

Prediction

Business Analytics

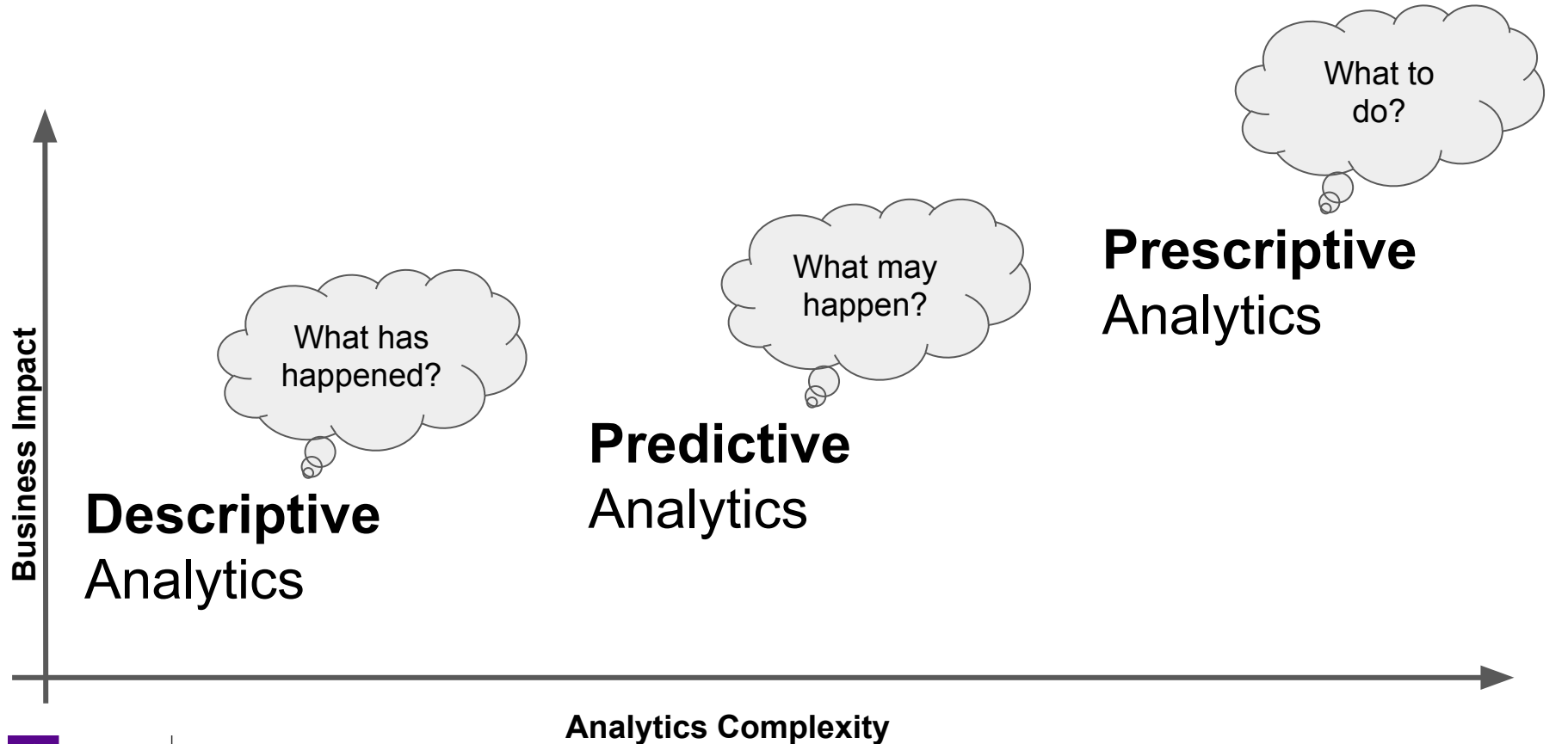
Business Analytics Definition

Business analytics refers to the application of data analysis and modeling techniques for understanding business situations and improving business decisions.

IMPLICATIONS:

- data → past business performance
- methods → statistics + mathematics + computational methods
- business decisions → actionable insight

Types of Analytics



Understanding Your Data

Exploratory analysis of data is useful for:

- understanding data properties
- detecting errors, ensuring data quality
- finding patterns in data
- determining relationships among variables
- checking assumptions
- mapping business problems into data mining tasks and suggesting modeling strategies



Lesson Objectives

1. Regression - Theory

- a. Linear Models
- b. Ordinary Least Squares
- c. Simple Linear Regression

2. Regression Applied

- a. Model strength
- b. Model interpretation
- c. Dummy variables
- d. Non-linear transformations

3. Classification

- a. Statistical classification
- b. Decision Trees



Linear Models

Regression models estimate the relationships among variables to predict outcomes.

Example: How does bike trip duration change as we introduce a new customer type, a new pricing scheme, or with different weather conditions.

In this week you will learn the basics of regression analysis and the specifics about linear regression models that example the relationship between numerical variables.

Business Case:

What is the influence of a variable (price, advertising, and etc.) on an outcome (market shares, sales, overall satisfaction)?

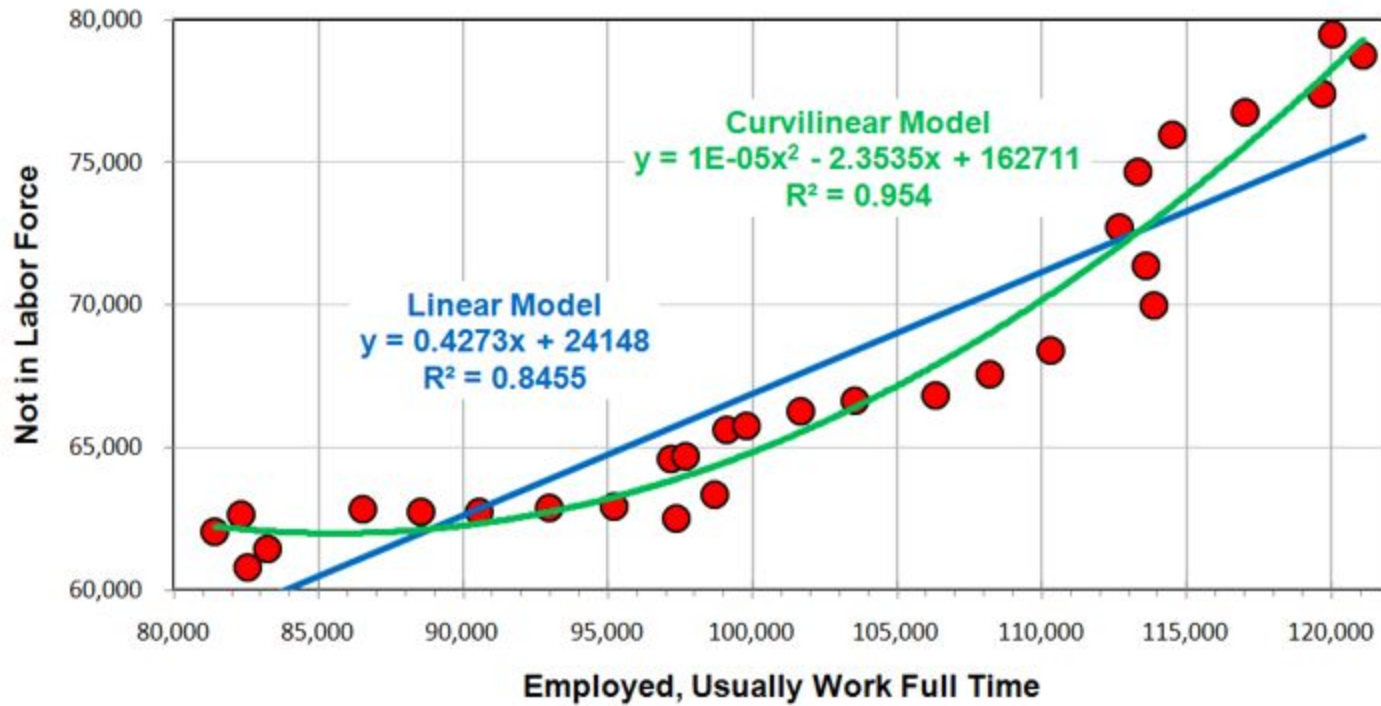
$$X \rightarrow Y$$

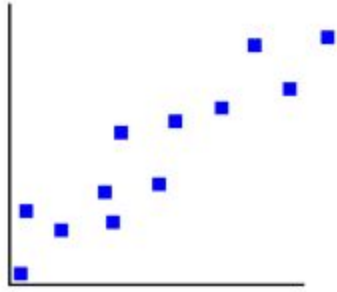
Independent variable (X) \rightarrow Dependant variable (Y)

$$y=mx+b$$

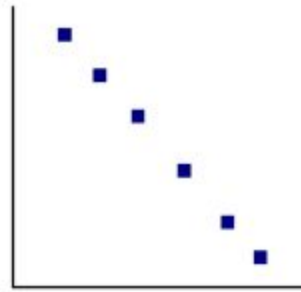
NOTE IN TERMINOLOGY

- *Y is know as the dependent variable the variable that regression model seek to predict or response variable*
- *X is the independent variable, predictor or explanatory variable.*

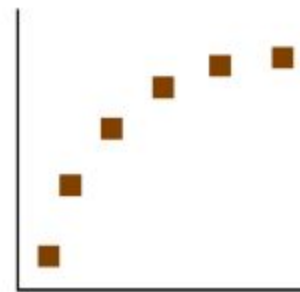




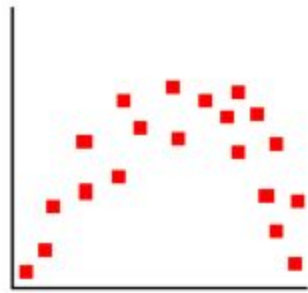
A



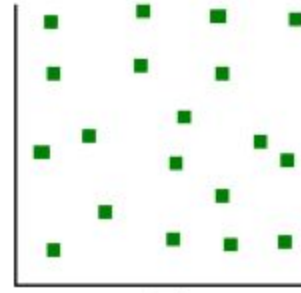
B



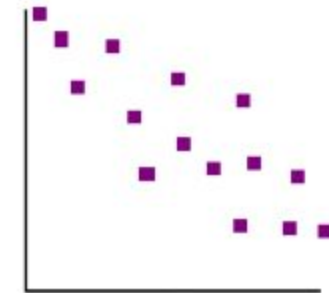
C



D



E



F

Simple Linear Regression

- The “workhorse” of statistical analysis is the simple linear regression.
- Used to determine the relationship between two variables.
 - Given one variable, a regression will provide the expected value of the other variable.
- The outcome of the regression → Y: response.
- The input variable → X: predictor.

$$Y_i = b_0 + b_1 X_i + \epsilon_i, i=1, \dots, n$$

where:

Y_i = i th observation of the dependent variable, Y

X_i = i th observation of the independent variable, X

b_0 = regression intercept term

b_1 = regression slope coefficient

ϵ_i = residual for the i th observation (also referred to as the disturbance term or error term)

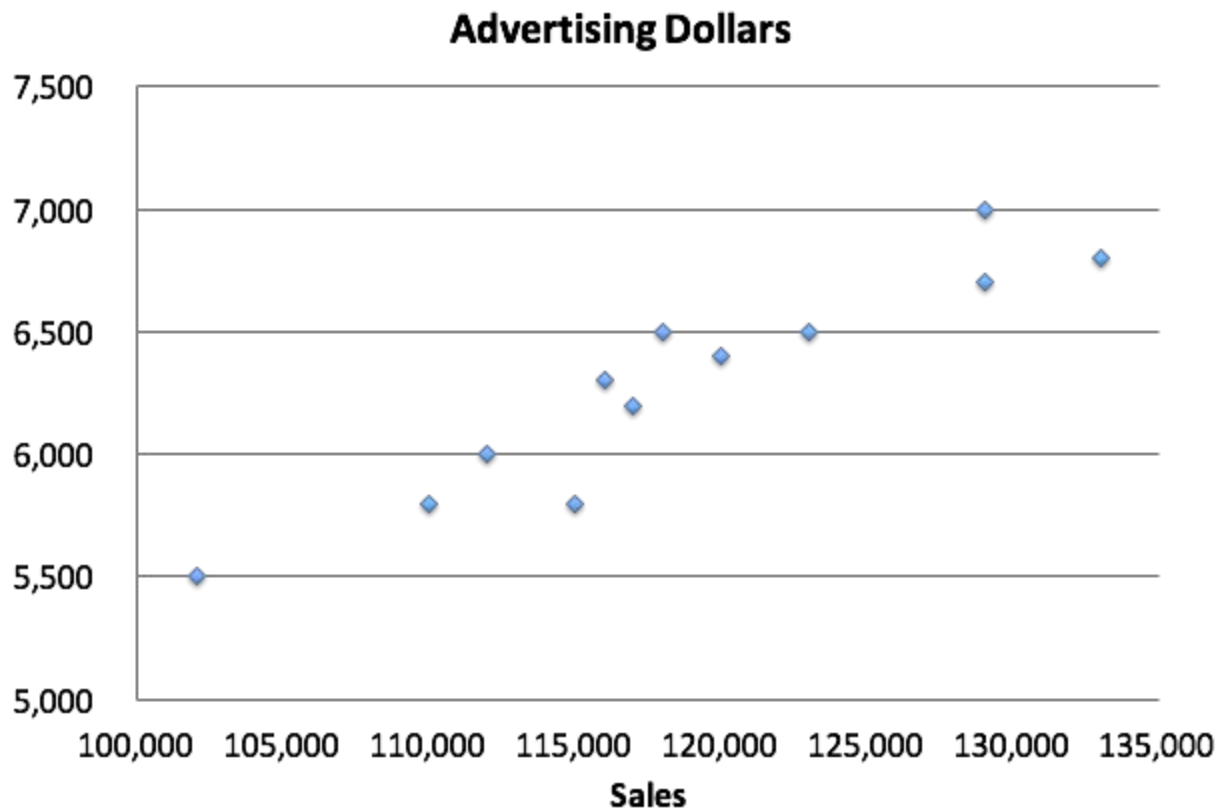


Sales vs Advertising

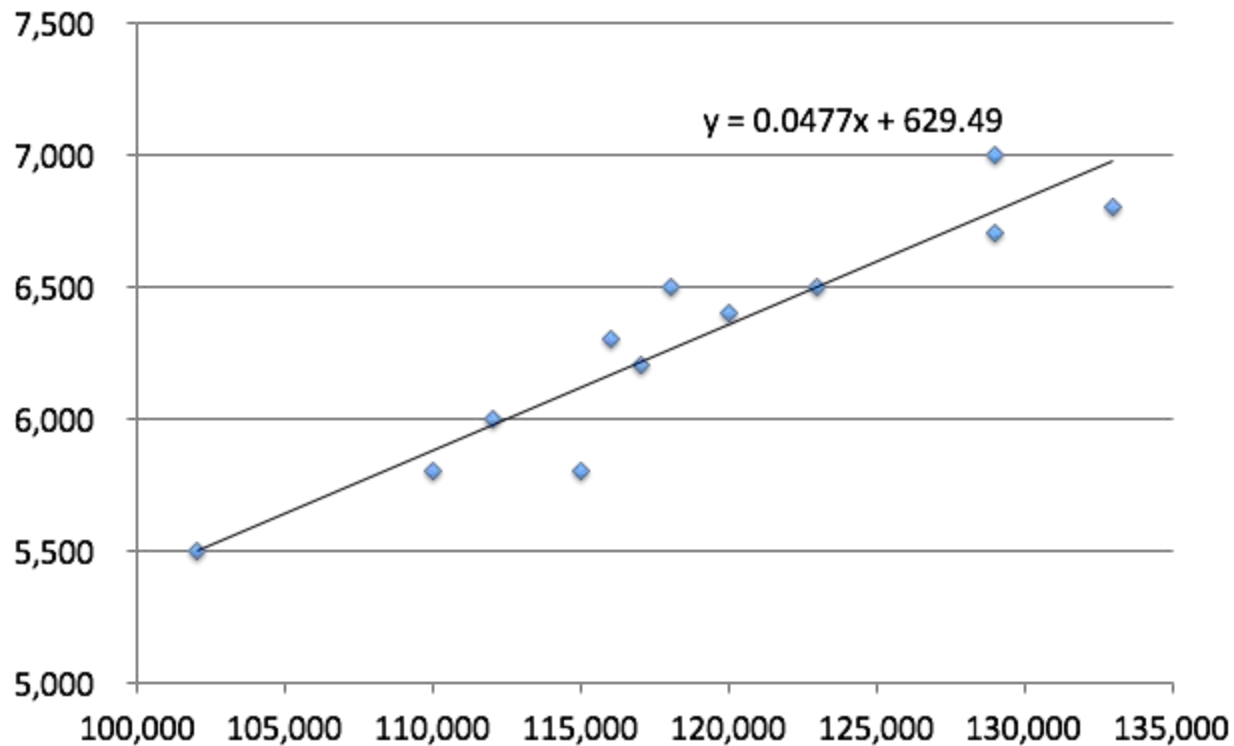
| Month | Sales | Advertising Dollars |
|-------|---------|---------------------|
| Jan | 102,000 | 5,500 |
| Feb | 110,000 | 5,800 |
| Mar | 112,000 | 6,000 |
| Apr | 115,000 | 5,800 |
| May | 117,000 | 6,200 |
| Jun | 116,000 | 6,300 |
| Jul | 118,000 | 6,500 |
| Aug | 129,000 | 7,000 |
| Sep | 123,000 | 6,500 |
| Oct | 120,000 | 6,400 |
| Nov | 129,000 | 6,700 |
| Dec | 133,000 | 6,800 |



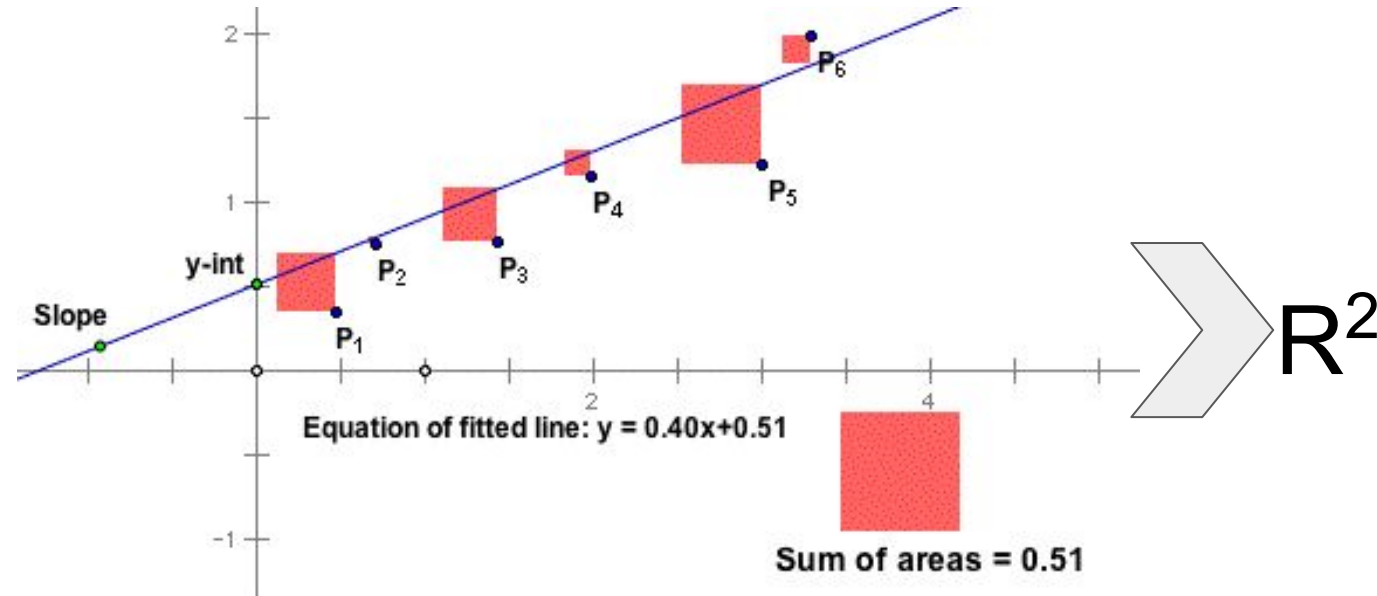
Scatter



Advertising Dollars



Ordinary Least Squares



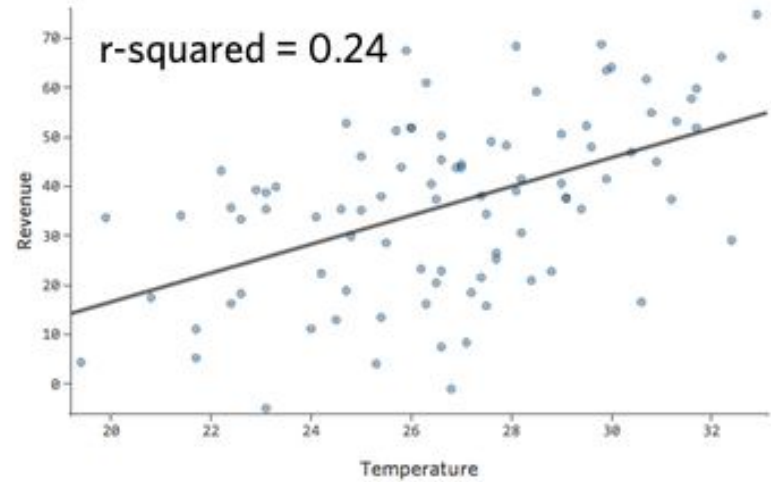
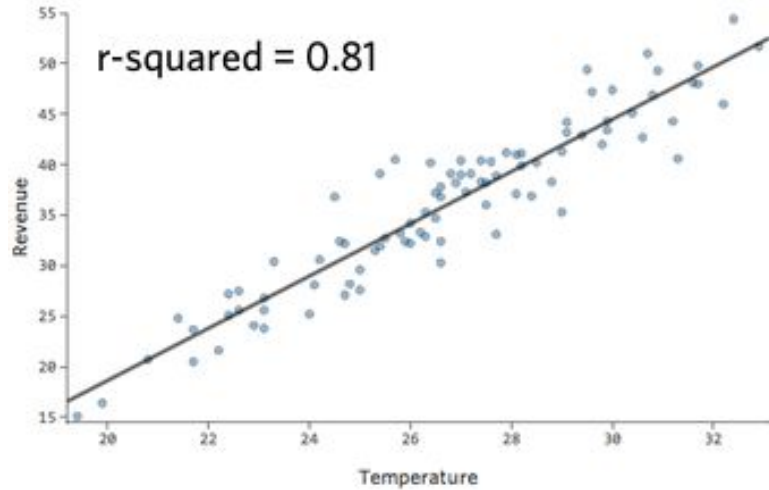
R^2 is an statistical measure of how close the data are to the fitted regression line.

It indicates the goodness of fit of the model.

R^2 definition: Explained variation / Total variation

R^2 is always between 0 and 100%:

- 0% → model explains none of the variability
- 100% → model explains all the variability



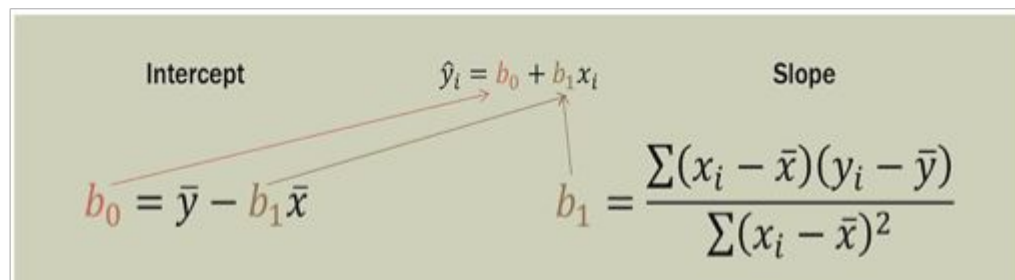
Calculating b1:

0.0477

| x-mean(x) | y-mean(y) | (x-mean(x))* (y-mean(y)) | (x-mean(x))^2 |
|------------|-----------|--------------------------|----------------|
| -16,666.67 | -791.67 | 13,194,444.44 | 277,777,777.78 |
| -8,666.67 | -491.67 | 4,261,111.11 | 75,111,111.11 |
| -6,666.67 | -291.67 | 1,944,444.44 | 44,444,444.44 |
| -3,666.67 | -491.67 | 1,802,777.78 | 13,444,444.44 |
| -1,666.67 | -91.67 | 152,777.78 | 2,777,777.78 |
| -2,666.67 | 8.33 | -22,222.22 | 7,111,111.11 |
| -666.67 | 208.33 | -138,888.89 | 444,444.44 |
| 10,333.33 | 708.33 | 7,319,444.44 | 106,777,777.78 |
| 4,333.33 | 208.33 | 902,777.78 | 18,777,777.78 |
| 1,333.33 | 108.33 | 144,444.44 | 1,777,777.78 |
| 10,333.33 | 408.33 | 4,219,444.44 | 106,777,777.78 |
| 14,333.33 | 508.33 | 7,286,111.11 | 205,444,444.44 |

Calculating b0:

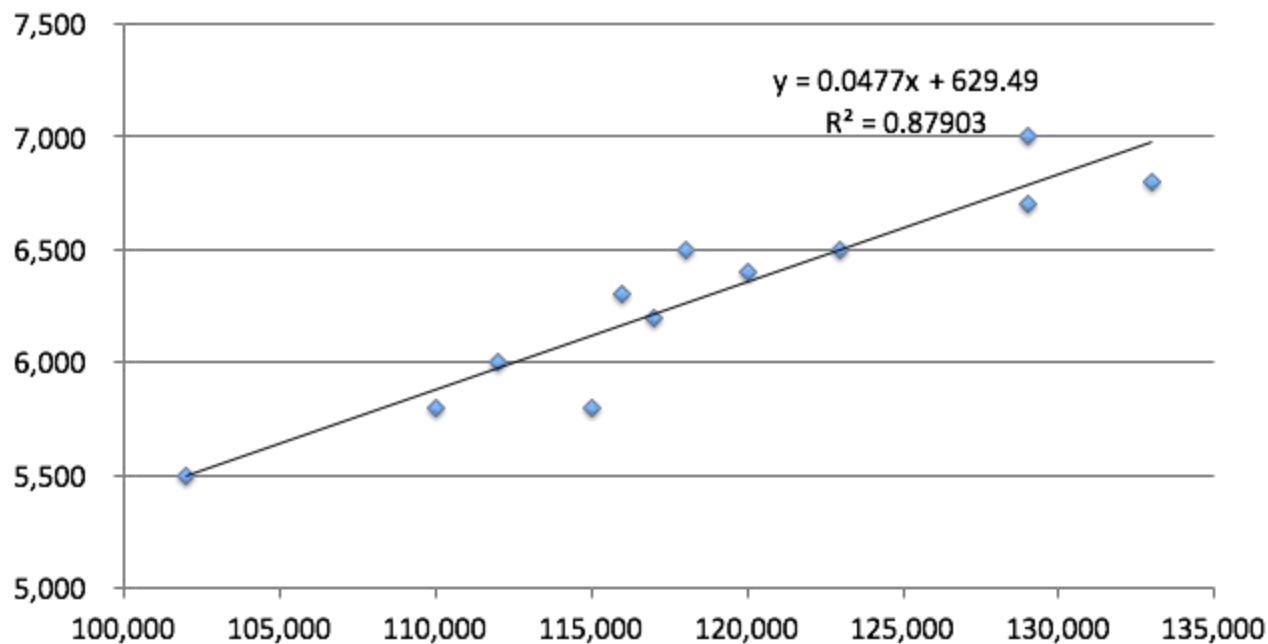
629.4926



$$R^2 = \left[\frac{\sum (xy) - (\sum x)(\sum y)}{n} \right]^2 \div \left[\left(\sum x^2 - \frac{(\sum x)^2}{n} \right) \left(\sum y^2 - \frac{(\sum y)^2}{n} \right) \right]$$



Advertising Dollars

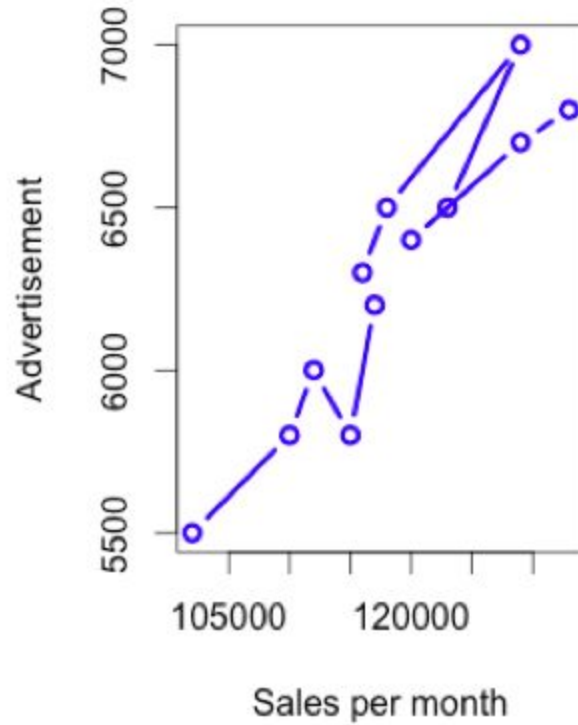
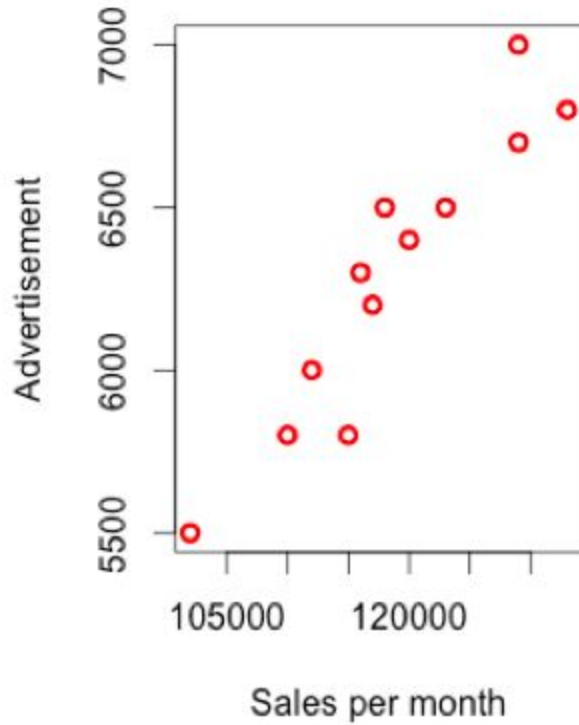


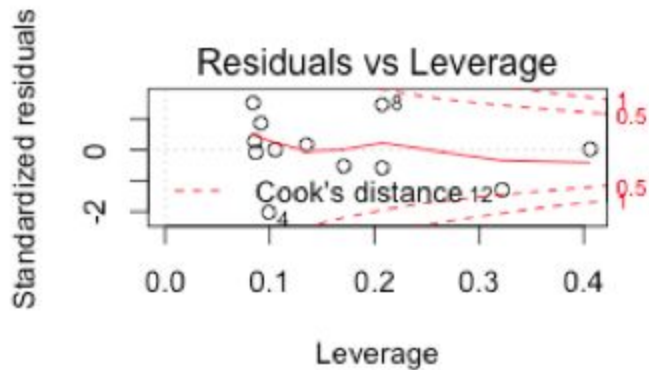
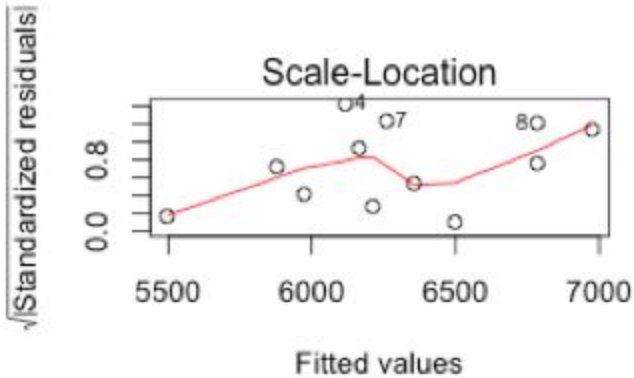
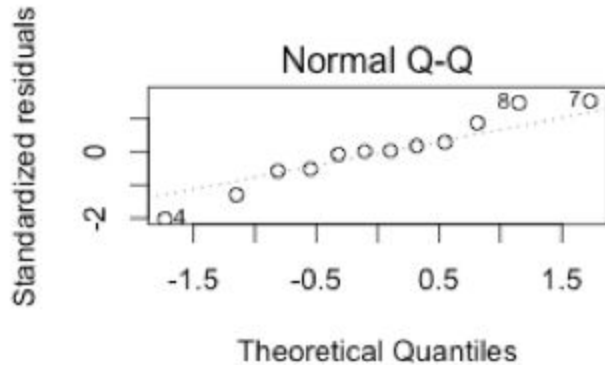
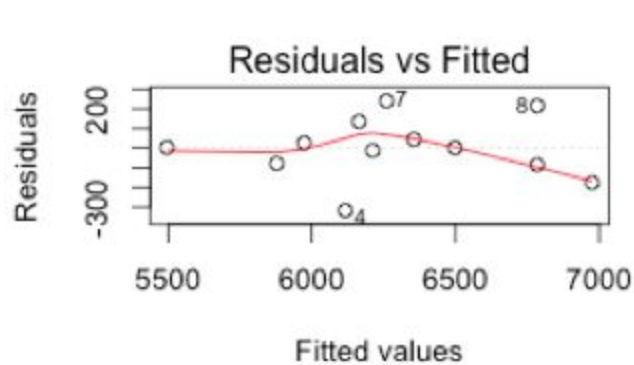
R Code for linear models

```
data <- read.csv("~/Google Drive/Business Analytics/Data/Sales vs Advertisement.csv",
  header=TRUE, stringsAsFactors=TRUE)
dim(data)
names(data)
x=data$Sales
y=data$Advertising.Dollars
par(mfrow=c(1,2))
plot(x,y , col="red", lwd=3,
  ylab="Advertisement", xlab="Sales per month")
plot(x,y, type="b", col="blue", lwd=3,
  ylab="Advertisement", xlab="Sales per month")

model<-lm(y ~ x)
model
summary(model)

par(mfrow=c(2,2))
plot(model)
```





R Session

Build a linear model for Zagat using “Food” as and predictors and “Price” as a response.

Build a linear model for Zagat using “Food” and “Decor” as and predictors and “Price” as a response. Hint use $lm(y \sim x1+x2)$

Build a linear model for Zagat using “Food”, “Decor”, and “Service” as and predictors and “Price” as a response. Hint use $lm(y \sim x1+x2+x3)$

