

Unit 5. Regression and Correlation

“ ‘Don’t let us quarrel,’ the White Queen said in an anxious tone. ‘What is the cause of lightning?’ ‘The cause of lightning,’ Alice said very decidedly, for she felt quite certain about this, ‘is the thunder-oh no!’, she hastily corrected herself. ‘I meant the other way.’ ‘It’s too late to correct it,’ said the Red Queen: ‘when you’ve once said a thing, that fixes it, and you must take the consequences.’ ”
- Carroll

Menopause heralds a complex interplay of hormonal and physiologic changes. Some are temporary discomforts (e.g., hot flashes, sleep disturbances, depression). Others are long-term changes that increase the risk of significant chronic health conditions, bone loss and osteoporosis in particular. Recent observations of an association between **depressive symptoms** and **low bone mineral density (BMD)** raise the intriguing possibility that alleviation of depression might confer a risk benefit with respect to bone mineral density loss and osteoporosis.

However, the finding of an association in a simple (one predictor) linear regression model analysis has multiple possible explanations, only one of which is causal. Others include, but are not limited to: (1) the apparent association is an artifact of the confounding effects of exercise, body fat, education, smoking, etc; (2) there is no relationship and we have observed a chance event of low probability (it can happen!); (3) the pathway is the other way around (low BMD causes depressive symptoms), albeit highly unlikely; and/or (4) the finding is spurious due to study design flaws (selection bias, misclassification, etc).

In settings where multiple, related predictors are associated with the outcome of interest, multiple predictor linear regression analysis allows us to investigate the joint relationships among the multiple predictors (depressive symptoms, exercise, body fat, etc) and a single continuous outcome (BMD).

In this example, we might be especially interested in using multiple predictor linear regression to isolate the effect of depressive symptoms on BMD, holding all other predictors constant (**adjustment**). Or, we might want to investigate the possibility of synergism or **interaction**.

Table of Contents

Topic		
	Learning Objectives	3
	1. <u>Review</u>	4
	a. Settings Where Regression Might be Considered	4
	b. Review - What is Statistical Modeling	7
	c. A General Approach for Model Development	8
	d. Review - Normal Theory Regression	9
	2. <u>R Illustration</u> - Fit a Simple Linear Regression Model	12
	3. <u>Multivariable Regression</u>	14
	a. Introduction	14
	b. Indicator and Design Variables	17
	c. Interaction Variables	20
	d. Look! Schematic of Confounding and Effect Modification	21
	e. The Analysis of Variance Table	22
	f. The Partial F Test	25
	g. Multiple Partial Correlation	27
	4. <u>Multivariable Model Development</u>	29
	a. Introduction	29
	b. Example – Framingham Study	30
	c. Suggested Criteria for Confounding and Interaction	37
	d. Additional Tips for Multivariable Analyses of Large Data Sets	38
	5. <u>Goodness-of-Fit and Regression Diagnostics</u>	40
	a. Introduction and Terminology.....	40
	b. Assessment of Normality.....	47
	c. Cook-Weisberg Test of Heteroscedasticity	51
	d. Method of Fractional Polynomials	52
	e. Ramsay Test for Omitted Variables	54
	f. Residuals, Leverage, & Cook's Distance	55
	g. Example – Framingham Study	57

<p>Datasets used (download from course website)</p> <p>janka.Rdata p53paper.Rdata framingham_1000.Rdata</p>	<p>Packages used (one time installation)</p> <p>ggplot2 Hmisc stargazer car gridExtra lmtest GGally summarytools</p> <p>Tip! Don't forget that R is case sensitive ...</p>
--	--

1. Learning Objectives

When you have finished this unit, you should be able to:

- Explain the concepts of association, causation, confounding, mediation, and effect modification;
- Construct and interpret a scatter plot with respect to: evidence of association, assessment of linearity, and the presence of outlying values;
- State the multiple predictor linear regression model and the assumptions necessary for its use;
- Perform and interpret the Shapiro-Wilk and Kolmogorov-Smirnov tests of normality;
- Explain the relevance of the normal probability distribution;
- Explain and interpret the coefficients (and standard error) and analysis of variance tables outputs of a single or multiple predictor regression model estimation;
- Explain and compare crude versus adjusted estimates (betas) of association;
- Explain and interpret regression model estimates of effect modification (interaction);
- Explain and interpret overall and adjusted R-squared measures of association;
- Explain and interpret overall and partial F-tests;
- Draft an analysis plan for a multiple predictor regression model analysis; and
- Explain and interpret selected regression model diagnostics: residuals, leverage, and Cook's distance.

1. Review

Simple linear regression and correlation were introduced in [BIOSTATS 540, Unit 12](#).

a. Settings Where Regression Might Be Considered

Example #1

Is the density of wood a predictor of hardness of timber?

Source:

Williams, E.J. (1959) Regression Analysis, New York: John Wiley & Sons

Wood density and timber hardness are two different things, with timber hardness being important in many of the products of wood processing. Wood density is pounds of weight per cubic foot of volume, while timber hardness is measure of force. One measure of the latter is the Janka Scale; it defines hardness as the number of pounds required to push a ball bearing into a timber sample using a machine press. So, as you might imagine, it might be of interest to estimate the parameters that define the relationship between the two so as to obtain a **prediction equation**. Thus, in this example, the predictor (explanatory variable) is wood density and the outcome (response variable) is the Janka Scale hardness score:

Y = hardness

X = density

Example #2

Does the expression of p53 change with parity and age?

Source:

Matthews et al. Parity Induced Protection Against Breast Cancer 2007.

P53 is a human gene that is a tumor suppressor gene. Malfunctions of this gene have been implicated in the development and progression of many cancers, including breast cancer. Matthews et al were interested in **exploring the relationship** of Y=p53 expression to parity and age at first pregnancy, *after adjustment for* selected risk factors for breast cancer, including: age at first mensis, family history of breast cancer, menopausal status, and history of oral contraceptive use.

- Among the initial analyses, a **simple linear regression** might be performed to obtain a thorough understanding of the relationship of p53 expression and age. Both the outcome (Y) and the predictor (X) are continuous.

Y = p53 expression

X = Age

- A **multiple linear regression** might then be performed to see if age and parity retain their predictive significance, after controlling for the other, known, risk factors for breast cancer. Thus, the analysis would consider one outcome variable (Y) and 6 predictor variables (X₁, X₂, X₃, X₄, X₅, X₆):

Y = p53
 X₁ = Age
 X₂ = Parity
 X₃ = Age at first mensis
 X₄ = Family history of breast cancer
 X₅ = Menopausal status
 X₆ = History of oral contraceptive use

Example #3

Does Air Pollution Reduce Lung Function?

Source:

Detels et al (1979) *The UCLA population studies of chronic obstructive respiratory disease. I. Methodology and comparison of lung function in areas of high and low pollution. Am. J. Epidemiol. 109: 33-58.*

Detels et al (1979) investigated the relationship of lung function to exposure to air pollution among residents of Los Angeles in the 1970's. Baseline and follow-up measurements of exposure and lung function were obtained. Also obtained were measurements of other variables that might confound or modify the effects of pollution on lung function: age, sex, height, weight, etc. Afifi, Clark and May (2004) consider portions of this data in their 2004 text, Computer-Aided Multivariate Analysis, Fourth Edition (Chapman & Hall)

- A **simple linear regression** might be performed to characterize the relationship between FEV and height:

Y = FEV, liters
 X = Height, inches

- A **multiple linear regression** might then be performed to determine the nature and strength of exposure to pollution for the prediction of lung function, taking into account the roles of other influences on lung function, such as age, height, smoking, etc. For example, the relationship of lung function to exposure to air pollution might be different for smokers and non-smokers; this would be an example of effect modification (interaction). It might also be the case that the relationship of lung function to exposure to air pollution is confounded by height. Here, we would have something like:

Y = FEV, liters
 X₁ = Exposure to air pollution
 X₂ = Height, inches
 X₃ = Smoking (1=yes, 0=no)

Example #4**Exercise and Glucose for the Prevention of Diabetes**Source:

Hulley et al (1998) *Randomized trial of estrogen plus progestin for secondary prevention of heart disease in postmenopausal women. The Heart and Estrogen/progestin Study. JAMA 280(7): 605-13.*

In the HERS study, Hulley et al. (1998) sought to determine if exercise, a modifiable behavior, might lower the risk of diabetes in non-diabetic women who are at risk of developing the disease. The question is a complex one because there are many risk factors for diabetes. Moreover, the type of woman who chooses to exercise may be related in other ways to risk of diabetes, apart from the fact of her exercise habit. For example, women who exercise regularly are typically younger and have lower body mass index (BMI); these characteristics also confer a risk benefit with respect to diabetes. Finally, the benefit of exercise may be mediated through a reduction of body mass index. Vittinghoff, Glidden, Shiboski and McCullogh (2005) consider portions of this data in their 2005 text, *Regression Methods in Biostatistics: Linear, Logistic, Survival and Repeated Measures Models* (Springer).

- A **multiple linear regression** was performed to assess the benefit of exercising at least three times/week, compared to no exercise, on blood glucose, after controlling for other factors associated with blood glucose levels. Thus, here we would have something like:

$Y = \text{Glucose, mg/dL}$

$X_1 = \text{Exercise (1=yes if 3x/week or more, 0 = no)}$

$X_2 = \text{Age, years}$

$X_3 = \text{Body Mass Index (BMI)}$

$X_4 = \text{Alcohol Use (1=yes, 0=no)}$

b. Review - What is Statistical Modeling

George E.P. Box, a very famous statistician, once said, “***All models are wrong, but some are useful.***”

Incorrectness of models notwithstanding, we do statistical modeling for very good reasons. Among them is an understanding of the natures and strengths of the relationships (if any) that might exist in a set of observations that vary.

For any set of observations, theoretically, lots of models are possible. So, how to choose? The **goal** of statistical modeling is to obtain a model that is simultaneously **minimally adequate** and a **good fit**. **The model should also make sense.**

Minimally adequate

- Each predictor is “important” in its own right
- Each extra predictor is retained in the model only if it yields a significant improvement (in fit and in variation explained).
- The model should not contain any redundant parameters (*more on this later*).

Good Fit

- Variance explained. The variability in the outcomes (the Y variable) explained is a lot
- Prediction. The outcomes predicted by the model are close to the observed outcomes.

The model should also make sense

- Biological sense. A preferred model is one based on “subject matter” considerations
- Useful. The preferred predictors are simple, measurable and convenient.

Sigh.

It is not possible to choose a model that is simultaneously minimally adequate and a perfect fit.
Model estimation and selection must achieve an appropriate balance.

c. A General Approach for Model Development

There are ***no*** rules ***nor a single best strategy***. Different study designs and research questions call for different strategies for building a regression model. **Tip**. Before you begin your model development, make a list of your study design, research aims, outcome variable, primary predictor(s), and covariates. As a general suggestion, the following approach has the advantages of providing a reasonably thorough exploration of the data and a relatively small risk of missing something important.

Preliminary – Be sure you have: (1) checked, cleaned and described your data, (2) screened the data for multivariable associations, and (3) thoroughly explored the bivariate relationships.

Step 1 – Fit the “maximal” model.

The maximal model is the large model that contains all the explanatory variables of interest as predictors. This model also contains all the covariates that might be of interest. It also contains all the interactions that might be of interest. Note the amount of the variability in the outcome that is explained.

Step 2 – Begin simplifying the model.

Inspect each of the terms in the “maximal” model with the goal of removing the predictor that is the least significant. Drop from the model the predictors that are the least significant, beginning with the higher order interactions (**Tip** -interactions are complicated and we are aiming for a simple model). Fit the reduced model. Compare the amount of variation explained by the reduced model with the amount of variation explained by the “maximal” model.

If the deletion of a predictor has little effect on the variation explained ...
Then leave that predictor out of the model.

And inspect each of the terms in the model again.

If the deletion of a predictor has a significant effect on the variation explained ...
Then put that predictor back into the model.

Step 3 – Keep simplifying the model.

Repeat step 2, over and over, until the model remaining contains nothing but significant predictor variables.

Beware of some important caveats

- Prioritize considerations of biology and what makes sense. In particular,
- Sometimes, you will want to keep a predictor in the model regardless of its statistical significance (an example is randomization assignment in a clinical trial)
- The order in which you delete terms from the model matters!

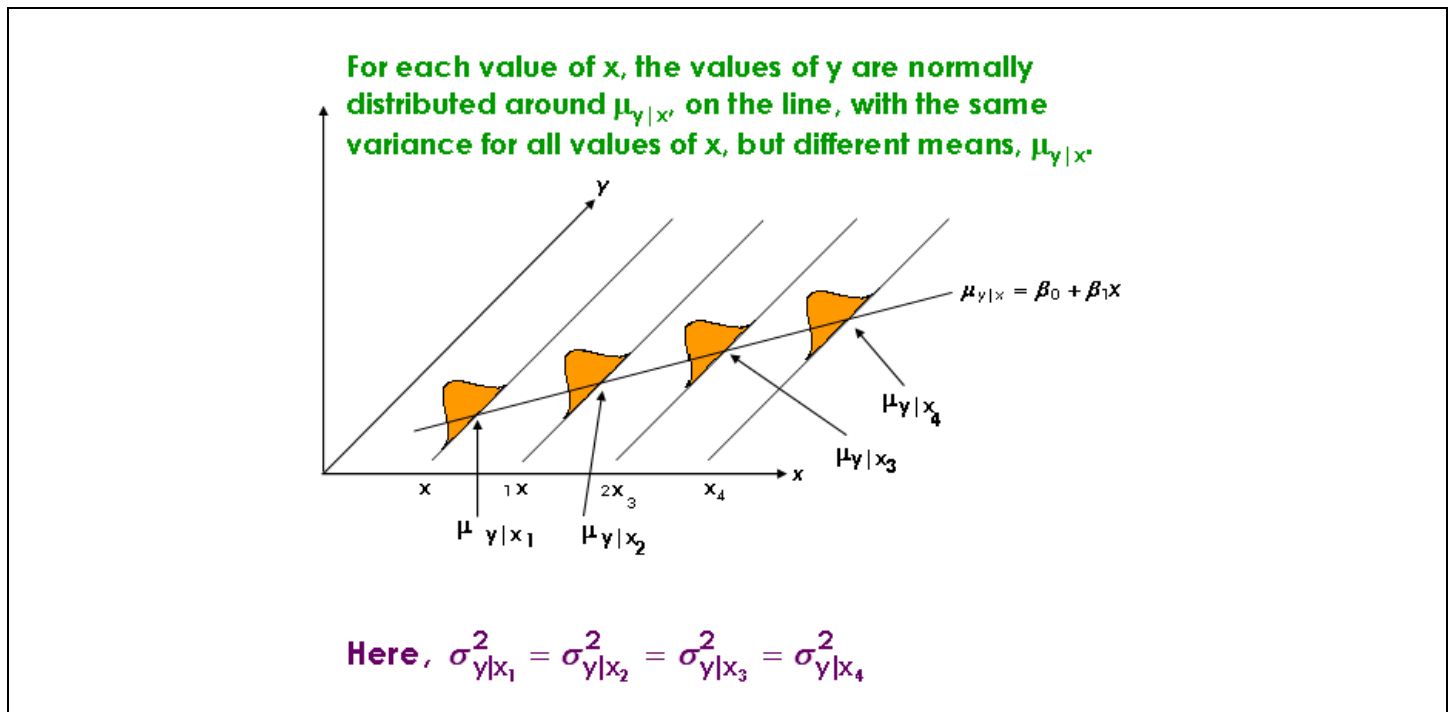
d. Review - Normal Theory Regression

Normal theory regression analysis can be used to model/investigate possibly complex relationships when:

- The outcome is a **single continuous variable (Y)** that is assumed to be **distributed normal**; *and*
- The outcome is potentially related to possibly **several predictors (X_1, X_2, \dots, X_p)** which can be **continuous or discrete**; *and*
- Some of the predictor variables might **confound** the prediction role of other explanatory variables; *and*
- Some of the predictor-outcome relationships may be different (are **modified** by) depending on the level of one or more different predictor variables (**interaction**)

Simple Linear Regression:

We're modeling the means of several subpopulations, each defined by a particular $X=x$. A simple linear regression model is one for which the mean μ (the average value) of **one continuous, and normally distributed, outcome** random variable Y (e.g. **Y=FEV** for forced expiratory volume) varies linearly with changes in **one continuous predictor** variable X (e.g. **X=Height**). It says that the subpopulation means $\mu_{Y|X=x}$ (the expected values of the outcome Y, as $X=x$ changes), lie on a straight line ("regression line").



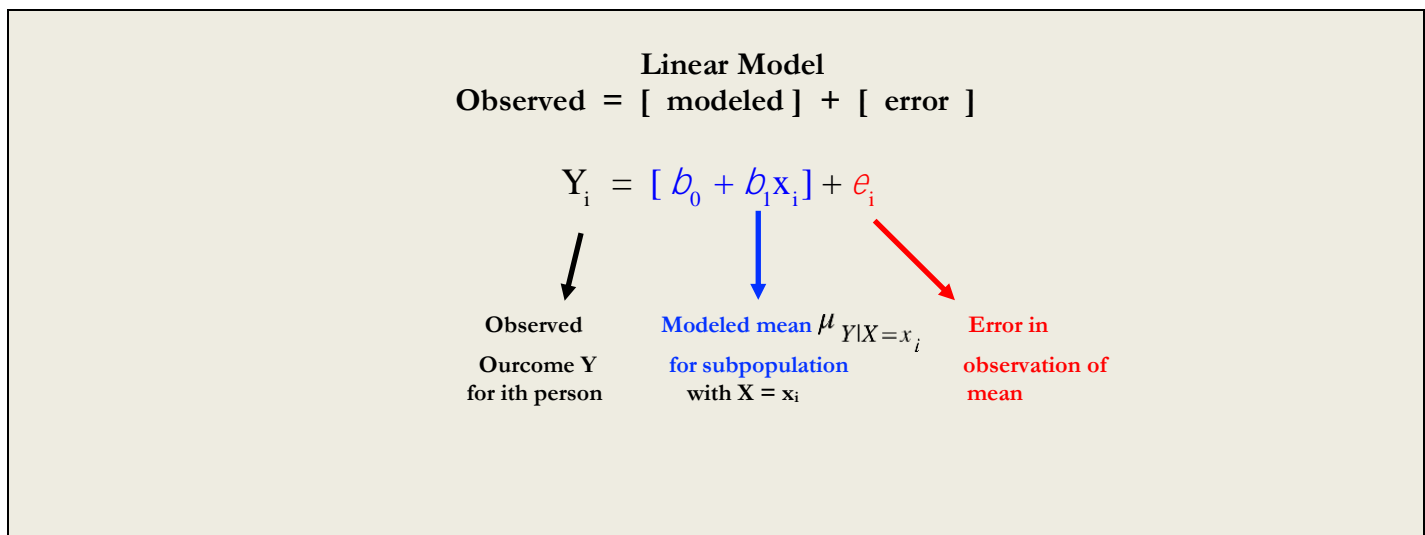
Assumptions of Simple Linear Regression

1. The outcomes Y_1, Y_2, \dots, Y_n are **independent**.
2. The values of the predictor variable X are fixed and measured without error.
3. At each value of the predictor variable $X=x$, the distribution of the outcome Y for the subpopulation with $X=x$ is **normal** with

$$\begin{aligned} \text{mean} &= \mu_{Y|X=x} = \beta_0 + \beta_1 x \\ \text{variance} &= \sigma_{Y|x}^2. \end{aligned}$$

Model

A linear model says “Observed = Model + Error.” These assumptions say that we are modeling the observed outcome for the i th subject as the sum of two pieces: 1) a model piece; plus 2) an error piece.



that is:

$$Y_i = [b_0 + b_1 x_i] + e_i$$

1. The errors e_1, e_2, \dots, e_n are **independent**.
2. Each error e_i is distributed is **normal** with

$$\begin{aligned} \text{mean} &= 0 \\ \text{variance} &= \sigma_{Y|x}^2. \end{aligned}$$

How to estimate β_0, β_1 : “Least Squares”, “Close” and Least Squares Estimation

It’s possible to draw lots of lines through an X-Y scatter of points! So, which one should we choose? “Least squares” estimation is one approach to choosing a line that is “closest” to the data. Least squares estimation says choose the values for $\hat{\beta}_0$ and $\hat{\beta}_1$ that, upon insertion, minimizes the total

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - [\hat{\beta}_0 + \hat{\beta}_1 X_i])^2$$

The total, $\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - [\hat{\beta}_0 + \hat{\beta}_1 X_i])^2$ has a variety of names:

- ◆ residual sum of squares, SSE or SSQ(residual)
- ◆ sum of squares about the regression line
- ◆ sum of squares due error (SSE)

Least Squares Estimation Solutions

Note – the estimates are denoted either using Greek letters with a caret or with Roman letters

Estimate of Slope $\hat{\beta}_1$ or b_1	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$
Intercept $\hat{\beta}_0$ or b_0	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

Analysis of Variance

Partitioning the Total Variance and all things sum of squares and mean squares

Source	df	Sum of Squares A measure of variability	Mean Square = Sum of Squares / df A measure of average/typical/mean variability
Regression due model	1	$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$	MSR = SSR/1
Residual due error	(n-2)	$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	MSE = SSE/(n-2) = $\hat{\sigma}_{Y X}^2$ Note: also called “mean squared error”
Total, corrected	(n-1)	$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$	

Nature — Population/ Sample — Observation/ Data — Relationships/ Modeling — Analysis/ Synthesis

2.

R Illustration: Fit a Simple Linear Regression Model

Preliminary - Set working directory (user edits)

```
setwd("/Users/cbigelow/Desktop/") # setwd( ) to set working directory = folder to read from and write to

Input R dataset janka.Rdata. Inspect.
library(tidyverse) # glimpse() in package {tidyverse}
load(file="janka.Rdata") # Assumes the data are in the working directory
janka$hardness <- as.numeric(janka$hardness)
glimpse(janka) # glimpse( ) to view dataset structure. Could also do str( ) in {base}

## Observations: 36
## Variables: 2
## $ density <dbl> 24.7, 24.8, 27.3, 28.4, 28.4, 29.0, 30.3, 32.7, 35.6, 3...
## $ hardness <dbl> 484, 427, 413, 517, 549, 648, 587, 704, 979, 914, 1070,...
```

janka

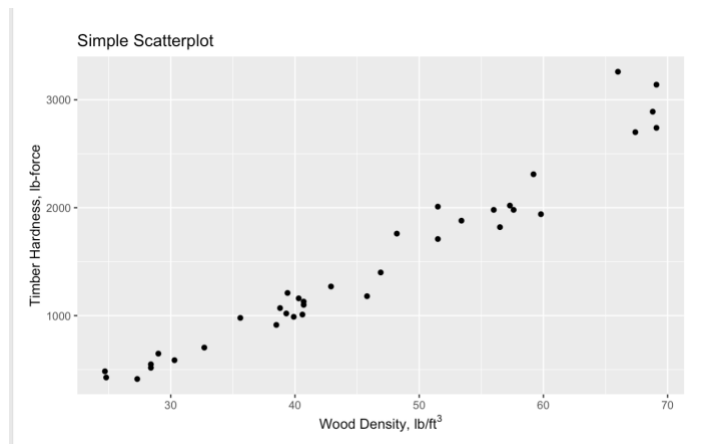
```
##   density hardness
## 1    24.7      484
## 2    24.8      427
## 3    27.3      413
## 4    28.4      517
... Rows omitted ...
## 34   68.8     2890
## 35   69.1     2740
## 36   69.1     3140
```

Descriptives using command stargazer() in package stargazer

```
library(stargazer)
stargazer::stargazer(data=janka, type="text", median=TRUE)
##
## =====
## Statistic N      Mean      St. Dev.  Min    Pctl(25) Median Pctl(75)  Max
## -----
## density  36  45.733    13.580  24.700  37.775  41.800  56.700  69.100
## hardness 36 1,469.472  801.517  413    962.8   1,195   1,980   3,260
## -----
```

Scatterplot using command ggplot() and option geom_point() in package ggplot2

```
library(ggplot2)
library(ggplot2)
ggplot(data=janka) + # required layer: data = to specify dataset
  aes(x=density,y=hardness) + # required layer: aes( ) to define x- and y-axis
  geom_point() + # required layer: geom_point( ) to produce XY scatterplot
  xlab(expression("Wood Density, lb/ft"^{3})) + # optional: Label the x-axis
  ylab("Timber Hardness, lb-force") + # optional: Label the y-axis
  ggtitle("Simple Scatterplot") # optional: provide a title
```



looks linear with no influential observations

Fit Simple Linear Regression. Obtain Coefficients Table. Obtain Analysis of Variance Table

KEY:
 # lm() fits the model. Example: MODELNAME <- lm(data=DATAFRAMENAME, YVARIABLE ~ XPREDICTOR)
 # summary() provides coefficients table and some other info. Example: summary(MODELNAME)
 # anova() produces anova table. Example: anova(MODELNAME)

```
model1 <- lm(data=janka, hardness~density)
summary(model1)
## lm(formula = hardness ~ density, data = janka)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -338.40  -96.98  -15.71   92.71  625.06
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept)  -1160.500     108.580  -10.69  0.000000000000207 ***
## density         57.507       2.279   25.24 < 0.000000000000002 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 183.1 on 34 degrees of freedom
## Multiple R-squared:  0.9493, Adjusted R-squared:  0.9478
## F-statistic: 637 on 1 and 34 DF, p-value: < 0.0000000000000022
```

Intercept = $\hat{\beta}_0 = b_0 = -1160.500$
 Slope = $\hat{\beta}_1 = b_1 = 57.507$

The fitted line is thus: Predicted hardness = hardness = $-1160.500 + 57.507 \cdot \text{density}$

```
anova(model1)
## Analysis of Variance Table
##
## Response: hardness
##              Df    Sum Sq Mean Sq F value      Pr(>F)
## density      1 21345674 21345674  636.98 < 0.0000000000000022 ***
## Residuals  34 1139366   33511
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

SSQ(model) = SSR = 21,345,674
 SSQ(residual) = SSE = 1,139,366
 Overall F-test of null: slope=0
 F = 636.98 with df = 1, 34
 p-value <<< .0001 REJECT null
 Conclude fitted line is significant

3. Multivariable Linear Regression

a. Introduction

In multiple linear regression, the number of explanatory (predictor) variables is > 1 . There is still just one outcome (response) variable Y , continuous and assumed distributed normal. The multiple predictors in a linear regression model can be any mix of continuous or discrete.

Definition

By convention, in multiple predictor linear regression, we say we have p predictors: X_1, X_2, \dots, X_p . A multiple linear regression model is a particular model of how the subpopulation means $\mu_{Y|X_1, X_2, \dots, X_p}$ (the average value) of **one continuous outcome random variable Y** (e.g. **$Y = \text{length of hospital stay}$**) varies, depending on the values of p predictor variables. These can be a mixture of continuous and discrete predictors (e.g. **$X_1 = \text{age}$, $X_2 = 0/1 \text{ history of vertebral fractures, etc.}$**). Because we now have p predictors instead of a single predictor X , a multiple predictor linear regression model says that the subpopulation means of the outcome variable Y , $\mu_{Y|X_1, X_2, \dots, X_p}$, as the profiles of predictors X_1, X_2, \dots etc change, lie on a “plane” (“regression plane”).

Example

P53 is a tumor suppressor gene that has been extensively studied in breast cancer research. Suppose we are interested in understanding the correlates of p53 expression, especially those that are known breast cancer risk variables. We might hypothesize that p53 expression is related to number of pregnancies and age at first pregnancy.

Y = p53 expression level
 X_1 = number of pregnancies (coded 0, 1, 2, etc)
 X_2 = age at first pregnancy ≤ 24 years (1=yes, 0=no)
 X_3 = age at first pregnancy > 24 years (1=yes, 0=no)

This is a multivariable linear model with number of predictors $p = 3$:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \text{error}$$

The General Multivariable Linear Model

Similarly, it is possible to consider a multivariable model that includes p predictors:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \text{error}$$

- $p = \#$ predictors, apart from the intercept
- Each $X_1 \cdots X_p$ can be either discrete or continuous.
- Data are comprised of n data points of the form $(Y_i, X_{1i}, \cdots, X_{pi})$
 Note: The subscript “ i ” is indexing the individual, while the subscripts 1, 2, ..., p are indexing the predictors
- For the i^{th} individual, we have a vector of predictor variable values that is represented $X'_i = [X_{1i}, X_{2i}, \dots, X_{pi}]$

Assumptions

The assumptions required are an extension of those for simple linear regression.

1. The sample size = n observations Y_1, Y_2, \cdots, Y_n are **independent**.
2. The values of the predictor variables $X_1 \cdots X_p$ are **fixed** and measured without error.
3. For each vector value of the predictor variable $\underline{X}=\underline{x}$, the distribution of values of Y is modeled as distributed **normal** distribution with mean equal to $\mu_{Y|\underline{X}=\underline{x}}$ and common variance equal to $\sigma_{Y|\underline{x}}^2$.
4. For each profile of values, x_1, x_2, \dots, x_p , of the p predictor variables $X_1 \cdots X_p$ (written using vector notation $\underline{X}=\underline{x}$), the distribution of values of Y modeled as distributed **normal** with

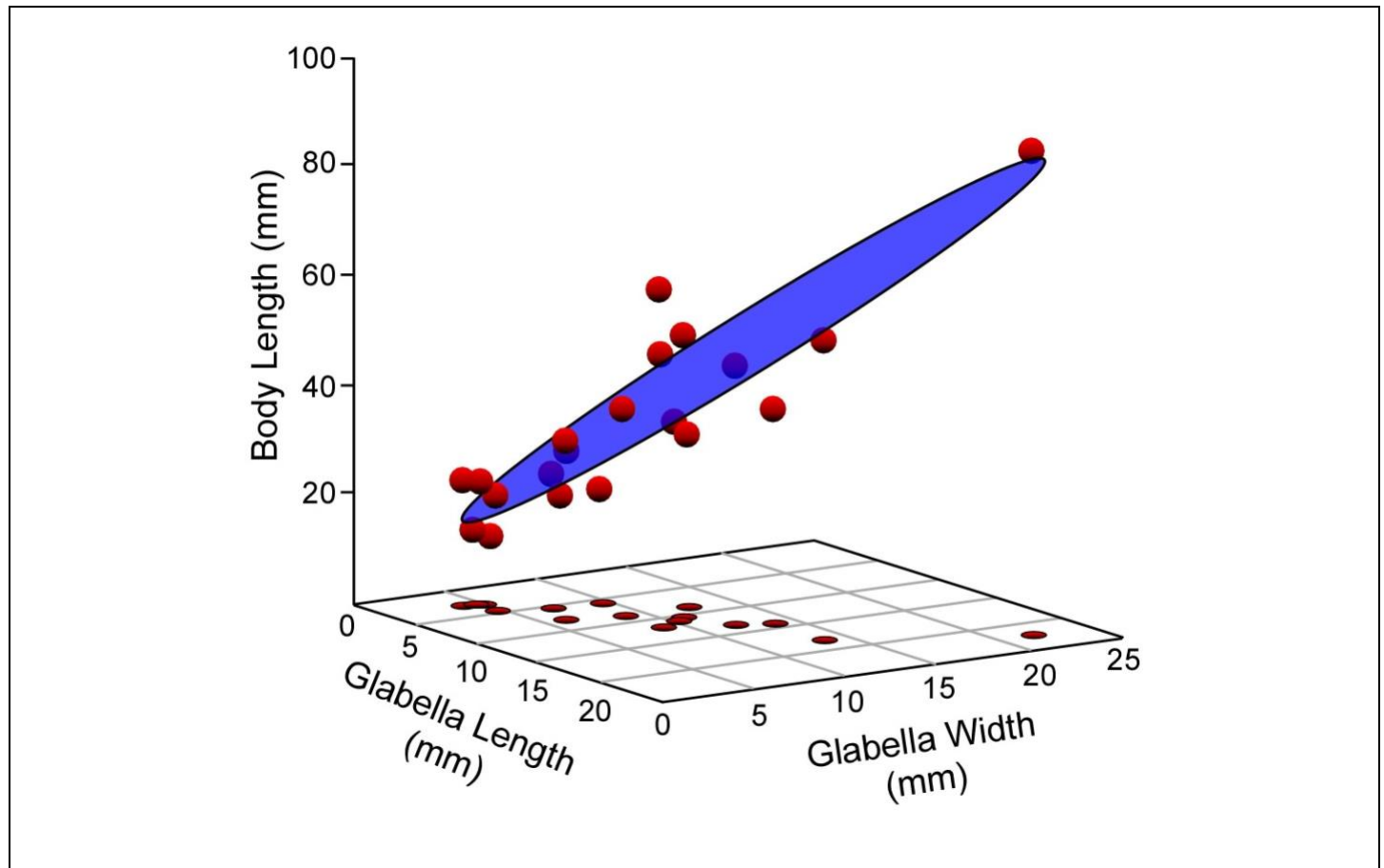
$$\text{mean} = \mu_{Y|\underline{X}=\underline{x}} = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

$$\text{variance} = \sigma_{Y|\underline{X}=\underline{x}}^2.$$

Model Fitting (Estimation)

When there are multiple predictors, the least squares fit is **multi-dimensional**. In the setting of just 2 predictors, it's possible (sort of anyway) to show a schematic of the fitted plane that results from least squares estimation.

Consider the picture below. The outcome (dependent variable) is Y =body length and there are two predictors: X_1 =glabella length and X_2 =glabella width. The purple ellipse is the least squares fit and is a **2-dimensional plane** in 3-dimensional space. It is analogous to the straight line fit that was explained in simple linear regression.



Source: www.palass.org

b. Indicator Variables (also called “dummy variables”) and Design Variables

Why Indicator Variables?

Example - Suppose you want to model some outcome (Y = duration of stay in ICU, in days) in relationship to a nominal predictor, type of surgery X . X might be lazily stored in the data using “1”, “2”, and “3” as placeholders for the names of the type of surgery; e.g., 1=medical therapy, 2=angioplasty, and 3=coronary bypass surgery).

Not really appreciating that “1”, “2”, and “3” are your lazy placeholders/names and not actually bona fide numbers, you might just forge on and fit a simple linear model. **Spoiler alert – the following would be incorrect):**

$$\text{days}_i = [b_0 + b_1 * (\text{type of surgery})_i] + e_i$$

The notion of slope representing the change in Y =days per 1 unit increase in X =type of surgery doesn’t work!

$$\begin{aligned} b_1 &= \text{D } Y \text{ per 1 unit increase in } X, \text{ by definition} \\ &= \text{Predicted change in duration of stay in ICU per 1 unit increase in TYPE OF SURGERY???} \\ &= \text{"makes no sense"} \end{aligned}$$

So, what to do? Answer: 1) we will NOT put X =type of surgery into the model; and 2) instead, we will substitute a set of what are called indicator variables, as described below.

Indicator Variables are Variables that are coded 0 or 1. They are very convenient.

Indicator variables are commonly used as predictors in multivariable regression models. We let

$$\begin{aligned} 1 &= \text{value of indicator when “trait” is present} \\ 0 &= \text{value of indicator when “trait” is not present} \end{aligned}$$

- ◆ The estimated regression coefficient β associated with an indicator variable has a straightforward interpretation, namely:
- ◆ β = predicted change in outcome Y that accompanies presence of “trait”
(estimated change in Y associated with unit change in trait: from “0=absent” to “1=present”)

Examples of Simple Indicator Variables

SEXF = 1 if individual is female
0 otherwise

TREAT = 1 if individual received experimental treatment
0 otherwise

Nature ——— Population/
Sample ——— Observation/
Data ——— Relationships/
Modeling ——— Analysis/
Synthesis

Design variables. To distinguish 2 groups, you need just one separator. This will be one indicator/dummy variable to distinguish the two possibilities (e.g., a 0/1 indicator to distinguish female sex at birth from male sex at birth. note – This is for illustration only; I understand that, in reality, there are yet additional possibilities at birth). To distinguish 3 groups, now you need 2 indicator/dummy variables; for example, for “low”, “medium”, and “high” you need one indicator/dummy variable to distinguish “medium” as being different from “low” and “high” and a 2nd indicator/dummy variable to distinguish “high” as being different from “low” and “medium”). And so on.

If a nominal predictor has k possible values, then you need (k-1) separators. This will be (k-1) indicator/dummy variables to distinguish the k levels. The set of 0/1 indicator variables that you create to distinguish all the separate groups are called **design variables**.

Returning to our Example (Y=duration of stay in ICU, X = type of surgery)

Our original predictor variable X is nominal with 3 possible values:

X = 1 if treatment is medical therapy
 2 if treatment is angioplasty
 3 if treatment is bypass surgery

So, we’ve agreed that we cannot put X = type of surgery into a regression model “as is” because the resulting estimated slope makes no sense. For three surgery types, we need 2 separators. Thus, we create TWO 0/1 indicator/dummy variables: 1) TR_ANG is a 0/1 indicator/dummy variable that “flags” angioplasty; and 2) TR_SUR is a 0/1 indicator/dummy variable that “flags” bypass surgery. Having obtained the required 2 separators (TR_ANG and TR_SUR), we do not need an indicator/dummy variable to “flag” the folks receiving medical therapy. The folks receiving medical therapy, in the presence of these two 0/1 indicator variable “flags”, serve as the “referent”. Specifically, observations for patients who received medical therapy will be uniquely identified because they have value = 0 for both of the 0/1 indicator/dummy variables TR_ANG and TR_SUR:

TR_ANG = 1 if treatment is angioplasty (X=2)
 0 otherwise

TR_SUR = 1 if treatment is bypass surgery (X=3)
 0 otherwise

A set of design variables comprised of (3-1) = 2 indicator variables summarize three possible values of treatment. The reference category is medical therapy.

Value of original X = Type of Surgery	Value of 0/1 Indicator TR_ANG	Value of 0/1 Indicator TR_SUR
X=“1” for “medical”, the “referent”	0	0
X=“2” for “angioplasty”	1	0
X=“3” for “surgery”	0	1

Guidelines for the Definition of Indicator and Design Variables

1) How do you want to define the referent group.

Often this choice will be straightforward. It might be one of the following categories of values of the nominal variable:

- The unexposed
- The placebo
- The standard
- The most frequent

2) K levels of the nominal predictors requires (K-1) separators to distinguish. Thus, need to define (K-1) indicator variables

When the number of levels of the nominal predictor variable = k , define $(k-1)$ indicator variables that will identify persons in each of the separate groups, apart from the reference group.

3) In general (this is not hard and fast), treat the $(k-1)$ design variables as a set. This means that you.

- Enter the set together; and
- Remove the set together; and
- In general, retain all $(k-1)$ of the indicator variables, even when only a subset are significant.

c. Interaction Variables

Previously, we've talked about "effect modification", sometimes called "synergism." It refers to the phenomenon that the nature of an X-Y relationship is *different (meaning the slope is different)*, depending on the level of some third variable which, for now, we'll call Z. In regression, we call this **interaction**.

How to create a predictor that will model the interaction of a continuous predictor X and a 0/1 predictor Z. The solution is straightforward. Use the product of X and Z. Here, I've named this new variable XZ.

$$\text{Interaction of predictor X with third variable Z} = \text{XZ} = \text{X} \times \text{Z}$$

Example: Y = length of stay
 X = age (years)
 Z = 0/1 indicator of history of vertebral fracture (Z=0 for NON fractures and Z=1 for fractures)
 XZ = [X] * [Z] = interaction of X and Z

Our full model is thus the following:

$$Y = b_0 + b_1Z + b_2X + b_3XZ$$

Key to the betas:

- b_0 = intercept for **referent** (the referent group are patients with Z = 0, the non-vertebral fracture folks)
- b_1 = **CHANGE in INTERCEPT** (associated with Z=1, that is - associated with vertebral fracture)
- b_2 = slope of change in Y per unit X for **referent** group
- b_3 = **CHANGE in SLOPE** associated with Z=1 (that is - associated with vertebral fracture)

Try it. What is the model of Y for **non-vertebral fractures** patients (Z=0)?

For the **non-vertebral fractures** patients, insertion of Z=0 yields

$$Y = b_0 + b_2X$$

$$\text{Intercept} = b_0$$

$$\text{Slope} = b_2$$

Try it. What is the model of Y for **vertebral fractures** patients (Z=1)?

For the **vertebral fractures** patients, insertion of Z=1 yields

$$Y = [b_0 + b_1] + [b_2 + b_3]X$$

$$\text{Intercept} = [b_0 + b_1]$$

$$\text{Slope} = [b_2 + b_3]$$

d. *Look!* Schematic of Confounding and Effect Modification

The use of indicator variables and interaction variables are helpful (but not without important caveats) in assessing confounding and effect modification.

Consider a similar regression setting:

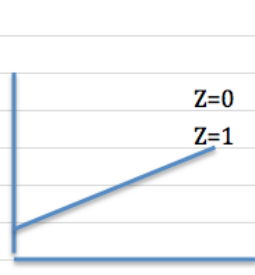
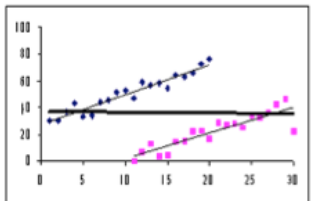
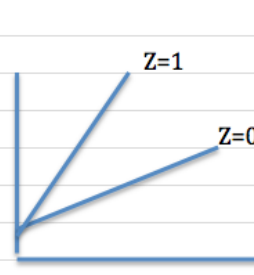
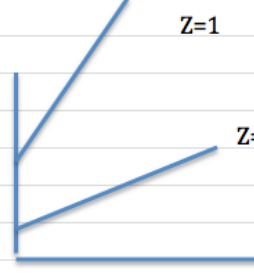
Y = length of hospital stay

X = duration of surgery, continuous

Z = a nominal predictor coded 0 for “no comorbidities” and coded 1 for “one or more comorbidities”.

Associated with Z=1 (the patient has comorbidities), relative to Z=0 (the referent patient with no comorbidities), the X-Y relationship might have a different intercept, or a different slope, or a different intercept and a different slope.

Take a look!

			
Coincident	Confounding (admittedly extreme!)	Effect Modification	Effect Modification
$Y = \beta_0 + \beta_1 X$	$Y = \beta_0 + \beta_1 X + \beta_2 Z$	$Y = \beta_0 + \beta_1 X + \beta_2 XZ$ where $XZ = X * Z$	$Y = \beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ$ where $XZ = X * Z$
Comorbidities=0: $Y = \beta_0 + \beta_1 X$	Comorbidities=0: $Y = \beta_0 + \beta_1 X$	Comorbidities=0: $Y = \beta_0 + \beta_1 X$	Comorbidities=0: $Y = \beta_0 + \beta_1 X$
Comorbidities=1: $Y = \beta_0 + \beta_1 X$	Comorbidities=1: $Y = (\beta_0 + \beta_2) + \beta_1 X$	Comorbidities=1: $Y = \beta_0 + (\beta_1 + \beta_2) X$	Comorbidities=1: $Y = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) X$
	$\beta_2 = \text{change}$ in intercept associated with presence of any comorbidities	$\beta_2 = \text{change}$ in slope of Y on X associated with presence of any comorbidities	$\beta_2 = \text{change}$ in intercept $\beta_3 = \text{change}$ in slope of Y on X

e. The Analysis of Variance Table

The ideas of the analysis of variance table introduced in BIOSTATS 540 (*Unit 12, Simple Linear Regression and Correlation*) apply here, as well. The total variability in the outcome (the total “pie”) is partitioned into its component sources (“wedges” of the pie)

1. **SST:** “Total” or “total, corrected”

- ◆ $SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$ is the variability of Y about \bar{Y}
- ◆ Degrees of freedom = df = (n-1).

2. **SSR** “Regression” or “due model”

- ◆ $SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$ is the variability of \hat{Y} about \bar{Y}
- ◆ Degrees of freedom = df = p = # predictors apart from intercept

3. **SSE:** “Residual” or “due error” refers to the

- ◆ $SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ is the variability of Y about \hat{Y}
- ◆ Degrees of freedom = df = (n-1) – (p)

Source	df	Sum of Squares	Mean Square
Model	p <small>p = # predictors in the model AFTER the intercept</small>	$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$	$MSR = SSR/p$
Residual	(n-1) - p	$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	$MSE = SSE/(n-1-p)$
Total, corrected	(n-1)	$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$	

Overall F Test

The overall F test also applies, yielding an overall F-test to assess the significance of the variance explained by the model. Note that the degrees of freedom is different here; this is because there are now “p” predictors instead of 1 predictor.

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

$$H_A: \text{At least one } \beta_i \neq 0$$

Eureka!!! When the null is true, the best model is “intercept only”

$$F_{\text{OVERALL}} = \frac{\text{mean square due model}}{\text{mean square due residual}} = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}/(p)}{\text{SSE}/(n-1-p)} \quad \text{with df} = p, (n-1-p)$$

Rejection of the null occurs for large values of F_{OVERALL} with accompanying small p-value. With rejection of the null, we conclude at least one predictor (Sigh - we don’t know which ones) has a slope that is statistically significantly different from zero.

Example - Consider a multiple linear regression analysis of the relationship of $Y = \text{p53 expression}$ to $\text{age at first pregnancy (pregnum)}$, 1^{st} pregnancy at age ≤ 24 (**early**), and 1^{st} pregnancy at age > 24 (**late**). The variables **early** and **late** are each 0/1. The referent group is nulliparous.

R illustration

The following assumes that you have downloaded **p53paper.Rdata** from the course website

```
load(file="p53paper.Rdata")
fit <- lm(p53 ~ pregnum + early + late, data=p53paper)
summary(fit)
```

```
##
## Call:
## lm(formula = p53 ~ pregnum + early + late, data = p53paper)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.86030 -0.57031  0.01611  0.51611  2.62100
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2.57031    0.24088   10.671 9.36e-16 ***
## pregnum       0.37641    0.20087    1.874  0.0656 .
## early         0.16076    0.55559    0.289  0.7733
## late        -0.06772    0.50174   -0.135  0.8931
## ---
## The fitted line is:   $\hat{p53} = 2.57 + 0.38 \cdot \text{pregnum} + 0.16 \cdot \text{early} - 0.07 \cdot \text{late}$ 
##
## Residual standard error: 0.9635 on 63 degrees of freedom
## (1 observation deleted due to missingness)
## Multiple R-squared:  0.203, Adjusted R-squared:  0.165
## F-statistic: 5.349 on 3 and 63 DF, p-value: 0.002402 The overall F-test of the null hypothesis of
## zero slopes on every predictor is rejected. Conclude at least one slope is statistically significantly
## different from zero. Upon inspection of the estimates, their standard errors, their t-values, what do you think?
```

```
anova(fit)

## Analysis of Variance Table
##
## Response: p53
##          Df Sum Sq Mean Sq F value    Pr(>F)
## pregnum   1 14.330  14.3301  15.4359 0.0002146 ***
## early     1  0.550   0.5497   0.5921 0.4444682
## late      1  0.017   0.0169   0.0182 0.8930686
## Residuals 63 58.487   0.9284
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

$H_0: \beta_{\text{PREGNUM}} = 0 \text{ and } \beta_{\text{EARLY}} = 0 \text{ and } \beta_{\text{LATE}} = 0$
 $H_A: \text{At least one slope } \neq 0$

$$F_{3,63} = \frac{\text{mean square due model}}{\text{mean square due residual}} = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}/(p)}{\text{SSE}/(n-1-p)}$$

$$F_{3,63} = \frac{\text{msq}(\text{Model})}{\text{msq}(\text{Residual})} = \frac{(14.330 + 0.550 + 0.017)/3}{(58.487)/63} = \frac{4.96557054}{0.9284} = 5.349$$

This matches "F-Statistic" p 23

The overall F-test of the null hypothesis of zero slopes on every predictor is rejected (p-value = .002; see previous page). Conclude at least one slope is statistically significantly different from zero. Important: all we can say at this point, however, is that the model that was fit explains statistically significantly more of the variability in $Y = p53$ than is explained by "no model" at all (the intercept only model).

f. The Partial F Test

The partial F test is used to choose between two models, where one model ("full") is an enhancement of the other model ("reduced"). This type of pairs of models are called hierarchical. We perform a partial F test on hierarchical models when we want to control for some predictors and then determine if the "extra" predictors are statistically significant, "above and beyond" the control variables.

What if we want to compare and choose between two models?

There are a variety of ways to do this. One way is to do a partial F Test. A partial F test is a statistical technique for comparing two models that are "hierarchical." It permits the assessment of associations while controlling for confounding.

Some more details of hierarchical models.

- **"Hierarchical"** means one model is an enhancement of the other. The smaller model has various names: "reduced", "reference", "smaller". When you enhance it, you keep all the predictors in the smaller model, but then you add some additional predictors. The larger (enhanced) model has various names: "full", "comparison", "larger"
- Thus, "hierarchical" means that all of the predictors in the smaller (reduced, reference) are contained in the larger (comparison) model.
- In the **Y = p53** example, we might be interested in comparing the following two hierarchical models:
 Predictors in smaller model = { **pregnum** }
 Predictors in larger model = { **pregnum** } + { **early** + **late** }
- "Hierarchical" is satisfied because all of the predictors (here there is just one - **pregnum**) that are contained in the smaller model are contained in the larger model.
- In a partial F test, we are assessing the nature and significance of the extra predictors, (**early** and **late**) for the prediction of **Y=p53**, adjusting for (controlling for) all of the variables in the smaller model (**pregnum**).

Thus, the comparison of the hierarchical models is addressing the following question:

What is the statistical significance of **early** and **late** for the prediction of **Y = p53**, after controlling for the association of **Y=p53** with the control variable **pregnum**?

Statistical Definition of the Partial F Test

Research Question: Does inclusion of the “*extra*” predictors explain significantly more of the variability in outcome compared to the variability that is explained by the predictors that are already in the model?

Partial F Test

H₀: Addition of $X_{p+1} \dots X_{p+k}$ is of no statistical significance for the prediction of Y after controlling for the predictors $X_1 \dots X_p$ meaning that:

$$\beta_{p+1} = \beta_{p+2} = \dots = \beta_{p+k} = 0 \text{ after adjustment for } X_1 \dots X_p$$

H_A: Not

$$F_{\text{PARTIAL}} = \frac{\{ \text{Extra regression sum of squares} \} / \{ \text{Extra regression df} \}}{\{ \text{Residual sum of squares larger model} \} / \{ \text{Residual df larger model} \}}$$

$$= \frac{[SSR(X_1 \dots X_p, X_{p+1} \dots X_{p+k}) - SSR(X_1 \dots X_p)] / [(p+k) - p]}{[SSE(X_1 \dots X_p, X_{p+1} \dots X_{p+k})] / [(n-1) - (p+k)]}$$

$$\begin{aligned} \text{Numerator df} &= (p+k) - (p) = k \\ \text{Denominator df} &= (n-1) - (p+k) \end{aligned}$$

H₀ true: The extra predictors are not significant in adjusted analysis	F statistic = small (close to 1) p-value = large
H₀ false: The extra predictors are significant in adjusted analysis	F statistic = large (bigger than 1) p-value = small

R illustration Example – continued.

<pre>reduced <- lm(data=p53paper, p53 ~ pregnum) full <- lm(data=p53paper, p53 ~ pregnum + early + late) anova(reduced, full)</pre>	<p>H₀: Controlling for pregnum, the additional predictors have $\beta_{\text{EARLY}} = 0$ and $\beta_{\text{LATE}} = 0$</p> <p>H_A: At least one extra predictor is of “ADDED” significance, after adjustment (controlling for) pregnum</p>
<pre>Analysis of Variance Table ## Model 1: p53 ~ pregnum ## Model 2: p53 ~ pregnum + early + late ## Res.Df RSS Df Sum of Sq F Pr(>F) ## 1 65 59.054 ## 2 63 58.487 2 0.56663 0.3052 0.7381</pre>	<p>telling us → F_{partial} = .3052 p-value = .7361</p>

The null hypothesis is NOT rejected (p-value = .74). Conclude that early and late are not statistically significant for the prediction of Y=p53 after adjustment for the control variable pregnum. Specifically, addition of early and late to the model does not explain statistically significantly more of the variability in Y=p53 beyond that explained by pregnum.

g. Multiple Partial Correlation

Beware. Partial F test \neq partial correlation

- The partial F test is a hypothesis test; whereas.
- A partial correlation is a statistic, measuring the what is explained (and expressed as a percent if squared)

Partial correlation. *“To what extent is Y correlated with X (or multiple X), after accounting for some control variable Z (or multiple control variables Z)?”*

In a partial correlation, we are removing the influence of the control variable (Z). A partial correlation is the correlation of (residuals of Y on Z) with the (residuals of X on Z). To appreciate what this means, consider:

- **Preliminary 1:** Regress the predictor X on the control variable Z
 - Obtain the residuals
 - These residuals represent the information in the predictor X that is independent of Z
- **Preliminary 2:** Now regress the outcome Y on the control variable Z
 - Obtain the residuals
 - These residuals represent the information in Y that is independent of Z
- The partial correlation of Y on X controlling for Z as the correlation between these two sets of residuals: (residuals of Y on Z) and (residuals of X on Z) give you a Z-controlled assessment of the relationship between X and Y, that is, *independent of Z*.

Partial Correlation

As a correlation

$R_{XY|Z}$ = Multiple Partial correlation (X,Y | controlling for Z)

= Correlation (residuals of X regressed on Z, residuals of Y regressed on Z)

As a squared correlation

$R^2_{XY|Z}$ = Multiple Squared Partial correlation (X,Y | controlling for Z)

$$= \frac{\text{SSR}(\text{due Model with Z and X}) - \text{SSR}(\text{due Model with Z alone})}{\text{SSE}(\text{due residual in Z only model})}$$

Putting this all together, and keeping track of the distinctions ...

F_{partial} = Partial F Test	R^2_{partial} = Partial Multiple Correlation Squared
Goal: Hypothesis test of significance of extra variables, after adjustment for the control variables.	Goal: Estimation of percent of variability in outcome Y that is explained by the extra variables, independent of the control variables.
Control variables: $X_1 \dots X_p$ Extra variables: $X_{p+1} \dots X_{p+k}$	Control variables: $X_1 \dots X_p$ Extra variables: $X_{p+1} \dots X_{p+k}$
F_{PARTIAL} hypothesis test compares mean squares to mean squares	R^2_{partial} multiple partial correlation squared compares sum of squares to sum of squares
The denominator has the FULL model	The denominator has the REDUCED model
$= \frac{[SSR(X_1 \dots X_p, X_{p+1}, \dots, X_{p+k}) - SSR(X_1 \dots X_p)] / [(p+k) - p]}{[SSE(X_1 \dots X_p, X_{p+1}, \dots, X_{p+k})] / [(n-1) - (p+k)]}$	$= \frac{SSR(\text{due Model with all}) - SSR(\text{due Model control only})}{SSE(\text{due residual in Z only model})}$

4. Multivariable Model Development

a. Introduction

Recall from page 7 The **goal** of statistical modeling is to obtain a model that is simultaneously **minimally adequate** and a **good fit**. **And the model should make sense.**

Recall. Some general guidelines (note – there is no single right answer)

Preliminary –

Be sure you have: (1) checked, cleaned and described your data, (2) screened the data for multivariate associations, and (3) thoroughly explored the bivariate relationships.

Step 1 –

Fit the “maximal” model..

Step 2 –

Begin simplifying the model.

Step 3 –

Keep simplifying the model.

Repeat step 2, over and over, until the model remaining contains nothing but significant predictor variables.

Then there is a Step 4 -

Perform regression diagnostics

We'll get to this later, *Section 5. Goodness-of-Fit and Regression Diagnostics*

b. Example

Framingham Study

Source:

Levy (1999) National Heart Lung and Blood Institute. Center for Bio-Medical Communication.
Framingham Heart Study

Description:

Cardiovascular disease (CVD) is the leading cause of death and serious illness in the United States. In 1948, the Framingham Heart Study, under the direction of the National Heart Institute (now known as the National Heart, Lung, and Blood Institute or NHLBI) was initiated. The objective of the Framingham Heart Study was to identify the common factors or characteristics that contribute to CVD by following its development over a long period of time in a large group of participants who had not yet developed overt symptoms of CVD or suffered a heart attack or stroke.

Here we use a subset of the data, $n=1000$.

Variable	Label	Codings
sbp	Systolic Blood Pressure (mm Hg)	
ln_sbp	Natural logarithm of sbp	$\ln_sbp = \ln(sbp)$
age	Age, years	
bmi	Body Mass index (kg/m ²)	
ln_bmi	Natural logarithm of bmi	$\ln_bmi = \ln(bmi)$
sex	Gender	1=male 2=female
female	Female Indicator	0 = male 1 = female
scl	Serum Cholesterol (mg/100 ml)	
ln_scl	Natural logarithm of scl	$\ln_scl = \ln(scl)$

Multiple Regression Variables:

Outcome $Y = \ln_sbp$

Predictor Variables: \ln_bmi , \ln_scl , age, sex

Research Question:

From among these 4 “candidate” predictors, what are the important “risk” factors and what is the nature of their association with $Y = \ln_sbp$?

Input Data. Check. Produce descriptives:

User edits

```
rm(list=ls())
setwd("/Users/cbigelow/Desktop/")

load(file="framingham_1000.Rdata")

framingham <- framingham_1000

summary(framingham)

##      sex      sbp      scl      age
## Men :443   Min.   : 80.0   Min.   :115.0   Min.   :30.00
## Women:557 1st Qu.:116.0   1st Qu.:197.0   1st Qu.:38.75
##          Median :128.0   Median :225.0   Median :45.00
##          Mean   :132.3   Mean   :227.8   Mean   :45.92
##          3rd Qu.:144.0   3rd Qu.:255.0   3rd Qu.:53.00
##          Max.   :270.0   Max.   :493.0   Max.   :66.00
##
##          bmi      id      ln_bmi      ln_sbp
## Min.   :16.40   Min.   : 1   Min.   :2.797   Min.   :4.382
## 1st Qu.:23.00   1st Qu.:1246 1st Qu.:3.135   1st Qu.:4.754
## Median :25.10   Median :2488 Median :3.223   Median :4.852
## Mean   :25.57   Mean   :2410 Mean   :3.230   Mean   :4.872
## 3rd Qu.:27.80   3rd Qu.:3605 3rd Qu.:3.325   3rd Qu.:4.970
## Max.   :43.40   Max.   :4697 Max.   :3.770   Max.   :5.598
## NA's :2      NA's :2
##      ln_scl
## Min.   :4.745
## 1st Qu.:5.283
## Median :5.416
## Mean   :5.410
## 3rd Qu.:5.541
## Max.   :6.201
## NA's :4
```

Clear the environment (workspace)
Tell R where to "read from" and "write to"

I'm lazy. So, I'm creating a shorter name

Inspect distributions of all study variables

There are 4 missing values of scl

There are 2 missing values of bmi, ln_bmi

There are 4 missing values of ln_scl

```
library(stargazer)
stargazer::stargazer(framingham, type="text", median=TRUE)
```

Another way to inspect distributions (more succinct).

```
##
## =====
## Statistic   N      Mean    St. Dev.   Min   Median   Max
## -----
## sbp        1,000  132.350   23.043     80    128    270
## scl         996  227.846   45.087    115    225    493
## age        1,000  45.922    8.545     30     45     66
## bmi         998  25.566    3.848    16.400 25.100 43.400
## id         1,000 2,410.031 1,363.439     1  2,487.5 4,697
## ln_bmi      998   3.230    0.147    2.797  3.223  3.770
## ln_sbp      1,000  4.872    0.163    4.382  4.852  5.598
## ln_scl      996   5.410    0.195    4.745  5.416  6.201
## -----
```

Nicer layout, slightly different info

Examination of the ranges of systolic bp, age, bmi look to be all plausible; no suggestion of significant errors in the data itself.

```
library(summarytools) # freq() in package {summarytools}
summarytools::freq(framingham$sex)

## Frequencies
## framingham$sex
## Type: Factor (unordered)
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##      Men    443    44.30      44.30    44.30    44.30
##      Women  557    55.70     100.00    55.70   100.00
##      <NA>     0    100.00     100.00    0.00   100.00
##      Total 1000    100.00     100.00   100.00   100.00

library(summarytools) # descr() in package {summarytools}
summarytools::descr(framingham$sbp, stats = c("n.valid", "mean", "sd", "min", "q1", "med", "q3", "max", "CV"),
transpose = TRUE) # option stats=c( ) to choose statistics to show

## Descriptive Statistics
## framingham$sbp
## N: 1000
##
##          N.Valid  Mean  Std.Dev  Min  Q1  Median  Q3  Max  CV
## -----
##      sbp  1000.00  132.35   23.04   80.00  116.00  128.00  144.00  270.00  0.17
```

Assess Normality of Candidate Dependent Variable = sbp. Shapiro-Wilk Test (Null: normality)
Histogram w Overlay Normal and QQ Plot

```
options(scipen=1000)
shapiro.test(framingham$sbp)

##
## Shapiro-Wilk normality test
##
## data:  framingham$sbp
## W = 0.92121, p-value < 0.0000000000000022
```

Interpretation: The null hypothesis of normality of the distribution of sbp is rejected ($p < .00001$)

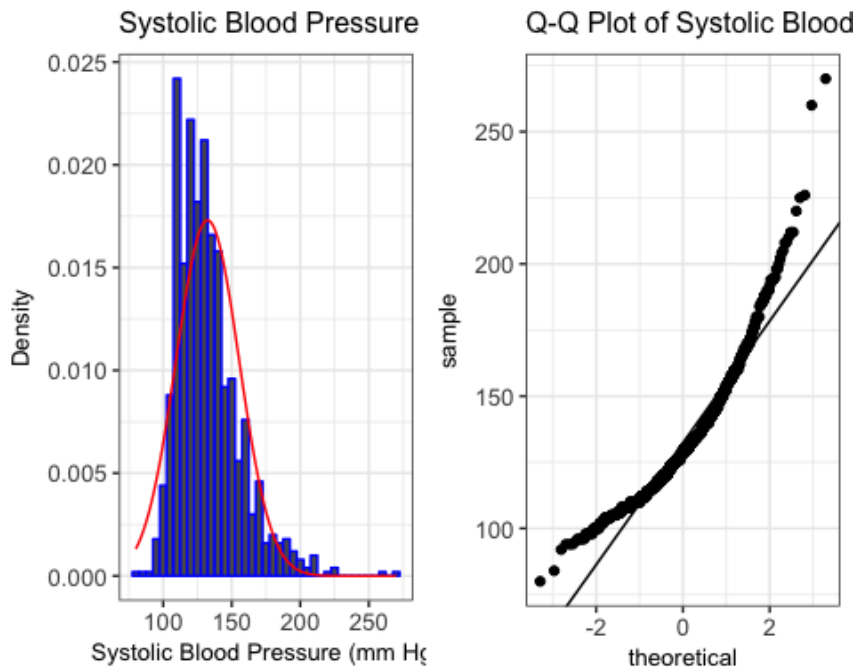
```
library(ggplot2)
library(gridExtra)

# p1 is panel 1 = histogram w overlay normal
p1 <- ggplot(data=framingham, aes(x=sbp)) +
  geom_histogram(binwidth=5, colour="blue",
    aes(y=..density..)) +
  stat_function(fun=dnorm,
    color="red",
    args=list(mean=mean(framingham$sbp),
      sd=sd(framingham$sbp))) +
  ggtitle("Systolic Blood Pressure (sbp)") +
  xlab("Systolic Blood Pressure (mm Hg)") +
  ylab("Density") +
  theme_bw() +
  theme(axis.text = element_text(size = 10),
    axis.title = element_text(size = 10),
    plot.title = element_text(size = 12))
```



```
# p2 is panel 2 = quantile-quantile plot
p2 <- ggplot(data=framingham, aes(sample=sbp)) +
  stat_qq() +
  geom_abline(intercept=mean(framingham$sbp), slope = sd(framingham$sbp)) +
  ggtitle("Q-Q Plot of Systolic Blood Pressure (sbp)") +
  theme_bw() +
  theme(axis.text = element_text(size = 10),
        axis.title = element_text(size = 10),
        plot.title = element_text(size = 12))

gridExtra::grid.arrange(p1, p2, ncol=2) # grid.arrange( ) in package {gridExtra} to lay out panels in figure
```



Interpretation: This confirms what the Shapiro Wilk test suggests. The null hypothesis of normality of the distribution of sbp is not supported.

Create “regression-friendly” indicator variables and interactions. Check.

```
library(summarytools)
library(Hmisc)
# Create 0/1 indicator/dummy variable using logical operator:
# If sex="Women" is TRUE, code new variable female=1. Otherwise, code new variable female=0
# option na.rm=TRUE ensures that missing values will not be considered and instead will be retained as missing.
framingham$female <- as.numeric(framingham$sex == "Women", na.rm=TRUE)
summarytools::cTable(framingham$sex, framingham$female, prop = 'n', totals = FALSE) # xtab check

## Cross-Tabulation
## Variables: sex * female
## Data Frame: framingham
## -----
##      female      0      1
## sex
## Men      443      0
## Women     0     557
## -----
```

female is the new indicator variable created and is coded 0/1
sex is the original variable used to create female

It worked!

```
Hmisc::label(framingham$female) <- "female01" # Label( ) in package {Hmisc} to label variables

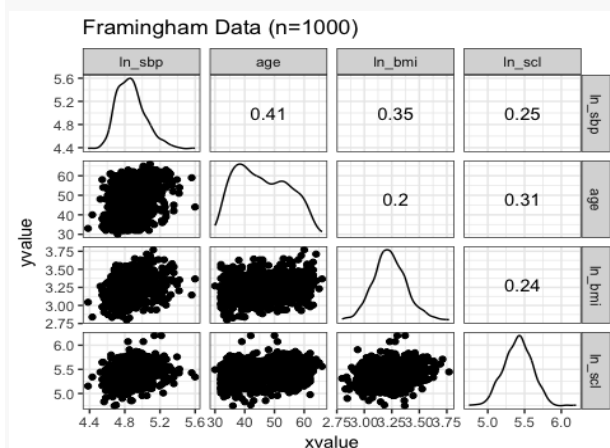
framingham$ageXfemale <- framingham$age*framingham$female
Hmisc::label(framingham$ageXfemale) <- "AGE x FEMALE interaction"

framingham$lnsclXfemale <- framingham$ln_scl*framingham$female
Hmisc::label(framingham$lnsclXfemale) <- "ln(scl) x FEMALE interaction"

framingham$lnbmiXfemale <- framingham$ln_bmi*framingham$female
Hmisc::label(framingham$lnbmiXfemale) <- "ln(bmi) x FEMALE interaction"
```

Examine Pairwise Relationships: 1) Y with X's; and 2) X's with X's

```
library(GGally)
GGally::ggscatmat(data=framingham,
                  columns=c("ln_sbp", "age", "ln_bmi", "ln_scl")) +
  ggtitle("Framingham Data (n=1000)") +
  theme_bw()
```



Create a dataset that has no missing values on any variables of interest. Name this dataset complete.

Then fit the following five (5) models named as follows

m_maximal: Contains all predictors

m_2: Drops 2 interactions - lnbmiXfemale and lnsclXfemale

m_3: One predictor model w predictor = ln_bmi

m_4: One predictor model w predictor = ln_scl

m_5: Three predictor model w predictors = age, female, and ageXfemale

```
library(stargazer)
```

```
# na.omit( ) to omit observations with anything missing; the resulting object named complete contains complete data only
# cols=c("var1", "var2", etc) to specify variables to keep
complete <- na.omit(framingham, cols=c("ln_sbp", "ln_bmi", "age", "female", "lnbmiXfemale", "lnsclXfemale", "ageXfemale"))

# Fit each model of interest to the SAME dataset comprised of complete data only
m_maximal <- lm(data=complete, ln_sbp ~ ln_bmi + ln_scl + age + female + lnbmiXfemale + lnsclXfemale + ageXfemale)
m_2 <- lm(data=complete, ln_sbp ~ ln_bmi + ln_scl + age + female + ageXfemale)
m_3 <- lm(data=complete, ln_sbp ~ ln_bmi)
m_4 <- lm(data=complete, ln_sbp ~ ln_scl)
m_5 <- lm(data=complete, ln_sbp ~ age + female + ageXfemale)
```

```
# stargazer( ) in package {stargazer} for nice display of models side by side
stargazer::stargazer(m_maximal,m_2,m_3,m_4,m_5,type="text",font.size="small", align=TRUE, omit.stat=c("f", "ser"))
```

Dependent variable:					
	(1)	(2)	ln_sbp (3)	(4)	(5)
ln_bmi	0.304*** (0.055)	0.271*** (0.032)	0.388*** (0.033)		
ln_scl	0.059 (0.037)	0.056** (0.025)		0.211*** (0.026)	
age	0.004*** (0.001)	0.004*** (0.001)			0.004*** (0.001)
female	-0.011 (0.304)	-0.217*** (0.051)			-0.327*** (0.051)
lnbmiXfemale	-0.051 (0.067)				
lnsclXfemale	-0.009 (0.050)				
ageXfemale	0.005*** (0.001)	0.005*** (0.001)			0.007*** (0.001)
Constant	3.396*** (0.234)	3.521*** (0.159)	3.618*** (0.106)	3.730*** (0.139)	4.701*** (0.039)
Observations	994	994	994	994	994
R2	0.267	0.266	0.123	0.064	0.203
Adjusted R2	0.261	0.262	0.122	0.063	0.200

Models 1 & 2 have nearly identical R² = % variance explained (26.7%, 26.6%). This suggests the extra predictors in model 1 are not needed. -> Model 2 is preferred (simpler!)

Note: *p<0.1; **p<0.05; ***p<0.01

```
# anova(reduced,full) to obtain Partial F Test
paste("Partial F-test, 2df: Null: lnbmiXfemale=0 lnsclXfemale=0")
anova(m_2, m_maximal)
```

```
[1] "Partial F-test, 2df: Null: lnbmiXfemale=0 lnsclXfemale=0"
Analysis of Variance Table
```

```
Model 1: ln_sbp ~ ln_bmi + ln_scl + age + female + ageXfemale
Model 2: ln_sbp ~ ln_bmi + ln_scl + age + female + lnbmiXfemale + lnsclXfemale + ageXfemale
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	988	19.314				
2	986	19.301	2	0.013173	0.3365	0.7144

Interpretation - This confirms that it is okay to DROP lnbmi_female and lnscl_female (Partial F = 0.34, p-value = .71) nsSo, model 2 is our "tentative" final model

Further work, regression diagnostics, are needed next (See, section 5. *Goodness-of-Fit and Regression Diagnostics*).

c. Suggested Criteria for Confounding and Interaction

A Suggested Statistical Criterion for Determination of Confounding

A variable Z might be judged to be a confounder of an X-Y relationship if **BOTH** of the following are satisfied:

- 1) Its inclusion in a model that already contains X as a predictor has adjusted significance level $< .10$ or $< .05$; and
- 2) Its inclusion in the model changes the estimated regression coefficient for X by 15-20% or more, relative to the model that contains only X as a predictor.

A Suggested Statistical Criterion for Assessment of Interaction

A “candidate” interaction variable might be judged to be worth retaining in the model if **BOTH** of the following are satisfied:

- 1) The partial F test for its inclusion has significance level $< .05$; and
- 2) Its inclusion in the model alters the estimated regression coefficient for the main effects by 15-20% or more.

d. Additional Tips for Multivariable Analysis of Large Data Sets

#1. State the Research Questions.

Aim for a focus that is explicit, complete, and focused, including:

- ◆ Statement of population
- ◆ Definition of outcome
- ◆ Specification of hypotheses (predictor-outcome relationships)
- ◆ Identification of (including nature of) hypothesized covariate relationships

#2. Define the Analysis Variables.

For each research question, note for each analysis variable, its hypothesized role.

- ◆ Outcome
- ◆ Predictor
- ◆ Confounder
- ◆ Effect Modifier
- ◆ Intermediary (also called intervening)

#3. Prepare a “Clean” Data Set Ready for Analysis (Data Management)

For each variable, check its distribution, especially:

- ◆ Completeness
- ◆ Occurrence of logical errors
- ◆ Within form consistency
- ◆ Between form consistency
- ◆ Range

#4. Describe the Analysis Sample

This description serves three purposes:

- 1) Identifies the population actually represented by the sample
- 2) Defines the range(s) of relationships that can be explored
- 3) Identifies, tentatively, the function form of the relationships

Methods include:

- ◆ Frequency distributions for discrete variables
- ◆ Mean, standard deviation, percentiles for continuous variables
- ◆ Bar charts
- ◆ Box and whisker plots
- ◆ Scatter plots

#5. Assessment of Confounding

The identification of confounders is needed for the correct interpretation of the predictor-outcome relationships. Confounders need to be controlled in analyses of predictor-outcome relationships.

Methods include:

- ◆ Cross-tabulations and single predictor regression models to determine whether suspected confounders are predictive of outcome and are related to the predictor of interest.
- ◆ This step should include a determination that there is a confounder-exposure relationship among controls.

#6. Single Predictor Regression Model Analyses

The fit of these models identifies the nature and magnitude of crude associations. It also permits assessment of the appropriateness of the assumed functional form of the predictor-outcome relationship.

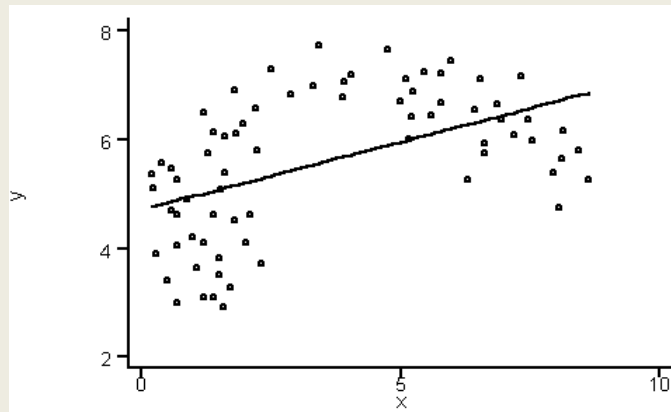
- ◆ Cross-tabulations
- ◆ Graphical displays (Scatter plots)
- ◆ Estimation of single predictor models



5. Goodness-of-Fit and Regression Diagnostics

a. Introduction and Terminology

Neither prediction nor estimation has meaning when the estimated model is a poor fit to the data:



What does this picture suggest?

- ◆ A better fitting relationship between X and Y is quadratic
- ◆ We notice different sizes of discrepancies; in particular:
- ◆ Some observed Y are close to the fitted line \hat{Y} (e.g. near X=1 or X=8)
- ◆ Other observed Y are very far from the fitted line \hat{Y} (e.g. near X=5)

Poor fits of the data to a fitted line can occur for several reasons and can occur even when the fitted line explains a large proportion (R^2) of the total variability in response:

- ◆ The wrong functional form (*more on this later*) was fit.
- ◆ Extreme values (outliers) exhibit uniquely large discrepancies between observed and fitted values.
- ◆ One or more important explanatory variables have been omitted.
- ◆ One or more model assumptions have been violated.

Consequences of a poor fit include:

- ◆ We learn the wrong biology.
- ◆ Comparison of group differences aren't "fair" because they are unduly influenced by a minority.
- ◆ Comparison of group means aren't "fair" because we used the wrong standard error.
- ◆ Predictions are wrong because the fitted model does not apply to the case of interest.

Available techniques of goodness-of-fit assessment are of two types:

1. **Systematic** - those that explore the appropriateness of the model itself

Have we fit the correct model?

Should we fit another model?

2. **Case Analysis** – those that investigate the influence of individual data points

Are there a small number of individuals whose inclusion in the analysis influences excessively the choice of the fitted model?

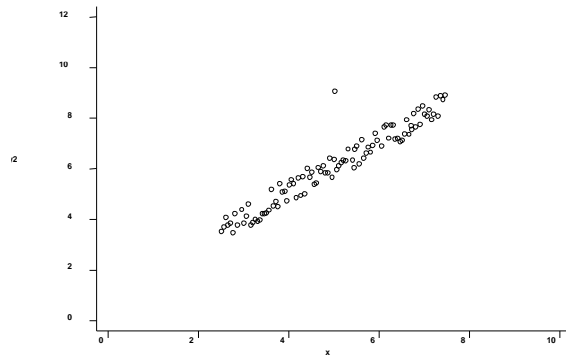
Goodness-of-Fit Assessment

Some Terminology - continued

Case Analysis

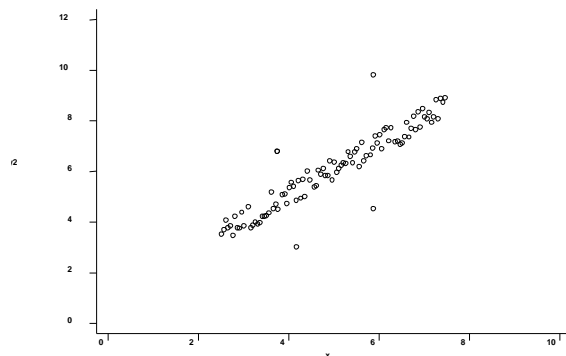
Residual:	<p>The residual is the difference between the observed outcome Y and the fitted outcome \hat{Y}.</p> $e = [Y - \hat{Y}]$ <p>It estimates the unobservable error ϵ.</p>
Outlier:	<p>An outlier is a residual that is <u>unusually</u> large.</p> <p><i>Note:</i> As before, we will rescale the sizes of the residuals via standardization so that we can interpret their magnitudes on the scale of SE units.</p>
Leverage:	<p>The leverage is a measure of the unusualness of the value of the predictor X.</p> <p>Leverage = distance (observed X, center of X in sample)</p> <p>Predictor values with high leverages have, potentially, a large influence on the choice of the fitted model.</p>
Influence:	<p>Measures of influence gauge the change in the fitted model with the omission of the data point.</p> <p>Example: Cook's Distance</p>

A Feel for Residual, Leverage, Influence Large residuals may or may not be influential



Large residual
Low leverage

The large residual effects a large influence.

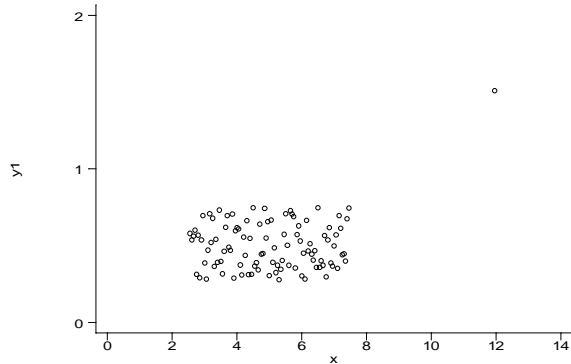


Large residual
Low leverage

Despite its size, the large residual effects only small influence.

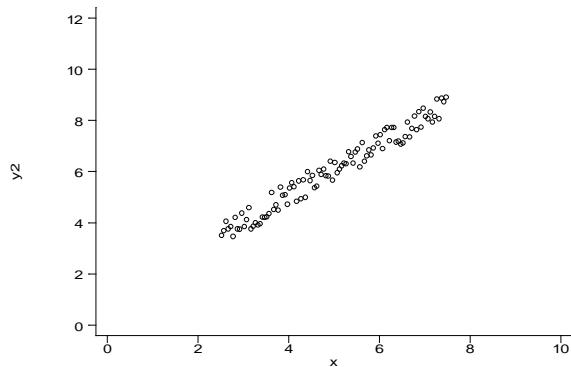
A Feel for Residual, Leverage, Influence

High leverage may or may not be influential



High leverage
Small residual

The high leverage effects a large influence.



High leverage
Small residual

Despite its size, the large leverage effects only small influence.

Thus, case analysis is needed to discover all of:

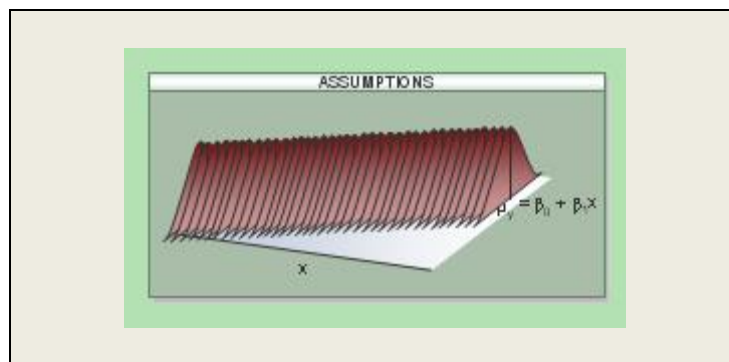
- ◆ high leverage
- ◆ large residuals
- ◆ large influence

b. Assessment of Normality

Recall what we are assuming with respect to normality:

- **Simple Linear Regression:**
Subpopulations of Y are defined by each level, $X = x$. At each level “ x ” of the predictor variable X , the outcomes Y_x are modeled as distributed normal with mean =
 $\mu_{Y|x} = \beta_0 + \beta_1 x$ and constant variance $\sigma_{Y|x}^2$
- **Multiple Linear Regression:**
At each vector level “ $\underline{x} = [x_1, x_2, \dots, x_p]$ ” of the predictor vector \underline{X} , the outcomes $Y_{\underline{x}}$ are modeled as distributed normal with mean = $\mu_{Y|\underline{x}} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$ and constant variance $\sigma_{Y|\underline{x}}^2$

This is what it looks like (courtesy of a picture on the web!)



Violations of Normality are sometimes, but not always, a serious problem

- **When not to worry:** Estimation and hypothesis tests of regression parameters are fairly robust to modest violations of normality
- **When to worry:** Predictions are sensitive to violations of normality
- **Beware:** Sometimes the cure for violations of normality is worse than the problem.

Some graphical assessments of normality and what to watch out for:

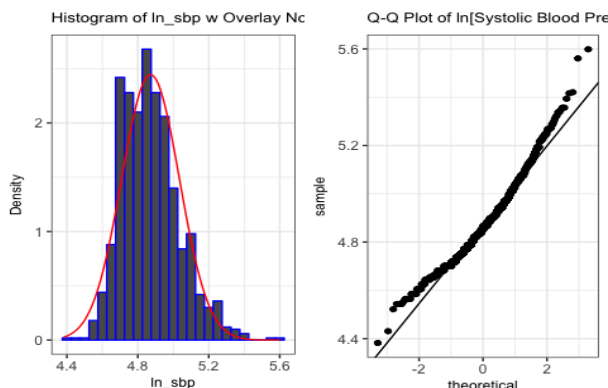
Method	What to watch out for:
1. Histogram of outcome variable Y and/or Histogram of residuals	Look for normal shape of the histogram.
2. Histogram of residuals (or studentized or jackknife residuals)	Look for normal shape of the histogram.
3. Quantile quantile plot of the quantiles of the residuals versus the quantiles of the assumed normal distribution of the residuals.	Normally distributed residuals will appear, approximately, linear.

Two Panel Graph: 1) Histogram w Overlay Normal + 2) QQ plot

```
library(ggplot2)
library(gridExtra)
# Left Panel
# ggplot(data= DATAFRAME, aes(x=VARIABLENAME)) + geom_hisotgram() + stat_function( ) + options
p1 <- ggplot(data=framingham, aes(x=ln_sbp)) +
  geom_histogram(binwidth=.05, colour="blue", # TIP - You may want to tweak binwidth =
    aes(y=..density..)) +
  stat_function(fun=dnorm,
    color="red",
    args=list(mean=mean(framingham$ln_sbp),
      sd=sd(framingham$ln_sbp))) +
  ggtitle("Histogram of ln_sbp w Overlay Normal") +
  xlab("ln_sbp") +
  ylab("Density") +
  theme_bw() +
  theme(axis.text = element_text(size = 9),
    axis.title = element_text(size = 9),
    plot.title = element_text(size = 10))

# Right Panel
p2 <- ggplot(data=framingham, aes(sample=ln_sbp)) +
  stat_qq() +
  geom_abline(intercept=mean(framingham$ln_sbp), slope = sd(framingham$ln_sbp)) +
  ggtitle("Q-Q Plot of ln[Systolic Blood Pressure (ln_sbp)]") +
  theme_bw() +
  theme(axis.text = element_text(size = 9),
    axis.title = element_text(size = 9),
    plot.title = element_text(size = 10))
```

```
gridExtra::grid.arrange(p1, p2, ncol=2)
```



Nature — Population/ Sample — Observation/ Data — Relationships/ Modeling — Analysis/ Synthesis

Skewness and Kurtosis Statistics for Assessing Normality:

	What to watch out for:
<p>Skewness - symmetry of the curve</p> <p>Standardization of the 3rd sample moment about the mean</p> $m_2 = E \left[(Y - \bar{m})^2 \right]$ $m_3 = E \left[(Y - \bar{m})^3 \right]$ <p>What is actually examined is $a_3 = \frac{m_3}{(m_2)^{3/2}}$</p> <p>because it is unitless</p> <p>$a_3 = 0$ indicates symmetry</p> <p>$a_3 < 0$ indicates lefthand skew (tail to left)</p> <p>$a_3 > 0$ indicates right hand skew (tail to right)</p>	<p>When <i>yvariable</i> is distributed normal:</p> <p>Skewness = 0</p> <p>Look for skewness between -2 and +2, roughly.</p>
<p>Kurtosis – flatness versus peakedness of the curve</p> <p>Standardization of the 4th sample moment about the mean</p> $m_2 = E \left[(Y - \bar{m})^2 \right]$ $m_4 = E \left[(Y - \bar{m})^4 \right]$ <p>Pearson kurtosis is $a_4 = \frac{m_4}{(m_2)^2}$</p> <p>$a_4 = 3$ when distribution is normal</p> <p>$a_4 < 3$ is “leptokurtic” (too little in the tails)</p> <p>$a_4 > 3$ is “platykurtic” (too much in the tails)</p>	<p>When <i>yvariable</i> is distributed normal:</p> <p>Kurtosis = 3</p>

Hypothesis Tests of Normality and what to watch out for:

Test Statistic	What to watch out for:
<p>1. <u>Shapiro Wilk (W)</u></p> <p>W is a measure of the correlation between the values in the sample and their associated normal scores (for review of Normal Scores, see BIOSTATS 540 Unit 7 - Normal Distribution)</p> <p>$W = 1$ under normality</p>	<p><u>Null Hypothesis H_0:</u> <i>yvariable</i> is distributed normal: <u>Alternative Hypothesis H_A:</u> Not.</p> <p>Violation of normality is reflected in $W < 1$ small p-value</p>
<p>2. <u>Kolmogorov-Smirnov (D). See also Lilliefors (K-S)</u></p> <p>This is a goodness of fit test that compares the distribution of the residuals to that of a reference normal distribution using a chi square test.</p> <p>Lilliefors utilizes a correction</p>	<p>Violation of normality is reflected in</p> <p>$D > 0$</p> <p>$K-S > 0$</p> <p>small p-value</p>

Guidelines

In practice, the assessment of normality is made after assessment of other model assumption violations.
The linear model is often more robust to violations of the assumption of normality.
The cure, is often worse than the problem. (e.g. – transformation of the outcome variable)

Consider doing a scatterplot of the residuals. Look for

- ◆ Bell shaped pattern
- ◆ Center at zero
- ◆ No gross outliers

c. Cook-Weisberg Test of Heteroscedasticity

Recall what we are assuming with respect to homogeneity of variance:

- **In Simple Linear Regression:**
At each level “x” of the predictor variable X, the outcomes Y are modeled as distributed normal with mean $\mu_{Y|x} = \beta_0 + \beta_1 x$ and constant variance $\sigma_{Y|x}^2$

Evidence of a violation of homogeneity (this is heteroscedasticity) is seen when

- There is increasing or decreasing variation in the residuals with fitted \hat{Y}
- There is increasing or decreasing variation in the residuals with predictor X

Some **graphical assessments** of homogeneity of variance and what to watch out for:

Method	What to watch out for:
1. Plot Residuals or standardized residuals or studentized residuals on the vertical — versus — Predicted outcomes \hat{Y} on the horizontal	Look for even band at zero
2. Plot Residuals or standardized residuals or studentized residuals on the vertical — versus — Predictor values X	Look for even band at zero

Hypothesis Test of homogeneity of variance is Cook-Weisberg

Cook-Weisberg Test	What to watch out for:
This test is based on a model of the variance as a function of the fitted values (or the predictor X). Specifically, it is a chi square test of whether the squared standardized residuals are linearly related to the fitted values (or the predictor X).	Evidence of violation of homogeneity of variance is reflected in Large test statistic > 0 small p-value

d. The Method of Fractional Polynomials

This method is beyond the scope of this course. However, it's helpful to understand the ideas.

Goal: The goal is to select a “good” functional form that relates Y to X from a collection of candidate models. Candidates are lower polynomials and members of the Box-Tidwell family.

Fractional Polynomials: Instead of $Y = \beta_0 + \beta_1 X$, we consider the following:

Instead of fitting a
simple linear relationship of the form
 $b_1 X$

We consider fitting a
fractional polynomial relationship of the form
 $b_1 X^{p_1} + b_2 X^{p_2} + b_3 X^{p_3} + \dots + b_m X^{p_m}$

where

m = number of powers (“degree”)

p₁, p₂, p₃, ..., p_m are choices from a special set of 8 candidate powers = { -2, -1, -0.5, 0, 0.5, 1, 2, 3 }

And where, when powers repeat

E.g. - when **p₂ = p₁** we consider $b_1 X^{p_1} + b_2 X^{p_1} \ln(X)$.

Example: Suppose **m=1** with **p₁ = 1**. This yields

$$Y = b_0 + b_1 X$$

Example: Next, suppose **m=2** with **p₁ = 0.5** and **p₂ = 0.5**. Because **p₂ = p₁** this yields

$$Y = b_0 + b_1 \sqrt{X} + b_2 \sqrt{X} \ln(X)$$

The Method of Fractional Polynomials - Continued

Guidelines

Competing models are assessed using a chi square statistic that compares the likelihoods of the data under each of the two models using what is called a “deviance” statistic. (*Stay tuned.* We will learn about the “deviance” statistic in Unit 7, Logistic Regression.)

The search for a "good" model by the method of fractional polynomials begins with an examination of all the models for which $m=1$. We choose the one model in this class that has the smallest deviance (think "left over variability that is not yet explained").

- ◆ We compare the best $m=1$ model to the specific model for which $m=1$ and $p_1=1$ because the latter is the simple linear model.
- ◆ Thus, we are asking whether it is really necessary to abandon the simple linear model.

Next, we compare the best $m=1$ model to the best $m=2$ model. And so on ...

- ◆ There's always a trade-off:

- 1) A smaller model has a lower goodness-of-fit 😞 but more generalizability 😊
- 2) A larger model has a higher goodness-of-fit 😊 but less generalizability 😞

- ◆ Our goal is to choose the smallest model for which the goodness-of-fit is acceptable.

e. Ramsey Test for Omitted Variables

A fitted model that fails to include an important explanatory variable is problematic.

- ◆ We are missing part of the story!
- ◆ Possibly, we have (incorrect) biased associations due to uncontrolled confounding.
- ◆ We may have violated some model assumptions.

Method of the Ramsey Test

- ◆ H_0 : Predicted values from the fitted model are unrelated to powers of the fitted model, after adjustment for the predictor variables in the model.

$$\text{corr}(\hat{Y}, \hat{Y}^p) = 0$$

- ◆ For example, we might fit the model $\hat{Y} = b_0 + b_1 \hat{Y} + b_2 \hat{Y}^2 + b_3 X + \text{error}$ and test the significance of \hat{b}_1 and \hat{b}_2 .
- ◆ The test statistic is an F statistic.

Guidelines

A large F statistic value is consistent with failure to include one or more explanatory variables.

Suggestion. Accompany this test with a visualization. Do also a scatterplot of the squared standardized residuals versus the leverage values. Omission of an important explanatory variables is suggested by

- ◆ Extreme values
- ◆ Any systematic pattern

f. Residuals, Leverage, and Cook's Distance

Residuals - There are multiple measures of “residual”.

Ordinary residual $e = (Y - \hat{Y})$	Standardized residual $e^* = \frac{e}{\sqrt{ms(residual)}} = \frac{e}{\sqrt{\hat{\sigma}_{Y x}^2}}$
Studentized residual $e^* = \frac{e}{\sqrt{ms(residual)}\sqrt{1-h}} = \frac{e}{\sqrt{\hat{\sigma}_{Y x}^2}\sqrt{1-h}}$	Jackknife residual, also called Studentized deleted residual $e^* = \frac{e}{\sqrt{ms(residual)_{-i}}\sqrt{1-h}} = \frac{e}{\sqrt{\hat{\sigma}_{Y \hat{x}}^2}\sqrt{1-h}}$

Which one or ones should we use?

- **Standardized** residuals can be (roughly) interpreted as z-scores.
- **Studentized** residuals can be (roughly) interpreted as t-scores from a Student's t (df=n-p-1) when regression assumptions hold.
- **Jackknife** residuals can be (roughly) interpreted as t-scores from a Student's t (df=n-p-2) when regression assumptions hold. These also have the advantage of correcting the magnitude of the $\sqrt{MS(residual)}$ when it is otherwise too big because of the effects of influential points.

Leverage, h_i :

Leverage is the distance of a predictor value $X=x$ from the center of the values of the predictor value $X = \bar{x}$. This distance is denoted h_i .

For simple linear regression,
$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

For simple linear regression, a “large” leverage value is $h_i \geq \frac{4}{n}$

Cook's Distance, d

Recall. Neither a large residual alone nor a high leverage alone is a guaranteed that an individual data point is influential. To see this, see again the pictures on pp 43-44.

Cook's distance to the rescue. Cook's distance provides a measure of the influence of an individual data point on the fitted model and is a function of the values of both the residual and leverage:

Cook's Distance
Change in estimated regression coefficient value, expressed in standard error units.

1) For simple linear regression
$$d = \frac{e^2 h}{2s^2 (1 - h)^2}$$

2) For multivariable linear regression models
$$d_i = \frac{(\hat{b}_{-i} - \hat{b})' (X'X) (\hat{b}_{-i} - \hat{b})}{p \hat{s}_{Y|X}^2}$$
 where

i indexes the individual for which measure of influence is sought
 \hat{b} = vector of estimated regression coefficients using the entire sample
 \hat{b}_{-i} = vector of estimated regression coefficients with omission of the i^{th} data point
 X = matrix of values of the predictor variables
 p = rank (X) = number of predictors + 1

How big should a Cook's Distance be to conclude the data point is influential?

Simple Linear Regression:

Cook's distance $d \geq 1$.

Multiple Linear Regression:

Cook's distance $\geq 2(p+1)/n$ where

n = sample size; and
 p = # predictors.

g. Example

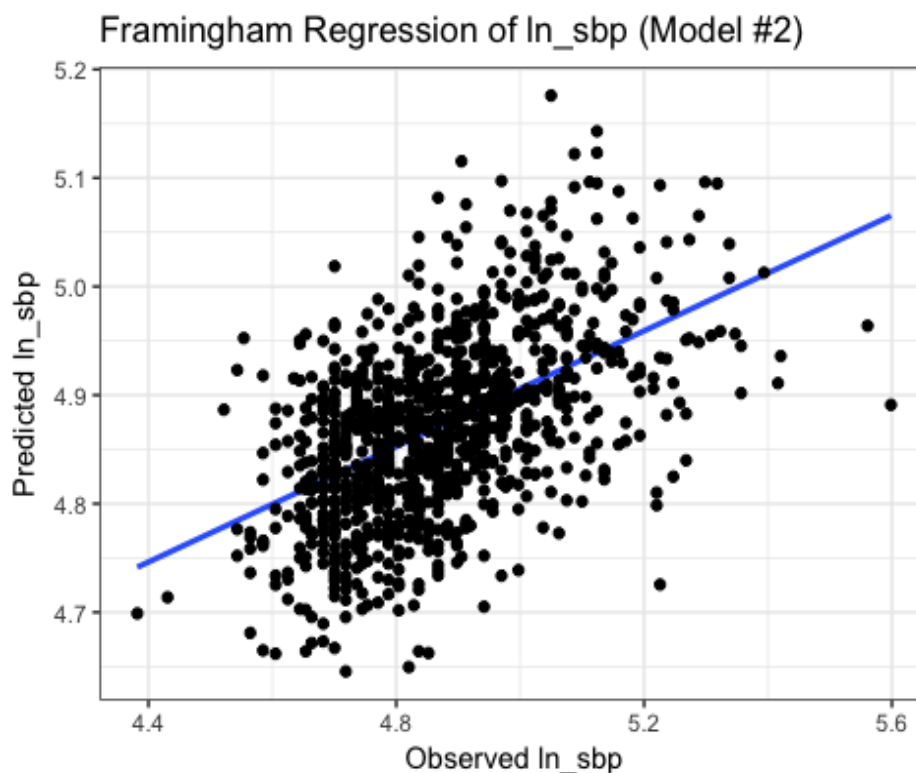
Framingham Study – model #2

Plot of Observed v Predicted. Look for: Points along a straight line (“all is well”)

```
library(ggplot2)
library(Hmisc)

m_best <- lm(data=complete, ln_sbp ~ ln_bmi + ln_scl + age + female + ageXfemale) # Fit model to complete data
complete$yhat <- predict(m_best) # Add predicted values to dataset
Hmisc::label(complete$yhat) <- "Predicted ln(sbp)"

ggplot(data=complete, aes(x=ln_sbp,y=yhat)) +
  geom_smooth(method=lm, se=FALSE) + # TIP - plot line first
  geom_point() + # Then plot your points on top
  xlab("Observed ln_sbp") +
  ylab("Predicted ln_sbp") +
  ggtitle("Framingham Regression of ln_sbp (Model #2)") +
  theme_bw()
```

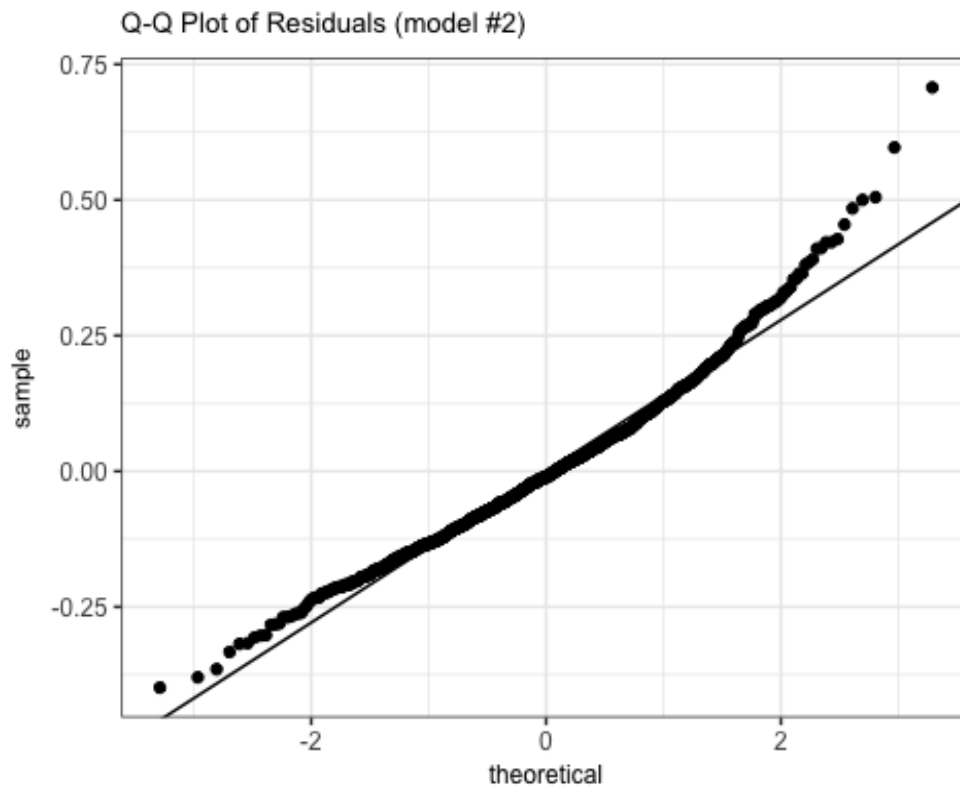


Interpretation – Not bad! Ideally, the scatter lies on the line defined by 45 degrees. We expect some widening of the confidence intervals at the ends of the range but not too much. What we see here is reasonable.

Normality of Residuals – QQ Plot and Shapiro Wilk Test. Null: Normality (“all is well”)

```
library(ggplot2)
complete$residuals <- residuals(m_best) # Add the residuals to the dataset

ggplot(data=complete, aes(sample=residuals)) +
  stat_qq() +
  geom_abline(intercept=mean(complete$residuals),
              slope = sd(complete$residuals)) +
  ggtitle("Q-Q Plot of Residuals (model #2)") +
  theme_bw() +
  theme(axis.text = element_text(size = 9),
        axis.title = element_text(size = 9),
        plot.title = element_text(size = 10))
```



```
options(scipen=1000)
shapiro.test(complete$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  complete$residuals
## W = 0.9775, p-value = 0.00000000028
```

Interpretation – Here too, we hope to see a scatter on the 45 degree line. Not bad!

Ramsay Test of Omitted Variables. Null: No omissions (“all is well”)

```
library(lmtest)
lmtest::resettest(m_best, power=2, type="regressor")
## RESET test
##
## data: m_best
## RESET = 0.42467, df1 = 5, df2 = 983, p-value = 0.8317
```

Interpretation - Ramsey test is NOT significant ($p=.83$) suggesting we're okay!

Assessment of Multicollinearity (“all is well” if VIF < 10)

```
library(car)
car::vif(m_best)

##      ln_bmi      ln_scl      age      female ageXfemale
##  1.115511  1.175531  2.378150  32.394888  34.116761
```

Interpretation - female and ageXfemale appear to be collinear suggesting some concern about the extent to which there is adequacy of range of age in the 2 genders.

Cook's Distances (flag observations for which Cook distance > $4/(n-p-1)$. Other definitions possible.

```
library(Hmisc)
library(ggplot2)
complete$ID <- as.numeric(row.names(complete)) # create study id using row.names( ) and as.numeric( )
Hmisc::label(complete$ID) <- "Observation Number"

complete$cooks <- cooks.distance(m_best) # Add cooks distances to the dataset

cutoff <- 4/((nrow(complete)-length(m_best$coefficients)-2)) # Solve for cutoff as equal to = 4 / (n-p-1).

ggplot(data=complete, aes(x=ID, y=cooks)) +
  geom_bar(stat="identity", position="identity") +
  xlab("Observation Number") +
  ylab("Cooks Distance") +
  geom_hline(yintercept=cutoff) +
  geom_text(aes(label=ifelse((cooks>cutoff), ID, "")), ID, ""), vjust=-0.2, hjust=0.5) +
  ggtitle("Cooks Distances > 4 / (n-p-1)") +
  theme_bw()
```

