

Fitted Values and Residuals

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 “ X ness” in the residuals.

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

- ▶ \hat{Y} is “made from X ”; $\text{corr}(X, \hat{Y}) = 1$.
- ▶ e is unrelated to X ; $\text{corr}(X, e) = 0$.

Another way to derive things (Optional)

The intercept:

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n e_i = 0 &\Rightarrow \frac{1}{n} \sum_{i=1}^n (Y_i - b_0 - b_1 X_i) \\ &\Rightarrow \bar{Y} - b_0 - b_1 \bar{X} = 0 \\ &\Rightarrow b_0 = \bar{Y} - b_1 \bar{X}\end{aligned}$$

Another way to derive things (Optional)

The slope:

$$\begin{aligned}\text{corr}(e, X) &= \sum_{i=1}^n e_i(X_i - \bar{X}) = 0 \\ &= \sum_{i=1}^n (Y_i - b_0 - b_1 X_i)(X_i - \bar{X}) \\ &= \sum_{i=1}^n (Y_i - \bar{Y} - b_1(X_i - \bar{X}))(X_i - \bar{X}) \\ \Rightarrow b_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = r_{xy} \frac{s_y}{s_x}\end{aligned}$$

Decomposing the Variance

How well does the least squares line explain variation in Y ?

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

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

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

Given that $\bar{e} = 0$, and $\bar{\hat{Y}} = \bar{Y}$ (why?) we get to:

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

Decomposing the Variance

$$\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{\substack{\text{Total Sum of} \\ \text{Squares} \\ \text{SST}}} = \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{\substack{\text{Regression SS} \\ \text{SSR}}} + \underbrace{\sum_{i=1}^n e_i^2}_{\substack{\text{Error SS} \\ \text{SSE}}}$$

SSR: Variation in Y explained by the regression line.

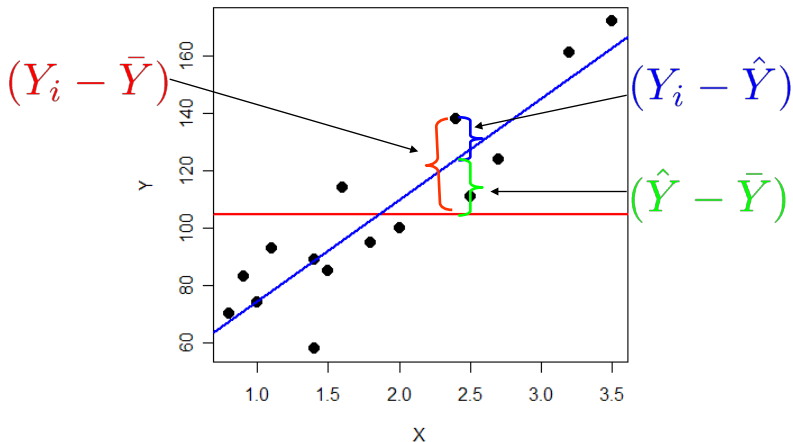
SSE: Variation in Y that is left unexplained.

$$\text{SSR} = \text{SST} \Rightarrow \text{perfect fit.}$$

Be careful of similar acronyms; e.g. SSR for “residual” SS.

Decomposing the Variance

$$\begin{aligned}(Y_i - \bar{Y}) &= \hat{Y}_i + e_i - \bar{Y} \\ &= (\hat{Y}_i - \bar{Y}) + e_i\end{aligned}$$



A Goodness of Fit Measure: R^2

The **coefficient of determination**, denoted by R^2 , measures goodness of fit:

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

- ▶ $0 < R^2 < 1$.
- ▶ The closer R^2 is to 1, the better the fit.

A Goodness of Fit Measure: R^2 (Optional)

An interesting fact: $R^2 = r_{xy}^2$ (i.e., R^2 is squared correlation).

$$\begin{aligned}R^2 &= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\&= \frac{\sum_{i=1}^n (b_0 + b_1 X_i - b_0 - b_1 \bar{X})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\&= \frac{b_1^2 \sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{b_1^2 s_x^2}{s_y^2} = r_{xy}^2\end{aligned}$$

No surprise: the higher the sample correlation between X and Y , the better you are doing in your regression.

Back to the House Data

SUMMARY OUTPUT

Regression Statistics	
Multiple R	0.909209967
R Square	0.826662764
Adjusted R Square	0.81332913
Standard Error	14.13839732
Observations	15

ANOVA					
	df	SS	MS	F	Significance F
Regression	1	12393.10771	12393.10771	61.99831126	2.65987E-06
Residual	13	2598.625623	199.8942787		
Total	14	14991.73333			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	38.86468274	9.09390389	4.275906499	0.000902712	19.23849785	58.53086763	19.23849785	58.53086763
X Variable 1	35.38596255	4.494082942	7.873900638	2.65987E-06	25.67708664	45.09483846	25.67708664	45.09483846

SSR

SST

SSE

$$R^2 = \frac{SSR}{SST} = 0.82 = \frac{12395}{14991}$$

Back to Baseball

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

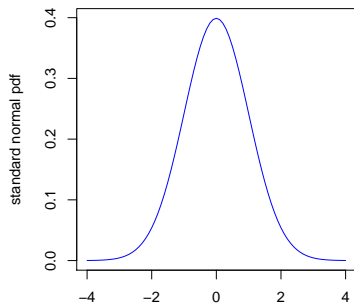
Review Break...

Before moving on we need to review the NORMAL distribution

The Normal Distribution



- ▶ A random variable is a number we are NOT sure about but we might have some idea of how to describe its potential outcomes. The Normal distribution is the most used probability distribution to describe a random variable
- ▶ The probability the number ends up in an interval is given by the area under the curve (pdf)

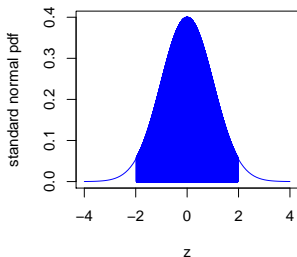
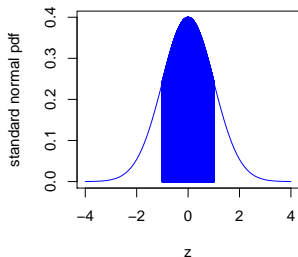


The Normal Distribution

- ▶ The standard Normal distribution has mean 0 and has variance 1.
- ▶ **Notation:** If $Z \sim N(0, 1)$ (Z is the random variable)

$$\Pr(-1 < Z < 1) = 0.68$$

$$\Pr(-1.96 < Z < 1.96) = 0.95$$



The Normal Distribution

Note:

For simplicity we will often use $P(-2 < Z < 2) \approx 0.95$

Questions:

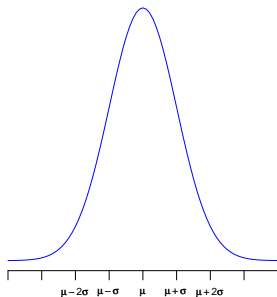
- ▶ What is $Pr(Z < 2)$? How about $Pr(Z \leq 2)$?
- ▶ What is $Pr(Z < 0)$?

The Normal Distribution

- ▶ The standard normal is not that useful by itself. When we say “the normal distribution”, we really mean a family of distributions.
- ▶ We obtain pdfs in the normal family by shifting the bell curve around and spreading it out (or tightening it up).

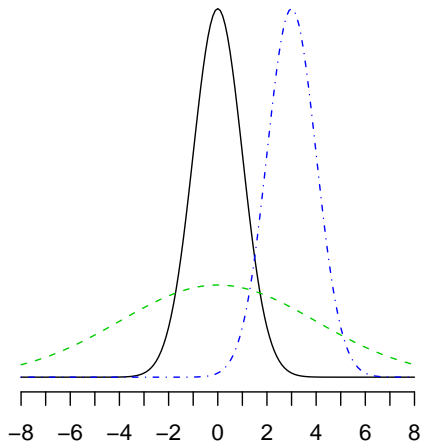
The Normal Distribution

- ▶ We write $X \sim N(\mu, \sigma^2)$. “Normal distribution with mean μ and variance σ^2 .”
- ▶ The parameter μ determines where the curve is. The center of the curve is μ .
- ▶ The parameter σ determines how spread out the curve is. The area under the curve in the interval $(\mu - 2\sigma, \mu + 2\sigma)$ is 95%.
 $Pr(\mu - 2\sigma < X < \mu + 2\sigma) \approx 0.95$



The Normal Distribution

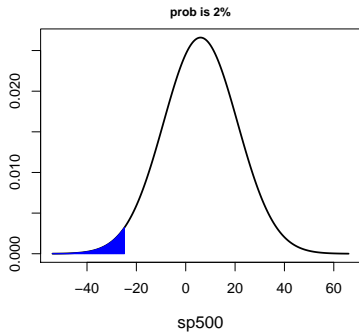
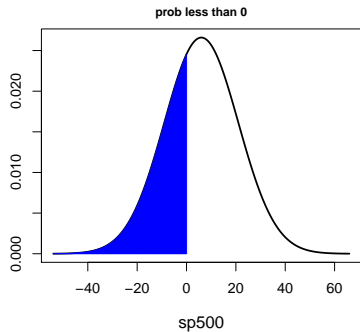
- ▶ **Example:** Below are the pdfs of $X_1 \sim N(0, 1)$, $X_2 \sim N(3, 1)$, and $X_3 \sim N(0, 16)$.
- ▶ Which pdf goes with which X ?



The Normal Distribution – Example

- ▶ Assume the annual returns on the SP500 are normally distributed with mean 6% and standard deviation 15%.
SP500 $\sim N(6, 225)$. (Notice: $15^2 = 225$).
- ▶ Two questions: (i) What is the chance of losing money on a given year? (ii) What is the value that there's only a 2% chance of losing that or more?
- ▶ Lloyd Blankfein: *"I spend 98% of my time thinking about 2% probability events!"*
- ▶ (i) $Pr(SP500 < 0)$ and (ii) $Pr(SP500 < ?) = 0.02$

The Normal Distribution – Example



- ▶ (i) $Pr(SP500 < 0) = 0.35$ and (ii) $Pr(SP500 < -25) = 0.02$
- ▶ In Excel: **NORMDIST** and **NORMINV** (homework!)

The Normal Distribution

1. Note: In

$$X \sim N(\mu, \sigma^2)$$

μ is the mean and σ^2 is the variance.

2. Standardization: if $X \sim N(\mu, \sigma^2)$ then

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

3. Summary:

$$X \sim N(\mu, \sigma^2):$$

μ : where the curve is

σ : how spread out the curve is

95% chance $X \in \mu \pm 2\sigma$.

The Normal Distribution – Another Example

Prior to the 1987 crash, monthly S&P500 returns (r) followed (approximately) a normal with mean 0.012 and standard deviation equal to 0.043. **How extreme was the crash of -0.2176?** The standardization helps us interpret these numbers...

$$r \sim N(0.012, 0.043^2)$$

$$z = \frac{r - 0.012}{0.043} \sim N(0, 1)$$

For the crash,

$$z = \frac{-0.2176 - 0.012}{0.043} = -5.27$$

How extreme is this zvalue? **5 standard deviations away!!**

Prediction and the Modeling Goal

A prediction rule is any function where you input X and it outputs \hat{Y} as a predicted response at X .

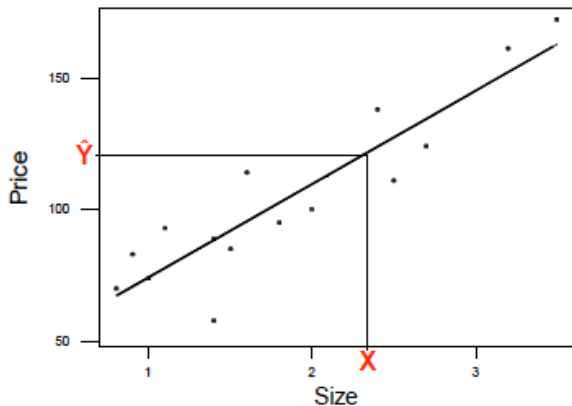
The least squares line is a prediction rule:

$$\hat{Y} = f(X) = b_0 + b_1X$$

Prediction and the Modeling Goal

\hat{Y} is not going to be a perfect prediction.

We need to devise a notion of **forecast accuracy**.



Prediction and the Modeling Goal

There are two things that we want to know:

- ▶ What value of Y can we expect for a given X ?
- ▶ How sure are we about this forecast? Or how different could Y be from what we expect?

Our goal is to measure the accuracy of our forecasts or **how much uncertainty there is in the forecast**. One method is to specify a range of Y values that are likely, given an X value.

Prediction Interval: probable range for Y -values given X

Prediction and the Modeling Goal

Key Insight: To construct a prediction interval, we will have to assess the likely range of error values corresponding to a Y value that has not yet been observed!

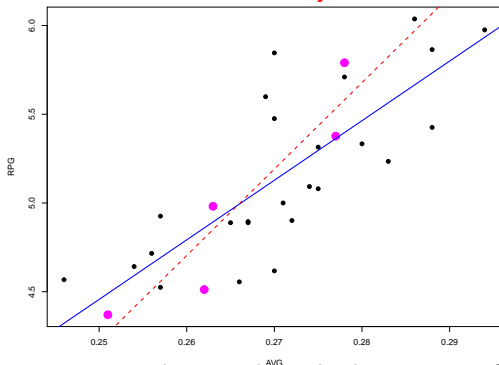
We will build a **probability model** (e.g., normal distribution).

Then we can say something like “with 95% probability the error will be no less than -\$28,000 or larger than \$28,000” .

We must also acknowledge that the “fitted” line may be fooled by particular realizations of the residuals.

Prediction and the Modeling Goal

- ▶ Suppose you only had the purple points in the graph. The dashed line fits the purple points. The solid line fits all the points. **Which line is better? Why?**



- ▶ In summary, we need to work with the notion of a “true line” and a probability distribution that describes deviation around the line.

The Simple Linear Regression Model

The power of statistical inference comes from the ability to make precise statements about the accuracy of the forecasts.

In order to do this we must invest in a **probability model**.

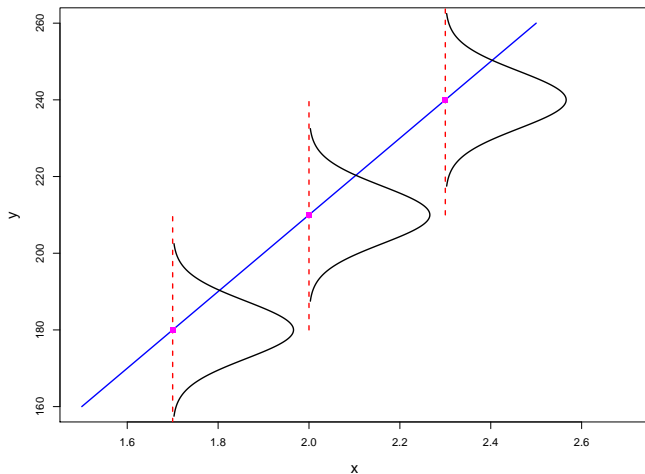
Simple Linear Regression Model: $Y = \beta_0 + \beta_1 X + \varepsilon$

$$\varepsilon \sim N(0, \sigma^2)$$

- ▶ $\beta_0 + \beta_1 X$ represents the “true line”; The part of Y that depends on X .
- ▶ The error term ε is independent “idiosyncratic noise”; The part of Y not associated with X .

The Simple Linear Regression Model

$$Y = \beta_0 + \beta_1 X + \varepsilon$$



The Simple Linear Regression Model – Example

You are told (without looking at the data) that

$$\beta_0 = 40; \beta_1 = 45; \sigma = 10$$

and you are asked to predict price of a 1500 square foot house.

What do you know about Y from the model?

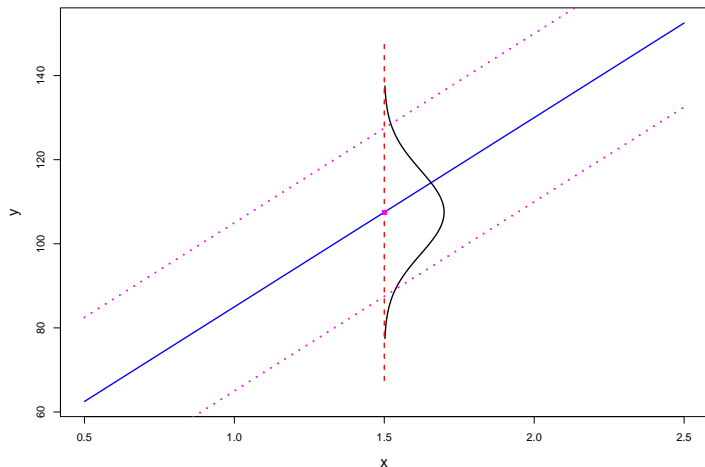
$$\begin{aligned} Y &= 40 + 45(1.5) + \varepsilon \\ &= 107.5 + \varepsilon \end{aligned}$$

Thus our prediction for price is $Y|X = 1.5 \sim N(107.5, 10^2)$

and a 95% *Prediction Interval* for Y is $87.5 < Y < 127.5$

Conditional Distributions

$$Y = \beta_0 + \beta_1 X + \varepsilon$$



The conditional distribution for Y given X is Normal:

$$Y|X = x \sim N(\beta_0 + \beta_1 x, \sigma^2).$$

Prediction Intervals (one more time!)

The model says that the mean value of a 1500 sq. ft. house is \$107,500 and that deviation from mean is within \approx \$20,000.

We are 95% sure that

- ▶ $-20 < \varepsilon < 20$
- ▶ $87.5 < Y < 127.5$

In general, the 95 % Prediction Interval is $PI = \beta_0 + \beta_1 X \pm 2\sigma$.

Conditional Distributions

Why do we have $\varepsilon \sim N(0, \sigma^2)$?

- ▶ $E[\varepsilon] = 0 \Leftrightarrow E[Y | X] = \beta_0 + \beta_1 X$
($E[Y | X]$ is “conditional expectation of Y given X ”).
- ▶ Many things are close to Normal (central limit theorem).
- ▶ It works! This is a very robust model for the world.

We can think of $\beta_0 + \beta_1 X$ as the “true” regression line.

Conditional Distributions

Regression models are really all about modeling the conditional distribution of Y given X .

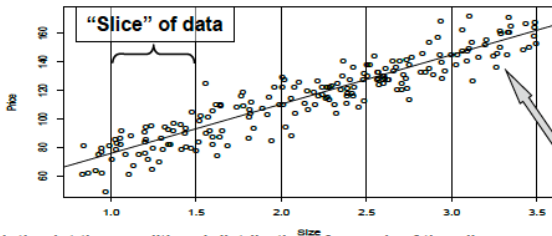
Why are conditional distributions important?

Given that I know X what kind of Y can I expect? Our model provides one way to think about this question.

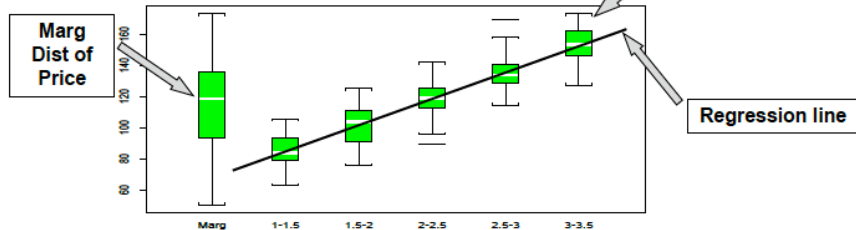
We can also look at this by “slicing” the cloud of points in the scatterplot to obtain the distribution of Y conditional on various ranges of X values.

Data Conditional Distribution vs Marginal Distribution

Let's consider a regression of house price on size:



Now let's plot the conditional distributions for each of the slices



Conditional Distribution and Marginal Distribution

Key Observations from these plots:

- ▶ Conditional distributions answer the forecasting problem: if I know that a house is between 1 and 1.5 1000 sq.ft., then the conditional distribution (second boxplot) gives me a point forecast (the mean) and prediction interval.
- ▶ The conditional means seem to line up along the regression line.
- ▶ The conditional distributions have much smaller dispersion than the marginal distribution.

Conditional Distribution vs Marginal Distribution

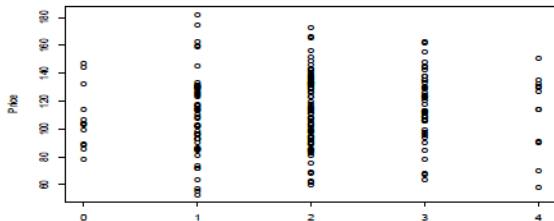
This suggests two general points:

- ▶ If X has no forecasting power, then the marginal and conditionals will be the same.
- ▶ If X has some forecasting information, then conditional means will be different than the marginal or overall mean and conditional standard deviation of Y given X will be less than the marginal standard deviation of Y .

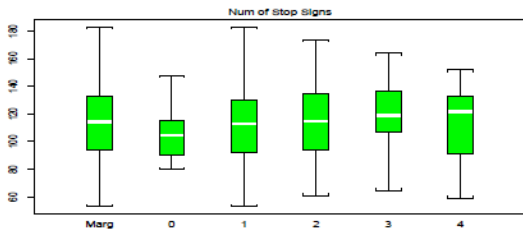
Conditional Distribution vs Marginal Distribution

Intuition from an example where X has no predictive power.

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!

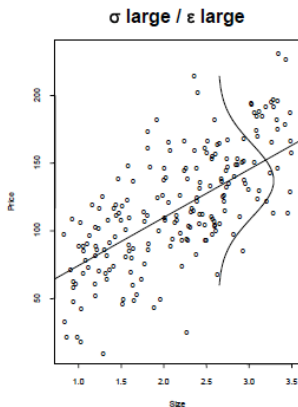
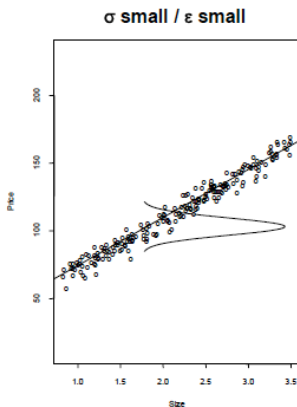


Conditional Distributions

The conditional distribution for Y given X is Normal:

$$Y|X \sim N(\beta_0 + \beta_1 X, \sigma^2).$$

σ controls dispersion:



Conditional vs Marginal Distributions

More on the conditional distribution:

$$Y|X \sim N(E[Y|X], \text{var}(Y|X)).$$

- ▶ The conditional mean is

$$E[Y|X] = E[\beta_0 + \beta_1 X + \varepsilon] = \beta_0 + \beta_1 X.$$

- ▶ The conditional variance is

$$\text{var}(Y|X) = \text{var}(\beta_0 + \beta_1 X + \varepsilon) = \text{var}(\varepsilon) = \sigma^2.$$

Remember our sliced boxplots:

- ▶ $\sigma^2 < \text{var}(Y)$ if X and Y are related.

Summary of Simple Linear Regression

Assume that all observations are drawn from our regression model and that errors on those observations are independent.

The model is

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

where ε is independent and identically distributed $N(0, \sigma^2)$.

- ▶ **independence** means that knowing ε_i doesn't affect your views about ε_j
- ▶ **identically distributed** means that we are using the same normal for every ε_i

Summary of Simple Linear Regression

The model is

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma^2).$$

The SLR has 3 basic parameters:

- ▶ β_0, β_1 (linear pattern)
- ▶ σ (variation around the line).

Key Characteristics of Linear Regression Model

- ▶ Mean of Y is **linear** in X .
- ▶ Error terms (deviations from line) are **normally distributed** (very few deviations are more than 2 sd away from the regression mean).
- ▶ Error terms have **constant variance**.

Estimation for the SLR Model

SLR assumes every observation in the dataset was generated by the model:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

This is a model for the conditional distribution of Y given X.

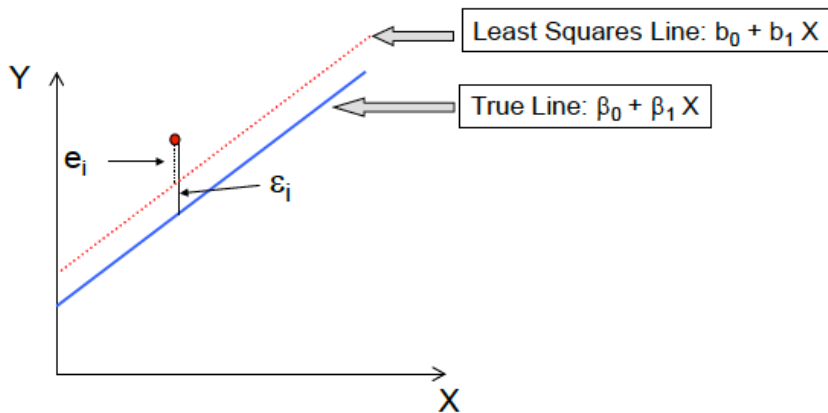
We use Least Squares *to estimate* β_0 and β_1 :

$$\hat{\beta}_1 = b_1 = r_{xy} \times \frac{s_y}{s_x}$$

$$\hat{\beta}_0 = b_0 = \bar{Y} - b_1 \bar{X}$$

Estimation for the SLR Model

NOTE!!: β_0 is not b_0 , β_1 is not b_1 and ε_i is not e



Estimation of Error Variance

We estimate s^2 with:

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{SSE}{n-2}$$

(2 is the number of regression coefficients; i.e. 2 for β_0 and β_1).

We have $n - 2$ degrees of freedom because 2 have been “used up” in the estimation of b_0 and b_1 .

We usually use $s = \sqrt{SSE/(n-2)}$, in the same units as Y . It's also called the **regression standard error**.

Degrees of Freedom

Degrees of Freedom is the number of times you get to observe useful information about the variance you're trying to estimate.

For example, consider $SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$:

- ▶ If $n = 1$, $\bar{Y} = Y_1$ and $SST = 0$: since Y_1 is “used up” estimating the mean, we haven't observed any variability!
- ▶ For $n > 1$, we've only had $n - 1$ chances for deviation from the mean, and we estimate $s_y^2 = SST / (n - 1)$.

In regression with p coefficients (e.g., $p = 2$ in SLR), you only get $n - p$ real observations of variability $\Rightarrow DoF = n - p$.

Estimation of Error Variance

Where is s in the Excel output?

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.909209967
R Square	0.826662764
Adjusted R Square	0.81332913
Standard Error	14.13839732
Observations	15

S

ANOVA

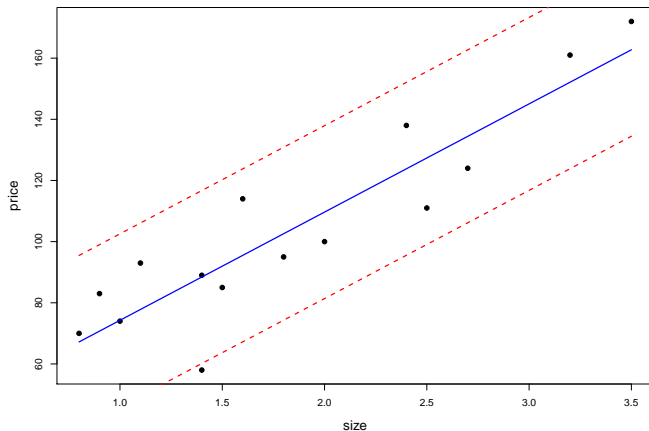
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	12393.10771	12393.10771	61.99831126	2.65987E-06
Residual	13	2598.625623	199.8942787		
Total	14	14991.73333			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	<i>Lower 95.0%</i>	<i>Upper 95.0%</i>
Intercept	38.88468274	9.09390389	4.275906499	0.000902712	19.23849785	58.53086763	19.23849785	58.53086763
X Variable 1	35.38596255	4.494082942	7.873900638	2.65987E-06	25.67708664	45.09483846	25.67708664	45.09483846

Remember that whenever you see “standard error” read it as estimated standard deviation: σ is the standard deviation.

One Picture Summary of SLR

- ▶ The plot below has the house data, the fitted regression line ($b_0 + b_1X$) and $\pm 2 * s...$
- ▶ From this picture, what can you tell me about β_0 , β_1 and σ^2 ?
How about b_0 , b_1 and s^2 ?



Sampling Distribution of Least Squares Estimates

How much do our estimates depend on the particular random sample that we happen to observe? Imagine:

- ▶ Randomly draw different samples of the same size.
- ▶ For each sample, compute the estimates b_0 , b_1 , and s .

If the estimates don't vary much from sample to sample, then it doesn't matter which sample you happen to observe.

If the estimates do vary a lot, then it matters which sample you happen to observe.

The Importance of Understanding Variation

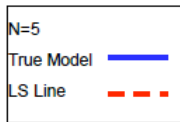
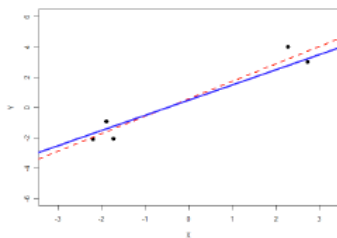
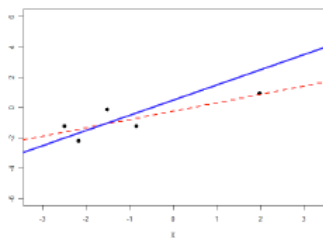
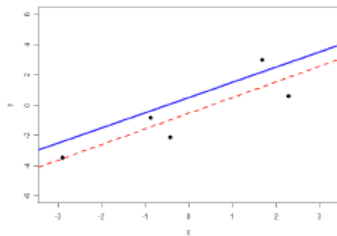
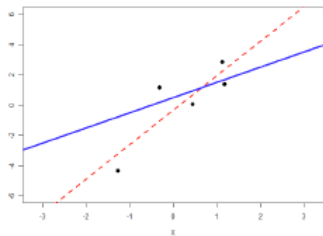
When **estimating** a quantity, it is vital to develop a notion of the **precision** of the estimation; for example:

- ▶ estimate the slope of the regression line
- ▶ estimate the value of a house given its size
- ▶ estimate the expected return on a portfolio
- ▶ estimate the value of a brand name
- ▶ estimate the damages from patent infringement

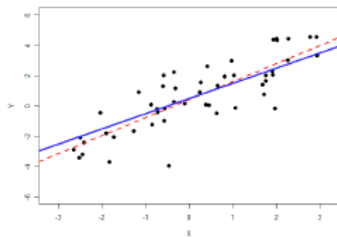
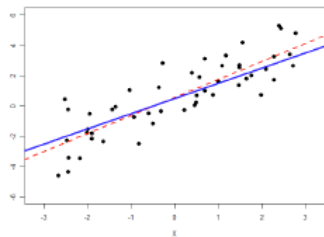
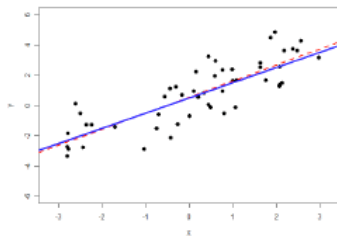
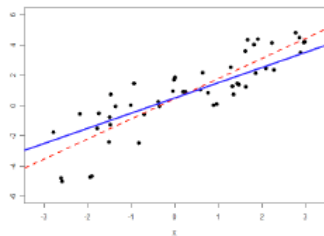
Why is this important?

We are making decisions based on estimates, and these may be very sensitive to the accuracy of the estimates!

Sampling Distribution of Least Squares Estimates



Sampling Distribution of Least Squares Estimates



N=50
True Model ———
LS Line - - - -

Sampling Distribution of Least Squares Estimates

LS lines are much closer to the true line when $n = 50$.

For $n = 5$, some lines are close, others aren't:

we need to get "lucky"

Sampling Distribution of b_1

The sampling distribution of b_1 describes how estimator $b_1 = \hat{\beta}_1$ varies over different samples with the X values fixed.

It turns out that b_1 is normally distributed (approximately):

$$b_1 \sim N(\beta_1, s_{b_1}^2).$$

- ▶ b_1 is unbiased: $E[b_1] = \beta_1$.
- ▶ s_{b_1} is the **standard error of b_1** . In general, the standard error is the standard deviation of an estimate. It determines **how close** b_1 is to β_1 .
- ▶ This is a number directly available from the regression output.

Sampling Distribution of b_1

Can we intuit what should be in the formula for s_{b_1} ?

- ▶ How should s figure in the formula?
- ▶ What about n ?
- ▶ Anything else?

$$s_{b_1}^2 = \frac{s^2}{\sum (X_i - \bar{X})^2} = \frac{s^2}{(n-1)s_x^2}$$

Three Factors:

sample size (n), error variance (s^2), and X -spread (s_x).

Sampling Distribution of b_0

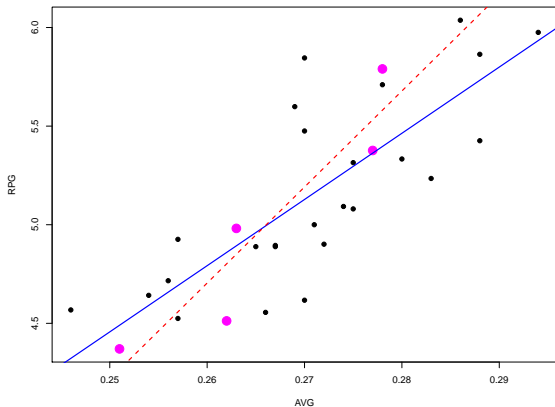
The intercept is also **normal** and **unbiased**: $b_0 \sim N(\beta_0, s_{b_0}^2)$.

$$s_{b_0}^2 = \text{var}(b_0) = s^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{(n-1)s_X^2} \right)$$

What is the intuition here?

Example: Runs per Game and AVG

- ▶ blue line: all points
- ▶ red line: only purple points
- ▶ Which slope is closer to the true one? How much closer?



Example: Runs per Game and AVG

Regression with all points

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.798496529
R Square	0.637596707
Adjusted R Square	0.624653732
Standard Error	0.298493066
Observations	30

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>significance F</i>
Regression	1	4.38915033	4.38915	49.26199	1.239E-07
Residual	28	2.494747094	0.089098		
Total	29	6.883897424			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	-3.936410446	1.294049995	-3.04193	0.005063	-6.587152	-1.2856692
AVG	33.57186945	4.783211061	7.018689	1.24E-07	23.773906	43.369833

$$s_{b_1} = 4.78$$

Example: Runs per Game and AVG

Regression with subsample

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.933601392
R Square	0.87161156
Adjusted R Square	0.828815413
Standard Error	0.244815842
Observations	5

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>significance F</i>
Regression	1	1.220667405	1.220667	20.36659	0.0203329
Residual	3	0.17980439	0.059935		
Total	4	1.400471795			

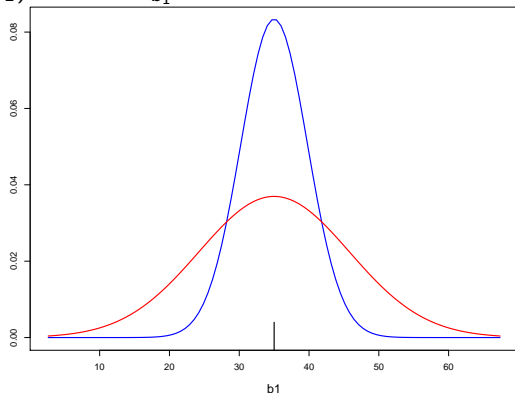
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	-7.956288201	2.874375987	-2.76801	0.069684	-17.10384	1.191259
AVG	48.69444328	10.78997028	4.512936	0.020333	14.355942	83.03294

$$s_{b_1} = 10.78$$

Example: Runs per Game and AVG

$$b_1 \sim N(\beta_1, s_{b_1}^2)$$

- ▶ Suppose $\beta_1 = 35$
- ▶ blue line: $N(35, 4.78^2)$; red line: $N(35, 10.78^2)$
- ▶ $(b_1 - \beta_1) \approx \pm 2 \times s_{b_1}$



Confidence Intervals

Since $b_1 \sim N(\beta_1, s_{b_1}^2)$, Thus:

- ▶ 68% Confidence Interval: $b_1 \pm 1 \times s_{b_1}$
- ▶ 95% Confidence Interval: $b_1 \pm 2 \times s_{b_1}$
- ▶ 99% Confidence Interval: $b_1 \pm 3 \times s_{b_1}$

Same thing for b_0

- ▶ 95% Confidence Interval: $b_0 \pm 2 \times s_{b_0}$

The confidence interval provides you with a set of plausible values for the parameters

Example: Runs per Game and AVG

Regression with all points

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.798496529
R Square	0.637596707
Adjusted R Square	0.624653732
Standard Error	0.298493066
Observations	30

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>significance F</i>
Regression	1	4.38915033	4.38915	49.26199	1.239E-07
Residual	28	2.494747094	0.089098		
Total	29	6.883897424			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	-3.936410446	1.294049995	-3.04193	0.005063	-6.587152	-1.2856692
AVG	33.57186945	4.783211061	7.018689	1.24E-07	23.773906	43.369833

$$[b_1 - 2 \times s_{b_1}; b_1 + 2 \times s_{b_1}] \approx [23.77; 43.36]$$

Testing

Suppose we want to assess whether or not β_1 equals a proposed value β_1^0 . This is called **hypothesis testing**.

Formally we test the null hypothesis:

$$H_0 : \beta_1 = \beta_1^0$$

vs. the alternative

$$H_1 : \beta_1 \neq \beta_1^0$$

Testing

That are 2 ways we can think about testing:

1. Building a test statistic... the **t-stat**,

$$t = \frac{b_1 - \beta_1^0}{s_{b_1}}$$

This quantity measures how many standard deviations the estimate (b_1) from the proposed value (β_1^0).

If the absolute value of t is greater than 2, we need to worry (why?)... we **reject** the hypothesis.

Testing

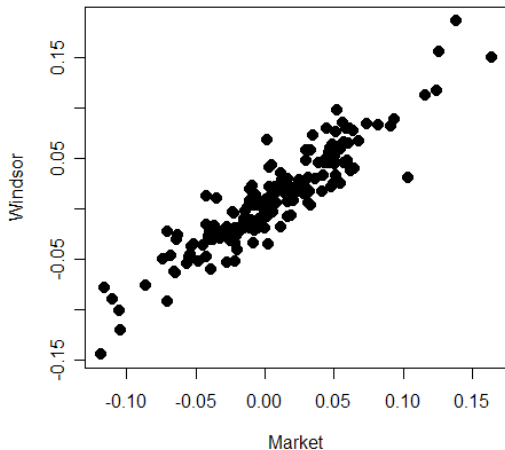
2. Looking at the **confidence interval**. If the proposed value is outside the confidence interval you **reject** the hypothesis.

Notice that this is equivalent to the t-stat. An absolute value for t greater than 2 implies that the proposed value is outside the confidence interval... therefore reject.

This is my preferred approach for the testing problem. You can't go wrong by using the confidence interval!

Example: Mutual Funds

Let's investigate the performance of the Windsor Fund, an aggressive large cap fund by Vanguard...



The plot shows monthly returns for Windsor vs. the S&P500

Example: Mutual Funds

Consider a CAPM regression for the Windsor mutual fund.

$$r_w = \beta_0 + \beta_1 r_{sp500} + \epsilon$$

Let's first test $\beta_1 = 0$

$H_0 : \beta_1 = 0$. Is the Windsor fund related to the market?

$H_1 : \beta_1 \neq 0$

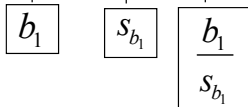
Example: Mutual Funds

Regression Statistics	
Multiple R	0.923417768
R Square	0.852700374
Adjusted R Square	0.851872848
Standard Error	0.018720015
Observations	180

ANOVA

	df	SS	MS	F	Significance F
Regression	1	0.3611	0.361099761	1030.421266	6.0291E-76
Residual	178	0.062378	0.000350439		
Total	179	0.423478			

	Coefficients	Standard Err	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	0.003646881	0.001409	2.587596412	0.010462425	0.000865657	0.006428	0.000866	0.006428
X Variable 1	0.935717012	0.02915	32.10017549	6.0291E-76	0.878193151	0.993241	0.878193	0.993241



- ▶ $t = 32.10\dots$ reject $\beta_1 = 0!!$
- ▶ the 95% confidence interval is $[0.87; 0.99]\dots$ again, reject!!

Example: Mutual Funds

Now let's test $\beta_1 = 1$. What does that mean?

$H_0 : \beta_1 = 1$ Windsor is as risky as the market.

$H_1 : \beta_1 \neq 1$ and Windsor softens or exaggerates market moves.

We are asking whether or not Windsor moves in a different way than the market (e.g., is it more conservative?).

Example: Mutual Funds

Regression Statistics	
Multiple R	0.923417768
R Square	0.852700374
Adjusted R Square	0.851872848
Standard Error	0.018720015
Observations	180

ANOVA					
	df	SS	MS	F	Significance F
Regression	1	0.3611	0.361099761	1030.421266	6.0291E-76
Residual	178	0.062378	0.000350439		
Total	179	0.423478			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	0.003646881	0.001409	2.587596412	0.010462425	0.000865657	0.006428	0.000866	0.006428
X Variable 1	0.935717012	0.02915	32.10017549	6.0291E-76	0.878193151	0.993241	0.878193	0.993241

$$b_1$$

$$s_{b_1}$$

$$\frac{b_1}{s_{b_1}}$$

- ▶ $t = \frac{b_1 - 1}{s_{b_1}} = \frac{-0.0643}{0.0291} = -2.205\dots$ reject.
- ▶ the 95% confidence interval is [0.87; 0.99]... again, reject, but...

Testing – Why I like Conf. Int.

- ▶ Suppose in testing $H_0 : \beta_1 = 1$ you got a t-stat of 6 and the confidence interval was

$$[1.00001, 1.00002]$$

Do you reject $H_0 : \beta_1 = 1$? Could you justify that to your boss? **Probably not!** (why?)

Testing – Why I like Conf. Int.

- ▶ Now, suppose in testing $H_0 : \beta_1 = 1$ you got a t-stat of -0.02 and the confidence interval was

$$[-100, 100]$$

Do you accept $H_0 : \beta_1 = 1$? Could you justify that to your boss? **Probably not!** (why?)

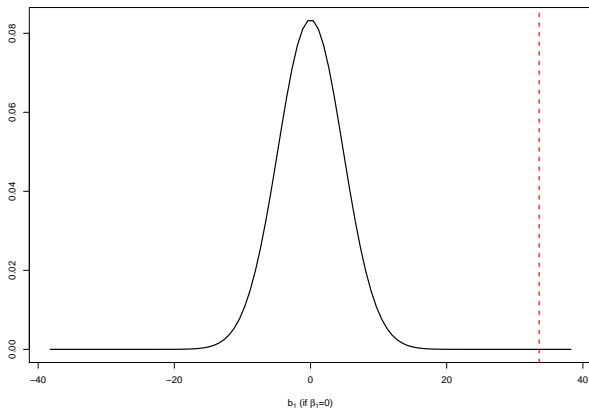
The Confidence Interval is your best friend when it comes to testing!!

P-values

- ▶ The p -value provides a measure of how **weird** your estimate is **if** the null hypothesis is true
- ▶ Small p -values are evidence against the null hypothesis
- ▶ In the AVG vs. R/G example... $H_0 : \beta_1 = 0$. How weird is our estimate of $b_1 = 33.57$?
- ▶ Remember: $b_1 \sim N(\beta_1, s_{b_1}^2)$... If the null was true ($\beta_1 = 0$),
 $b_1 \sim N(0, s_{b_1}^2)$

P-values

- ▶ Where is 33.57 in the picture below?



The p -value is the probability of seeing b_1 equal or greater than 33.57 in absolute terms. Here, $p\text{-value}=0.000000124!!$

Small p -value = bad null

P-values

- $H_0 : \beta_1 = 0 \dots$ p-value = 1.24E-07... reject!

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.798496529
R Square	0.637596707
Adjusted R Square	0.624653732
Standard Error	0.298493066
Observations	30

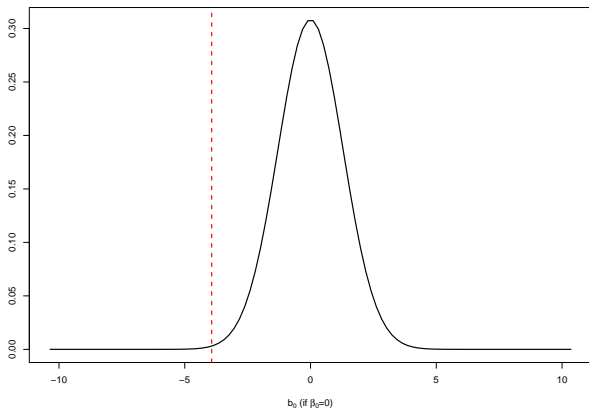
ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>significance F</i>
Regression	1	4.38915033	4.38915	49.26199	1.239E-07
Residual	28	2.494747094	0.089098		
Total	29	6.883897424			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	-3.936410446	1.294049995	-3.04193	0.005063	-6.587152	-1.2856692
AVG	33.57186945	4.783211061	7.018689	1.24E-07	23.773906	43.369833

P-values

- ▶ How about $H_0 : \beta_0 = 0$? How weird is $b_0 = -3.936$?



The p -value (the probability of seeing b_1 equal or greater than -3.936 in absolute terms) is **0.005**.

Small p -value = bad null

P-values

- ▶ $H_0 : \beta_0 = 0 \dots$ **p-value = 0.005**... we still reject, but not with the same strength.

SUMMARY OUTPUT

<i>Regression Statistics</i>	
Multiple R	0.798496529
R Square	0.637596707
Adjusted R Square	0.624653732
Standard Error	0.298493066
Observations	30

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>significance F</i>
Regression	1	4.38915033	4.38915	49.26199	1.239E-07
Residual	28	2.494747094	0.089098		
Total	29	6.883897424			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	-3.936410446	1.294049995	-3.04193	0.005063	-6.587152	-1.2856692
AVG	33.57186945	4.783211061	7.018689	1.24E-07	23.773906	43.369833

Testing – Summary

- ▶ Large t or small p -value mean the same thing...
- ▶ p -value < 0.05 is equivalent to a t -stat > 2 in absolute value
- ▶ Small p -value means something weird happen if the null hypothesis was true...
- ▶ Bottom line, small p -value \rightarrow REJECT! Large $t \rightarrow$ REJECT!
- ▶ But remember, always look at the confidence interval!

Forecasting

The **conditional forecasting problem**: Given covariate X_f and sample data $\{X_i, Y_i\}_{i=1}^n$, predict the “future” observation y_f .

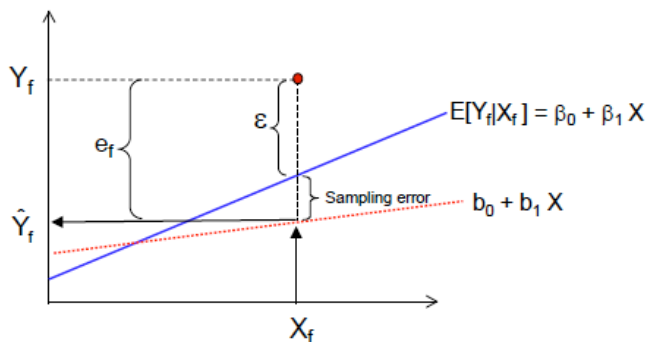
The solution is to use our LS fitted value: $\hat{Y}_f = b_0 + b_1 X_f$.

This is the easy bit. The hard (**and very important!**) part of forecasting is assessing uncertainty about our predictions.

Forecasting

If we use \hat{Y}_f , our **prediction error** is

$$\begin{aligned}e_f &= Y_f - \hat{Y}_f = Y_f - b_0 - b_1 X_f \\&= (\beta_0 + \beta_1 X_f + \epsilon) - (b_0 + b_1 X_f) \\&= (\beta_0 - b_0) + (\beta_1 - b_1) X_f + \epsilon\end{aligned}$$



Forecasting

The most commonly used approach is to assume that $\beta_0 \approx b_0$, $\beta_1 \approx b_1$ and $\sigma \approx s$... in this case, the error is just ϵ hence the 95% plug-in prediction interval is:

$$b_0 + b_1 X_f \pm 2 \times s$$

It's called "plug-in" because we just plug-in the estimates (b_0 , b_1 and s) for the unknown parameters (β_0 , β_1 and σ).

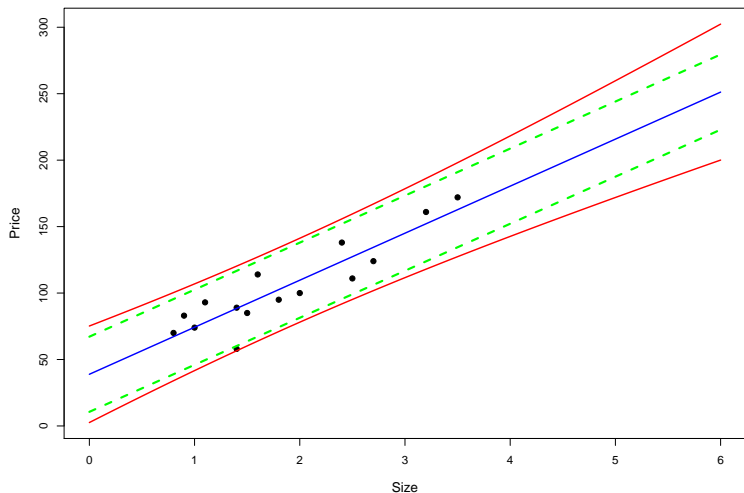
Forecasting

Just remember that you are uncertain about b_0 and b_1 ! As a practical matter if the confidence intervals are big you should be careful!! Some statistical software will give you a larger (and correct) predictive interval.

A large predictive error variance (high uncertainty) comes from

- ▶ Large s (i.e., large ε 's).
- ▶ Small n (not enough data).
- ▶ Small s_x (not enough observed spread in covariates).
- ▶ Large difference between X_f and \bar{X} .

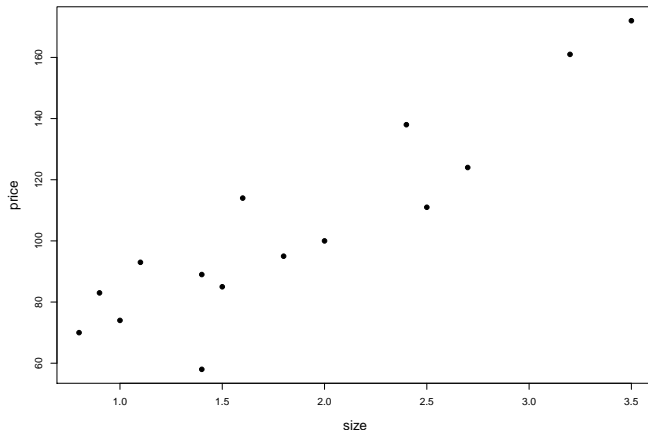
Forecasting



- ▶ **Red lines:** prediction intervals
- ▶ **Green lines:** “plug-in” prediction intervals

House Data – one more time!

- ▶ $R^2 = 82\%$
- ▶ Great R^2 , we are happy using this model to predict house prices, right?



House Data – one more time!

- ▶ But, $s = 14$ leading to a predictive interval width of about US\$60,000!! How do you feel about the model now?
- ▶ As a practical matter, s is a much more relevant quantity than R^2 . Once again, *intervals* are your friend!

