



# Case Studies using the R Programming Language

Katherine Bennett, Miguel de los Reyes, Raymond Gao,

Sophia Hurr, Hannah Gahagan, Nikhil Miland, Mridu Nanda,

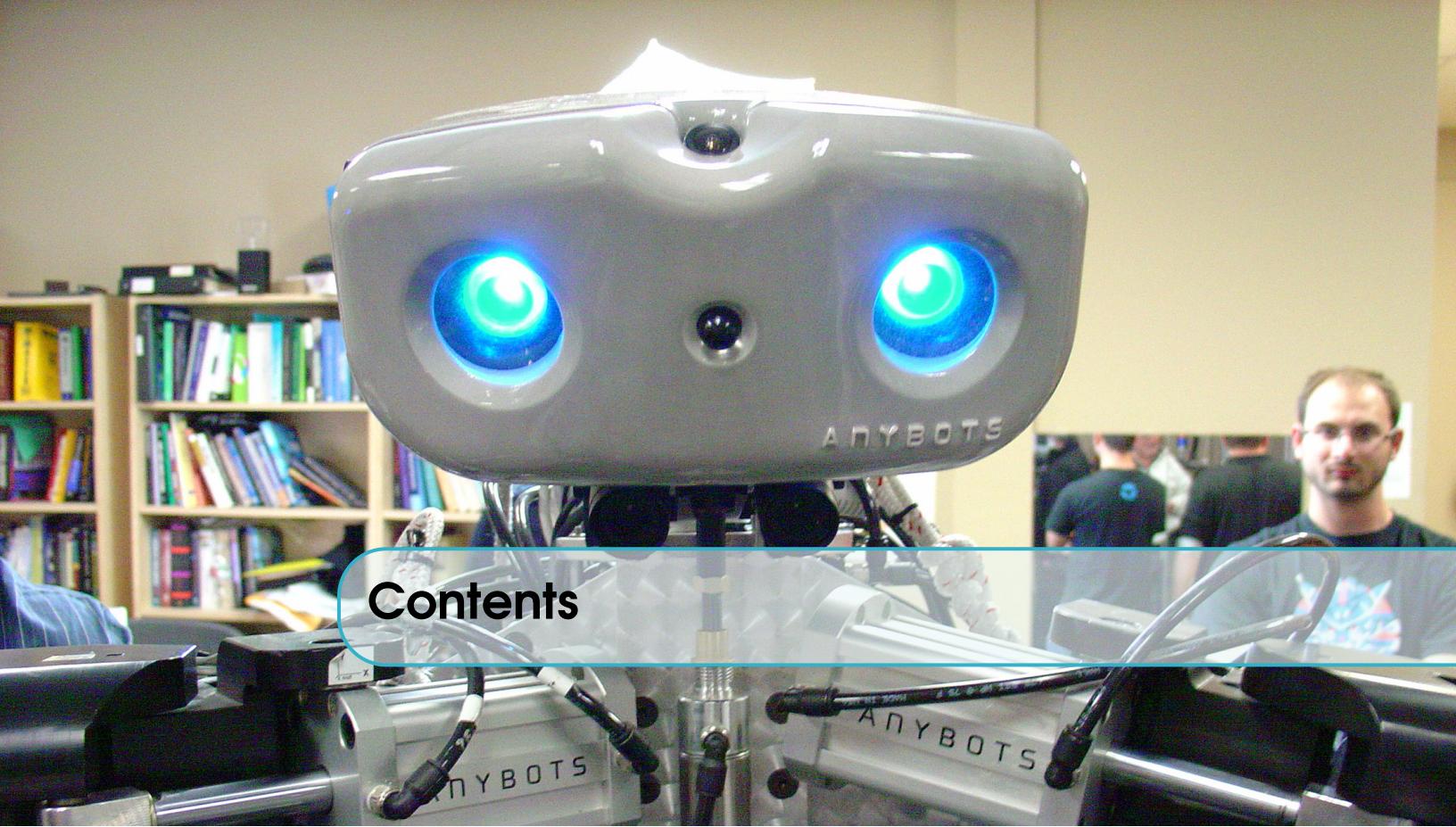
Christa Parrish, Ishaan Rao, Grayson York

RESEARCH IN COMPUTATIONAL SCIENCE

NORTH CAROLINA SCHOOL OF SCIENCE AND MATHEMATICS, DURHAM, NC

This case studies manual was prepared under the direction of Mr. Robert Gotwals,  
faculty mentor, The North Carolina School of Science and Mathematics, Durham  
NC.

*First release, November 2016*







# Introduction

Gotwals will write an Intro.





# 1. ANOVA

## 1.1 Introductory Reading

The most common type of linear models are analysis of variance (ANOVA) and linear regression models. R is built to handle linear models, and as such it is easy to work with ANOVA. ANOVA is a statistical method used to compare two or more means. These values can be used to determine whether a significant correlation exists between variables. We will be using one-way ANOVA, which compares the means between the groups of interest, and determines whether those means are significantly different from each other. If one-way ANOVA returns a significant result, there are at least two group means that are significantly different from each other. It is important to note that ANOVA is an omnibus test statistic, meaning that although ANOVA tells the user there is a significant difference in means, it is unable to show which means in particular differ. ANOVA can only be used with certain types of data configurations. To perform ANOVA, data must have a continuous response variable and at least one categorical factor with two or more levels, in lay terms, ANOVA can only be used with numeric values that can be ordered sequentially, with a certain number of possible responses. (i.e. data comparing object's weights, and each new weight is a level) ANOVA is easier to use if data is from approximately normally distributed populations with equal variances between factor levels, however, as R is built to work with statistics, ANOVA procedures will generally work without incident unless one or more of the distributions or variances are highly skewed. A basic understanding of statistics how to use R is recommended before performing ANOVA. The ANOVA user will create and order factors, make box plots, and combine and stack data. (**CRANR**) The following functions will be useful in using ANOVA.

1. aov

2. `as.factor`
3. `ls.str`
4. `data.frame`
5. `stack`
6. `TukeyHSD`
7. `levels`

## 1.2 Objectives

In this assignment, we will use ANOVA to analyze data from ELISA HIV Optical Density Readings. ANOVA identifies the causes of variation, and sorts out the corresponding components of variation with associated degrees of freedom. In this case, we are interested in seeing how HIV Optical Density Readings are related to their lot. The type of ANOVA we are using is one-way between groups, as we are comparing one grouping (the optical density readings) to define the groups (lots).

By the end of this lesson, the reader should be able to:

1. Understand the logic behind one-way analysis of variance.
2. Perform one-way analysis of variance in R for any data.
3. Appropriately interpret results of analysis of variance tests.

### 1.2.1 ANOVA

We will look at variances in data using R's ANOVA functions.

### 1.2.2 Visualization

We will also use our fitted models and our ANOVA data to create box plots and line plots. Much of the information gleaned from ANOVA is presented via numeric such as the sum of square, degrees of freedom, and the mean. ANOVA presents the null hypothesis that there is no difference in means of the treatments, and once this hypothesis is proven incorrect, the question arises of how the treatments differ. The post-hoc test also allows the user to find the differences in means, and specifically categorizes the lower and upper means of the data.

Figure 1.1 shows a box plot of the data before running ANOVA or TukeyHSD. Constructing a box plot before analyzing the data may prove helpful in deciding what kind of analysis would be preferable.

## 1.3 Building the Model

Looking at the data set of interest for ANOVA, note that that data may be in numeric form. Use the function `is.numeric` to review the data. If this function prints TRUE, use the function `as.factor` to change data from numeric to factor. In order to make it

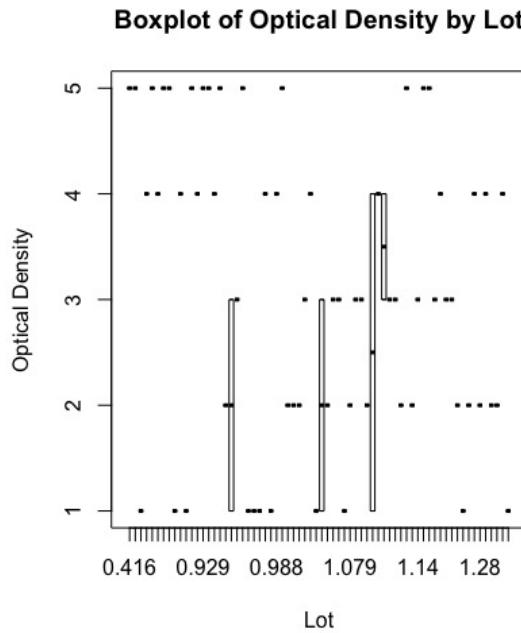


Figure 1.1: Boxplot

easier to manipulate data later on, it is recommended that the factored data is renamed.

---

```
1 lot = as.factor(elisa$Lot)
2 summary(lot)
3 levels(as.factor(elisaOptical))
4 optical = (as.factor(elisa$Optical))
5 levels(as.factor(elisa$Run))
6 run = (as.factor(elisa$Run))
```

---

Now your factor data type is in non-numerical variables. Each different variable of the factor is called a level. For example, our factor type data for the lots of various ELISA HIV Optical Density Readings has five levels, 1, 2, 3, 4 and 5. Moving on, we can make a box plot of our data, so that we can visually compare data. Making a summary of this data allows us to look at the residuals and standard error of the plot. As mentioned before, this type of comparison will work better with data that has a relatively normal distribution. Factors can also be created within factors. We used our factor Optical to categorize and label by the density. For each level within the main factor, create labels. These labels will designate unordered factors. To make this data easier to read, unordered factors should be ordered. Choose ways to order the data that is relative to the averages in the data. For example, some levels were Low density(below 1) and some were High density(above 1.5) so we ordered the data starting at 0, then continuing with 1, 1.5, 2. Choose labels that reflect the data, noting that the first label will relate to the lower numbers. Next,

classify this newly ordered factor, and now label the data in lay terms.

---

```
1 elisa$Optical.type<-ordered(cut(elisa$Optical, c(0,1,1.5,2), labels = c("Low
  ", "Medium", "High")))
2 class(elisa$Optical.type)
3 elisa$Optical.type
```

---

Note, that if you wish to remove one group of subjects (i.e. Lot 1) R will keep this removed group as a level, possibly skewing your data. Use the function drop.levels to remove a selected group as a level. Rename this reconfigured data as to not overwrite your previous data. In order to use ANOVA, the data must be in a specific format. To properly configure data, combine the data from the two factors intended to be compared with the function data.frame. Next, stack these combined groups. Finally use this stacked data and the function aov to perform ANOVA. Make a summary of these results.

---

```
1 finalelisa <- droplevels(newelisa)
2 summary(finalelisa$Lot)
3 Combined_Groups <- data.frame(cbind(lot,optical))
4 Combined_Groups
5 Stacked_Groups <- stack(Combined_Groups)
6 Stacked_Groups
7 Anova_Results <- aov(values ~ ind, data = Stacked_Groups)
```

---

This method can be used for a simple one-way ANOVA. However some data can compare multiple sets of data against each other. To do this, we can use a post-hoc test. Post-hoc tests compare outcome measurements between multiple groups. With post-hoc analysis the reader can examine differences between pairs of groups after global analysis. Use the function TukeyHSD to perform a pairwise post-hoc analysis on the ANOVA results.

---

```
1 summary(Anova_Results)
2 TukeyHSD(Anova_Results)
3 summary(TukeyHSD(Anova_Results))
```

---

### 1.3.1 Programming Hints

It is recommended to look at data before analysis in order to better understand the levels and attributes of data. In addition, we recommend the use of `summary()` and the console to view the attributes of objects.

## 1.4 Deliverable

Using ANOVA, students should be able to take data with sequentially ordered number data, make a box plot, factor the data as unordered and ordered factors, configure the data such that the function aov can be performed, and perform

TukeyHSD. Using the box plot, students should be able to identify residuals, coefficients of the linear regression, and residual standard error. Using aov, students should be able to identify degrees of freedom, the sum of the squares, the mean of the squares, and the F ratio (the ratio of two mean square values) for their data.

## 1.5 Teaching Code

Begin ANOVA by formatting the data in a .csv file, then uploading it to RStudio. This script is set up to analyze publicly available data about HIV Optical density readings, but can be adapted for any properly set up data. As the reader works his or her way through the script, be sure to liberally comment the purpose of each command.

---

```
1 # Your name here
2 # Date
3 # ANOVA - ELISA HIV Optical density readings data
4
5 # Clean up and read in data
6 rm(list=ls())
7 setwd("C:/Your/Directory")
8 elisa = read.csv("ELISAHIV.csv")
9
10 # Create levels for lot, optical, and run
11 # Make sure to view/check your data
12 levels(as.factor(elisa$Lot))
13 lot=as.factor(elisa$Lot)
14 # etc.
15
16
17 # Create a boxplot showing optical density by lot
18 # Make sure to label your axes!
19 boxplot(..., ...)
20
21
22 # Perform a linear regression of optical density by lot
23 summary(lm(..., data=elisa))
24
25 # Create an ordered factor using Optical.type and verify its type
26 elisa$Optical.type<-ordered(...)
27 class(elisa$Optical.type)
28
29 # Remove 1 as a possible level
30 # Note: even if you eliminate a group of subjects (ex. 1),
31 # because it's a factor, R keeps 1 as a possible level for lot
32
33 # Use droplevels() to remove 1
34
35 # Create stacked results to run aov()
36
37 # Run TukeyHSD on the ANOVA results
```

---

## 1.6 Example Student Code

---

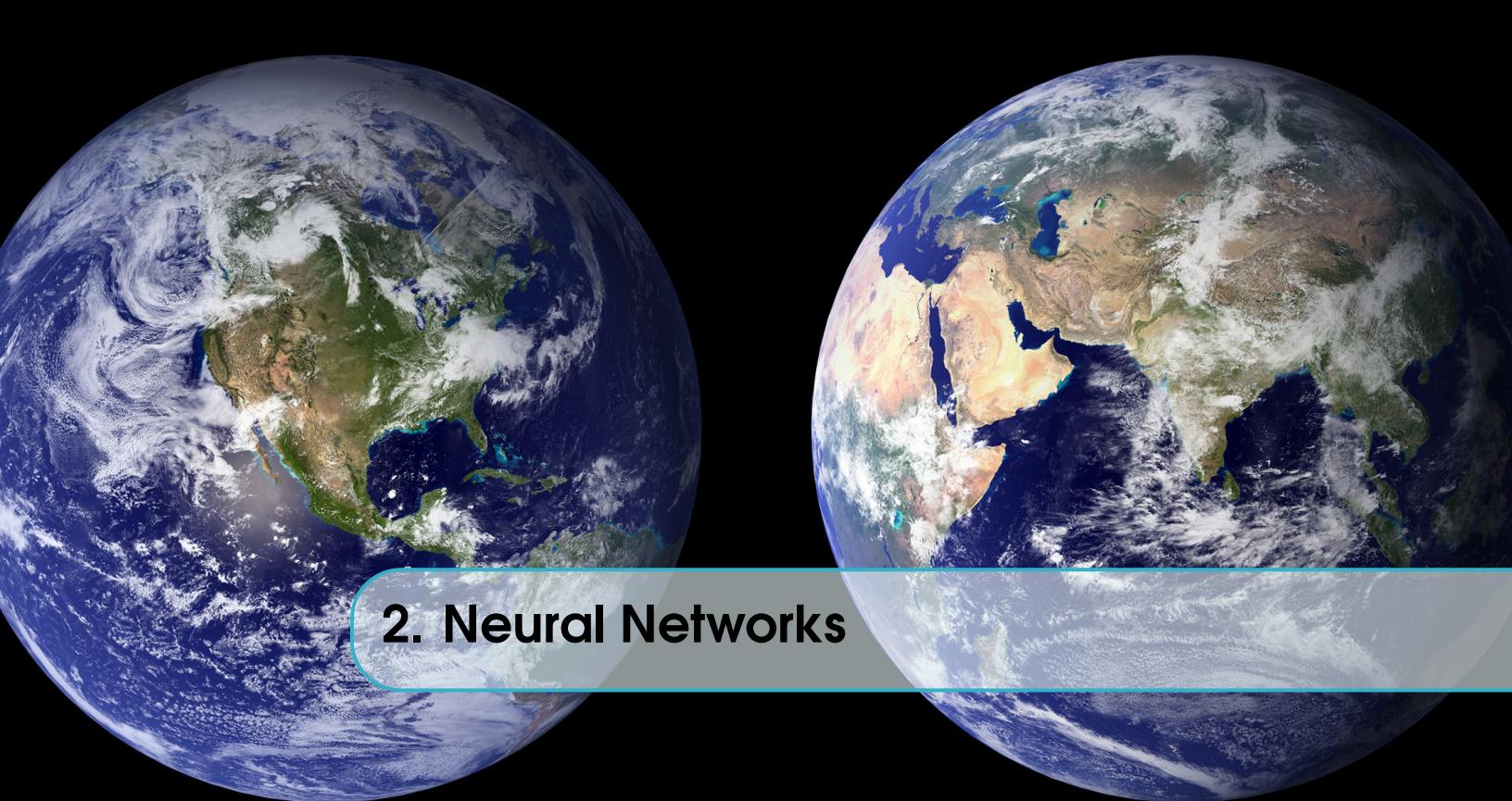
```
1 # KEY
2 # ANOVA - ELISA HIV Optical density readings data
3
4 # Clean up and read in data
5 rm(list=ls())
6 setwd("C:/Example/Code")
7 elisa = read.csv("ELISAHIV.csv")
8
9 # Look at variables and create levels
10 ls.str(elisa)
11 levels(as.factor(elisa$Lot))
12 lot=as.factor(elisa$Lot)
13 summary(lot)
14 levels(as.factor(elisa$Optical))
15 optical = (as.factor(elisa$Optical))
16 levels(as.factor(elisa$Run))
17 run = (as.factor(elisa$Run))
18
19 # Create a boxplot
20 boxplot(elisa$Lot~elisa$Optical, xlab="Lot", ylab="Optical Density",
21 main="Boxplot of Optical Density by Lot")
22 summary(lm(elisa$Optical~elisa$Lot, data=elisa))
23
24 # Perform regression
25 summary(optical)
26 elisa$Optical.type<-ordered(cut(elisa$Optical, c(0,1,1.5,2),
27 labels=c("Low","Medium","High")))
28 class(elisa$Optical.type)
29 elisa$Optical.type
30
31 # Remove 1 as a level
32 # We've verified that optical is an ordered factor
33 # Note: even if you eliminate a group of subjects (ex. 1),
34 # because it's a factor, R keeps 1 as a possible level for lot
35 newelisa <- elisa[1:2,]
36 summary(newelisa$Lot)
37 # We remove 1 as a possible level using droplevels()
38 finalelisa <-droplevels(newelisa)
39
40 # We create stacked results to run aov()
41 summary(finalelisa$Lot)
42 Combined_Groups <- data.frame(cbind(lot,optical))
43 Combined_Groups
44 Stacked_Groups <- stack(Combined_Groups)
45 Stacked_Groups
46 Anova_Results <- aov(values ~ ind, data = Stacked_Groups)
47 summary(Anova_Results)
48
49 # We also run TukeyHSD on the ANOVA results
50 TukeyHSD(Anova_Results)
51 summary(TukeyHSD(Anova_Results))
```

---

## 1.7 Further Readings

1. Seefeld, Kim and Linder, Ernst. *Statistics Using R with Biological Examples*, University of New Hampshire, Durham, NH Department of Mathematics and Statistics(2007)
2. Julian J. Faraway. *Practical Regression and Anova using R*, University of Michigan, (2002)





## 2. Neural Networks

### 2.1 Introduction to Neural Networks

In the field of computational research, the use of computers is generally restricted to discreet, exact commands that are meant to be executed in a certain manner, which provide the computational power required for complex calculations. In contrast, neural networks are an example of machine learning (ML), wherein the machine is capable of learning behavior or interpreting patterns from a given data set. The field of predictive analysis uses this paradigm of computational interpretation. For example, it is extremely challenging to develop a system that is purely programmed to recognize handwriting. Yet, computers at the United States Postal Service are able to decipher human words and reroute letters and packages without the need for humans. These computers contain a neural network that has been trained to recognize handwriting and predict characters as they appear on mail. ([nielsen](#))

Neural networks were designed using the human brain as an inspiration. The neural network is composed of perceptrons, which are similar to the smallest functional units of the brain: the neuron. Neurons act as the wiring of the brain by passing along signals that they receive to other neurons. When a neuron communicates with another neuron, it uses neurotransmitters to either increase or decrease the probability of the next neuron firing. Thus, a neuron can be excitatory (increases the chance of the next neuron firing) or inhibitory (decreases the chance of the next neuron firing). Similarly, neural networks utilize perceptrons that can receive multiple inputs and produce a single output. The inputs are binary values (0 or 1), and affect the perceptron differently depending on their excitatory or inhibitory ability. The amount of excitation or inhibition is captured by weights, or values that the inputs are multiplied by. Finally, the aggregate sum of these inputs and their

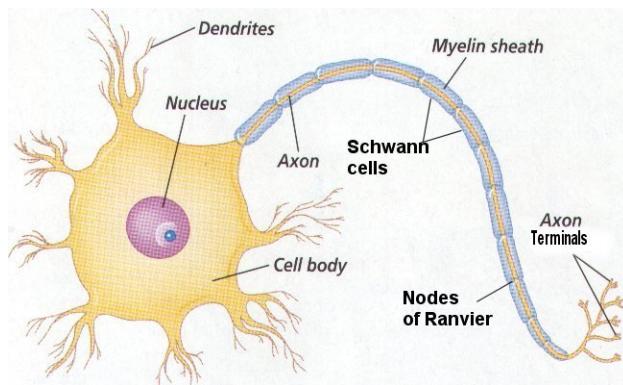


Figure 2.1: The Parts of an Anatomical Neuron

weights is used to calculate if the perceptron in question will fire an output to the next layer of perceptrons.

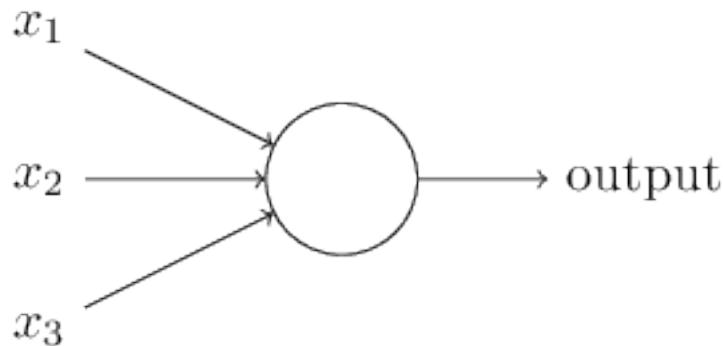


Figure 2.2: A neuron in a neural network, with binary inputs  $x_1$ ,  $x_2$ , and  $x_3$ .

Neurons that only utilize binary inputs, however, are not efficient in learning. For example, neurons in the brain release different amount of neurotransmitters to indicate different levels of excitation. Similarly, most modern neural networks generate an output that is between 0 and 1, inclusive. This type of function is known as the sigmoid function, and the resulting neuron is known as the sigmoid neuron.

### 2.1.1 Neural Networks with Sigmoid Neurons

With practice, humans can arrange sigmoid neurons to generate a functional unit that can take inputs and produce a certain computational output. Functionally, sigmoid neurons act as NAND gates, which are universal to computational procedures. However, the true value of the neural network lies in data and input values that do not have any discernible pattern. In such cases, neural networks are capable of learning patterns by themselves. Therefore, neural networks can be trained with an initial data set, after which they can predict the output values of new information. Although a single sigmoid neurons can be trained, it is not adequate to use such a unit to capture the entirety of complex patterns found in data sets. Rather, a

network of sigmoid neurons is created to assess input values and generate proper outputs. Neural networks are most often arranged in layers. The first layer is known as the input layer, whilst the last layer is known as the output layer. All the neural layers in the middle are known as hidden layers. The sample neural network shown

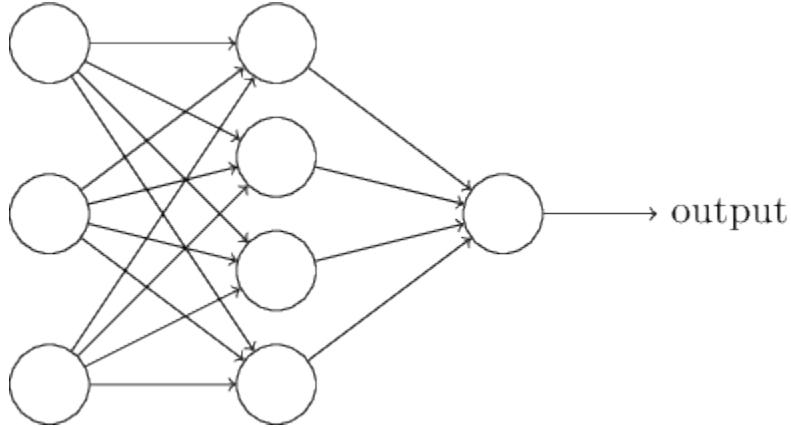


Figure 2.3: A neural network with three layers: an input layer, a hidden layer, and an output layer.

in figure 2.3 takes four inputs in the first layer and produces a singular output. Since all the components involved are sigmoid neurons, the value of the output of the neural network can take any value between 0 and 1.

### 2.1.2 Training in a Neural Network

In general, every sigmoid neuron has a certain set of inputs  $X = \{x_1, x_2, \dots, x_n\}$ . Each input is associated with a weight, or the amount of excitation or inhibition each input causes, which are expressed as the weight set  $W = \{w_1, w_2, \dots, w_n\}$ . The value of the previous inputs is also influenced by a neuron-specific bias  $b$ , which acts as a value of how reactive the neuron is to the inputs. Finally, the sigmoid neuron produces an output  $O$  between 0 and 1 depending on the influence of previous inputs as follows:

$$\begin{aligned}\sigma(x) &= \frac{1}{1 + e^{-x}} \\ O(W, X) &= \sigma \left( \sum_i^n w_i x_i + b \right) \\ O(W, X) &= \sigma (W \cdot X + b) \\ O(W, X) &= \frac{1}{1 + e^{-W \cdot X - b}}\end{aligned}$$

Therefore, by computing the dot product of set  $W$  and set  $X$ , offsetting the dot product by the neural bias, and then passing that value through the sigmoid function, we generate an output between 0 and 1. If the sigmoid neuron in question

is part of a hidden layer, it will pass on its output to other neurons in the next layer.

To train a network, we begin with a neural network that is set up with well-chosen (but random) weights and begin sending input values through it. The output is then compared to the value we were expecting, and computational corrections are made in increments to the weights and biases in the neural network using gradient descent algorithms. Over time, the neural network becomes more accurate and has the ability to predict future output values from input values.

### 2.1.3 Measuring Predictive Error

To train our neural network computationally, we need to develop a cost function, that quantifies the amount of training our neural network has undergone. That is, this function has to return a value that quantifies the difference between prediction and reality. During training, every predicted value is compared to the actual value of the data through subtraction and squared. Thus, for a prediction  $y$  and actual value  $\hat{y}$ ,

$$C(y, \hat{y}) = (\hat{y} - y)^2$$

The cost function is computed for every data point used to train the neural network. The overall error of the neural network is simply the average of all the cost function values for every data point that was used to train the network. However, since we are squaring the cost function, we divide the average by 2 to decrease the severity of predictive error. This error function is most commonly known as the Mean-Squared Error function, and is predominantly used to train neural networks. For a set of predictions  $Y = \{y_1, y_2, \dots, y_n\}$  and a set of actual values  $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n\}$ :

$$\begin{aligned} J(Y, \hat{Y}) &= \sum_{i=0}^n \frac{C(y_i, \hat{y}_i)}{2n} \\ J(Y, \hat{Y}) &= \sum_{i=0}^n \frac{(\hat{y}_i - y_i)^2}{2n} \end{aligned}$$

## 2.2 Installation

The neuralnet package can be found at <https://cran.r-project.org/web/packages/neuralnet/index.html>. Download the package that works for your operating system and CPU architecture. The following code can be used to install the package from the R command line. This example illustrates the installation of a windows system binary file. However, all other installations will also be the same, with FILE\_PATH leading to the package downloads. Alternatively, a script file can be generated to install the neuralnet package.

Note that safari automatically decompresses .tgz files to tar files, and this will cause errors when running the above code. The installation may throw a warning reminding the user that the package may be out-of-date when compared to the R version being used.

Figure 2.4: Code for Windows

---

```
1 install.packages("FILE_PATH.zip", repos=NULL)
2 package 'neuralnet' successfully unpacked with MD5 sums checked
3 library(neuralnet)
```

---

Figure 2.5: Code for Mac / OSX

---

```
1 install.packages("FILE_PATH.tgz")
2 package 'neuralnet' successfully unpacked with MD5 sums checked
3 library(neuralnet)
```

---

## 2.3 Objectives of Case Study

There are multiple objectives that the student must accomplish over the course of this neural networks case study.

1. Understand neural networks, their functional units, and their overall purpose in Machine Learning. It is important to understand how neural networks work, beginning from an understanding of sigmoid neurons and culminating in an overall appreciation for the function of neural networks in predictive analysis.
2. Have the ability to prepare data for use in neural networks in R. This objective includes the ability to identify and remove rows with empty data, normalizing the data to allow for learning over a smaller number of iterations, and fitting data curves to create the predictive function.
3. Be able to implement the neuralnet package in R. An understanding of the functionality that the neuralnet package provides, including creating the neural network, training it with a data set, and using it for predictive analysis, is imperative to this objective.
4. Perform cross-validation on the neural network to confirm understand its learning progress. This involves using statistical measures to see how the neural network is functioning during training and prediction.

## 2.4 Case Study - Boston Housing Database

The Boston Housing Database is a database that contains the housing data from the boston area. It contains data on housing sales, including the area of the house sold, the location of the house, and other price-determining factors.

### 2.4.1 Loading and Separating the Data

Create a R script file and enter the following code. This code will set up the environment for the neural network. We load the data into our scripting environment, modify it to fit the neuralnet package, and split it into a training and testing dataset.

---

```

1 # Name:      neural.R
2 # Author:    First Last
3 # Date:     Month Date Year
4
5 # Load necessary libraries into R
6 library(neuralnet) # Contains the neuralnet package
7 library(MASS)       # Contains the database we will use
8
9 # Clear the current environmental variables
10 rm(list=ls())
11
12 # We set the seed for the random generator
13 # We use random data values for training and for testing
14 # Use this seed to replicate the results in this chapter
15 set.seed(500)
16
17 # Import the database into R from the MASS library
18 data <- Boston
19
20 # Define the number of neurons in each hidden layer
21 # The input layer depends on the number of columns in the data frame
22 # The output layer depends on variables we ask the network to predict
23 hiddenLayers = c(5, 3)

```

---

Unfortunately, the neuralnet package cannot handle NULL values in the data. Therefore, we must remove such extraneous data from the dataset.

---

```

1 # Unlike most R libraries, neuralnet is unable to parse NA values
2 # We check for NA values in the data set and print their amount
3 # If we have rows with missing data, we will have to remove them
4 apply(data, 2, function(x) sum(is.na(x)))

```

---

We will use 75% of the data to train the network and 25% to test the efficacy of the neural network. We split the dataset using the following code.

---

```

1 # Generate a list of indices that are randomly sampled
2 # The indices point to rows in the data
3 # 75 percent of the data set indices is selected in this list
4 index <- sample(1:nrow(data), round(0.75 * nrow(data)))
5
6 # The training data set is composed of 75 percent of the test data
7 train <- data[index,]
8
9 # The testing data set is composed of 25 percent of the test data
10 test <- data[-index,]

```

---

### 2.4.2 Linear Model for Comparison

First, we run the data through a simple linear model to see how well a linearly scaled system can predict the value of interest in the data set. Our data frame is constructed so that the last column, named `data$medv`, is what is predicted. Every other column is the value of a pixel in the image.

---

```

1 # Create a regression line using prog as the prognostic variable
2 lm.fit <- glm(medv ~ ., data=train)
3 # Print the summary of our fit
4 summary(lm.fit)
5 # Predict values using the fit and the testing data set
6 lm.pr <- predict(lm.fit, test)
7 # Calculate the Mean-Square Error of the linear model
8 lm.MSE <- sum((lm.pr - test$medv)^2) / nrow(test)

```

---

### 2.4.3 Neural Network Setup and Execution

We will compare the linear model to the neural network graphically after generating results from the neural network. We now work to set up the data for the neural network. Neural networks generally do not work well with non-uniform distributions of data. Therefore, we will scale every column to distribute the values correctly. When we scale the data set, we will lose the true values. Therefore, it is important to "unscale" the data once we generate output from the network.

---

```

1 # Create a list of maximums from each data column
2 maxs <- apply(data, 2, max)
3 # Create a list of minimums from each data column
4 mins <- apply(data, 2, min)
5
6 # Scale the values in each column to its respective limits
7 scaled <- as.data.frame(scale(data, center=mins, scale=maxs - mins))
8
9 # Redefine a training data set using the scaled values
10 train_ <- scaled[index,]
11 # Redefine a testing data set using the scaled values
12 test_ <- scaled[-index,]

```

---

At this point, we can generate a neural network and execute it. We begin by constructing a rule for the neural network to follow. This rule simply describes the output variables and the input variables that are thought to influence the output.

---

```

1 # Every column affects the prognostic variable
2 # We extract the name of every column in the data set
3 n <- names(train_)
4
5 # Using the normal fit construct (prog ~) does not work in this package
6 # Instead, we generate the command using the following code
7 # Will generate something like this: prog ~ c1 + c2 + c3 + ...
8 f <- as.formula(paste("medv ~", paste(n[!n %in% "medv"], collapse=" + ")))
9

```

---

```

10 # The following code will generate the neural network
11 # This step can take some time
12 nn <- neuralnet(f, data=train_, hidden=hiddenLayers, linear.output=T)
13
14 # We can plot the neural network to view its structure
15 plot(nn)

```

---

Just like the linear model that we use to compare, we will calculate the Mean-Squared Error of the neural network. For this, we extract the predictions of the neural network and "unscale" them.

---

```

1 # Compute predictions from our neural network
2 nn.pr_ <- compute(nn, test_[,1:13])
3
4 # Unscale the data output from the neural network
5 nn.pr <- nn.pr_$net.result * (max(data$medv) - min(data$medv)) + min(data$
   medv)
6 test.r <- (test_$medv) * (max(data$medv) - min(data$medv)) + min(data$medv)
7
8 # Root Mean Square Error computation for the neural network
9 nn.MSE <- sum((test.r - nn.pr)^2) / nrow(test_)

```

---

#### 2.4.4 Graphical Representation

Now, we simply use descriptive and graphical methods to highlight the difference between a linear fit model and the neural network.

---

```

1 # A higher number will suggest a lower correlation between the prediction
  and reality
2 print(paste(lm.MSE, nn.MSE))
3
4 # Set up a plot to hold two graphs
5 par(mfrow=c(1, 2))
6
7 # Plot the neural network's real vs. predicted values
8 plot(test$medv, nn.pr, col="red", main="Real vs. Predicted NN", pch=18, cex
      =0.7, xlab="Actual Value", ylab="Neural Network")
9 # Generate a line and legend for the plot
10 abline(0, 1, lwd=2)
11 legend("bottomright", legend="LM", pch=18, col="red", bty="n", cex=0.95)
12
13 # Plot the linear model's real vs. predicted values
14 plot(test$prog, lm.pr, col="blue", main="Real vs. Predicted LM", pch=18, cex
      =0.7, xlab="Actual Value", ylab="Linear Fit")
15 # Generate a line and legend for the plot
16 abline(0, 1, lwd=2)
17 legend("bottomright", legend="LM", pch=18, col="blue", bty="n", cex=0.95)
18
19 # Generate a plot with both scatter plots for better comparison
20 par(mfrow=c(1, 1))
21 plot(test$medv, nn.pr, col="red", main="Real vs. Predicted", pch=18, cex
      =0.7, xlab="Actual Value", ylab="Predicted Value")

```

---

```

22 points(test$medv, lm.pr, col="blue", pch=18, cex=0.7)
23 abline(0, 1, lwd=2)
24 legend("bottomright", legend=c("NN", "LM"), pch=18, col=c("red", "blue"))

```

---

A proper execution of the neural network should show that the Mean-Squared Error of the generalized linear fit model is much higher than that of the neural network. Furthermore, the graphs should portray a tighter clustering of the neural network data points around the linear fit.

## 2.5 Student Assignment - Temperature Dataset

The student is tasked with creating a neural network that interprets a dataset that contains the monthly temperatures from January to November over multiple years. The goal of the network is to predict the temperatures of the month of December. The following code should be used to parse the data:

---

```

1 # Name:           neuralStudent.R
2 # Author:        First Last
3 # Date:          Month Date Year
4
5 ##### PARSING #####
6
7 monthlyavg <- c(-4.07, -1.75, 2.17, 8.06, 13.90, 18.77, 21.50, 20.50, 16.27,
8   9.74, 2.79, -2.36)
9 rawdata <- read.csv("berkeleyusdata.txt", header = FALSE)
10 tempdevs <- rawdata$V4
11 temps <- rawdata$V4 + monthlyavg
12 months <- rawdata$V3
13 mat <- matrix(temps, ncol=12)
14 datafornnet <- data.frame(mat)
15 colnames(datafornnet) <- c("Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug",
16   "Sep", "Oct", "Nov", "Prog")
15 data <- datafornnet
16 rm(monthlyavg, rawdata, tempdevs, temps, months, mat, datafornnet)

```

---

Thus, the data set is stored in a variable called `data`. The neural network can now be built by the students. The name of the variable being predicted is `Prog`.

---

```

1 # Clear environment variables
2 rm(list = ls())
3
4 # Import the appropriate libraries
5 library(neuralnet)
6
7 # We do not need to set a seed
8 # We can if the pseudo-random system creates better results
9 # set.seed(500)
10
11 # Vary the number of neurons in the hidden layers to get better results
12 # This is an example a student may use
13 hiddenLayers = c(5, 3)

```

```

14
15 # Ensure that the data has no NaN values
16 apply(data, 2, function(x) sum(is.na(x)))
17
18 # Split the data into 75% and 25%
19 # One dataset is for training and the other is for testing
20 index <- sample(1:nrow(data), round(0.75 * nrow(data)))
21 train <- data[index,]
22 test <- data[-index,]
23
24 # Create a generalized linear fit to compare with the neural network
25 lm.fit <- glm(Prog ~ ., data=train)
26 # Print the summary of the fit
27 summary(lm.fit)
28 # Predict values using the fit
29 lm.pr <- predict(lm.fit, test)
30 # Calculate the Mean-Squared Error of the model
31 lm.MSE <- sum((lm.pr - test$Prog)^2) / nrow(test)
32
33 # Calculate mins and maxes for scaling
34 maxs <- apply(data, 2, max)
35 mins <- apply(data, 2, min)
36 # Scale using the range of the data
37 scaled <- as.data.frame(scale(data, center=mins, scale=maxs-mins))
38 # Create the scaled training and test data sets
39 train_ <- scaled[index,]
40 test_ <- scaled[-index,]
41
42 # Create the neural network
43 n <- names(train_)
44 f <- as.formula(paste("Prog ~", paste(n[!n %in% "Prog"], collapse=" + ")))
45 nn <- neuralnet(f, data=train_, hidden=hiddenLayers, linear.output=T)
46 plot(nn)
47
48 # Calculate the scaled predictions and errors
49 nn.pr_ <- compute(nn, test_[,1:11])
50 nn.pr <- nn.pr_$net.result * (max(data$Prog) - min(data$Prog)) + min(data$Prog)
51 test.r <- (test_$Prog) * (max(data$Prog) - min(data$Prog)) + min(data$Prog)
52 nn.MSE <- sum((test.r - nn.pr)^2) / nrow(test_)
53
54 # Print the resulting mean-squared errors
55 print(paste(lm.MSE, nn.MSE))
56
57 # Generate the graphics for displaying data
58 par(mfrow=c(1,2))
59 plot(test$Prog, nn.pr, col="red", main="Real vs. Predicted NN", pch=18, cex=0.7, xlab="Actual Value", ylab="Neural Network")
60 abline(0, 1, lwd=2)
61 legend("bottomright", legend="LM", pch=18, col="red", bty="n", cex=0.95)
62 plot(test$Prog, lm.pr, col="blue", main="Real vs. Predicted LM", pch=18, cex=0.7, xlab="Actual Value", ylab="Linear Fit")
63 abline(0, 1, lwd=2)
64 legend("bottomright", legend="LM", pch=18, col="blue", bty="n", cex=0.95)
65 par(mfrow=c(1,1))
66 plot(test$Prog, nn.pr, col="red", main="Real vs. Predicted", pch=18, cex

```

```
=0.7, xlab="Actual Value", ylab="Predicted Value")  
67 points(test$Prog, lm.pr, col="blue", pch=18, cex=0.7)  
68 abline(0, 1, lwd=2)  
69 legend("bottomright", legend=c("NN", "LM"), pch=18, col=c("red", "blue"))
```

---

## 2.6 Downloading The Dataset

Download the file located at <https://github.com/nosyarg/textbookdata/blob/master/berkeleyusdata.txt>. This data comes from the Berkeley Earth Project.



### 3. Shiny Package Tutorial

#### 3.1 Introductory Reading

Shiny is a package used in R, that enables the user to easily create interactive applications. Shiny can be used to create both simple histograms of disease data or interactive maps of the crime rate in different cities. With Shiny the possibilities are endless. In this lab, you will learn how to create a basic Shiny application. Look at the picture in Figure 6.1. ([gallery](#)) This is a sample Shiny app depicting a graph of many famous movies. The graph plots movies based on their rating and the number of reviews it has on Rotten Tomatoes. On the left of the graph you find various inputs, such as Oscars and genres, that can customize the output of the graph. These are the basic components of a Shiny app: inputs and outputs. Here are some examples of Shiny input functions:

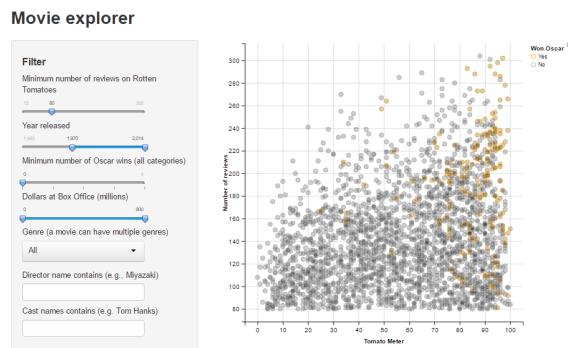


Figure 3.1: Sample Shiny app

- checkboxInput: a single check box
- dateInput: a calendar
- fileInput: file upload
- numericInput: enter number
- radioButtons: set of radio buttons
- selectInput: box with multiple choices
- sliderInput: a slider bar
- textInput: enter text

The users of a Shiny app will interact solely with the inputs that you, the developer, will allow them to. The above inputs can collect various types of data input from the user, to determine many different types of outputs. Here are some of the output functions in Shiny:

- imageOutput: creates image
- plotOutput: creates a plot
- tableOutput: creates a table
- textOutput: creates text

Understanding how inputs and outputs relate to each other is a key requirement for developing a Shiny app. In the model students will create today, the inputs will be dropdown boxes (selectInput) and the output will be a simple scatterplot (plotOutput).

## 3.2 Objectives

Students will be expected to learn the basic components of a Shiny App. Students will create their own application after learning the fundamentals of the package. This application will illustrate the progression of Olympic swimming medalists' times since the modern Olympics began.

- By the end of this lesson the student will understand the various parts of a Shiny app, while building a simple population prediction application.
- Students will also be able to create their own Shiny apps from scratch, and will be expected to create a Shiny app displaying Olympic swimming times over various events and years.

### 3.3 Building the Model

Shiny applications have a wide variety of uses. Today students will learn how to create a basic population prediction model.

To begin the application, open a new R script file in RStudio, clear the environment, set a working directory, and import the libraries necessary for this project. To clear the environment and set a directory use the following functions. Make sure you set a directory on your computer that the dataset is located in.

---

```
1 rm(list=ls())
2 setwd("C:/Users/Ishaan/Documents/NCSSM 2015-2016")
```

---

The libraries are as follows:

---

```
1 library(ggplot2)
2 library(scales)
3 library(ggmap)
4 library(dplyr)
5 library(car)
6 library(gcookbook)
7 library(shiny)
```

---

If you don't have one of these libraries installed on your computer, you can simply type

---

```
1 install.packages("insertlibraryname")
```

---

in the console to download them.

Next, we need to import the dataset of all the population data in NC for the next five years into a variable of your choice (I used popData). This can be done with the `read.csv` function. Make sure the data file is in directory you specified above.

---

```
1 popData <- read.csv("Population.csv")
```

---

Shiny applications have two components: a user-interface definition and a server script. The user-interface is used to control the design and layouts of the Shiny application. The server script for the Shiny application controls the algorithms and code that tell the Shiny app how to manipulate the inputs to produce outputs.(**shiny**) Three functions — `headerPanel`, `sidebarPanel`, and `mainPanel` — define the various regions of the user-interface. The application will be called "Population Statistics" so we specify that as the title when we create the header panel. The other panels are empty for now. We create a UI function with the aforementioned components, like this:

---

```
1 ui <- shinyUI(fluidPage(
2   headerPanel("Population Statistics"),
```

---

```
3 sidebarPanel(),
4 mainPanel()
5 ))
```

---

Next, define a skeletal server implementation. This can be done by creating a function called server with two parameters: input and output.

---

```
1 server <- function(input, output) {
2
3 }
```

---

The server function will be empty to begin with, but it will eventually be used to set relationships between the inputs and outputs of the application. At the end of your code, we need to call the shiny function so the application will run. To do that use the following syntax:

---

```
1 shinyApp(ui = ui, server = server)
```

---

The skeletal structure of the simplest Shiny application possible is shown in figure 3.2. Run the Shiny application by selecting all of your code and clicking the run button in the top right corner of RStudio. The blank Shiny application should open in a new window as shown in figure 3.2. To begin the model, we need to display the inputs

Population Statistics

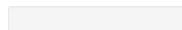


Figure 3.2: Blank Shiny App

of our program in the sidebar panel of the application. Our program is taking age group as an input and displaying a plot showing the predicted populations of that age group in North Carolina over the next five years. To do this, we will create a selectInput with all of the age groups as choices, like so:

---

```
1 sidebarPanel(
2   selectInput(
3     inputId = "age",
4     label = "Age:",
5     choices = c(
6       "Total", "0-4 years", "5-9 years",
7       "10-14 years", "15-19 years", "20-24 years",
8       "25-29 years", "30-34 years", "35-39 years",
```

```

9      "40-44 years", "45-49 years", "50-54 years",
10     "55-59 years", "60-64 years", "65-69 years",
11     "70-74 years", "75-79 years", "80-84 years",
12     "85+ years")
13   )
14 )

```

---

The inputId is simply a name that we give to each input, so we can access them in the server function. The label is the text that will appear on the application above the input box.

We also need to tell the application that our plot will be displayed in the main panel of the UI. This involves creating a plotOutput, like so:

```

1 mainPanel(
2   plotOutput(outputId = "myPlot")
3 )

```

---

Just like in `selectInput`, we need to give the output an ID, so we can use that name to access it in the server function. Once the model is run, the application should look like figure 3.3. Next we need to define the server-side of the application which



Figure 3.3: Adding Inputs

will accept inputs and compute outputs. Our server function will contain the code necessary to do this. We first tell the function that we are creating a plot as an output, and we can create the plot using the `renderPlot` command. The command `renderPlot` is what makes the application react. This means when a user changes the input, the output will instantly change. This is what makes Shiny so interactive. Within the `renderPlot` function we first select what data we will use and then plot that data using the `ggplot` function. To select our data, create a variable called `dataInput`, and take a subset of the `.csv` file we imported into R. We want to take only the data that the user specified, so instead of importing all the age group data, we only import the age group that the user selected. We can access that specific age group using `input$age`. The `ggplot` function takes in various inputs, such as the data, the variables, and the title. The `aes` function within `ggplot` indicates which axis is given to which variable, while the `geom_point` function draws the points on the graph. The `ggtitle` function creates a title over the plot. The completed `server` function should look like this:

---

```

1 server <- function(input, output) {
2   output$myPlot <- renderPlot({
3     dataInput <- subset(popData, (Age==input$age), select=c(Year,
4       Population, State))
5
6     ggplot(
7       data = dataInput,
8       aes(x=Year, y=Population)
9     ) + geom_point() + ggtitle("Population")
10   })
11 }
12 }
```

---

Finally, run the entire code including the `shinyApp` command. Your screen should look like the screen displayed in Figure 3.4. The application should have a working select box, a title, and a plot displaying the population data.

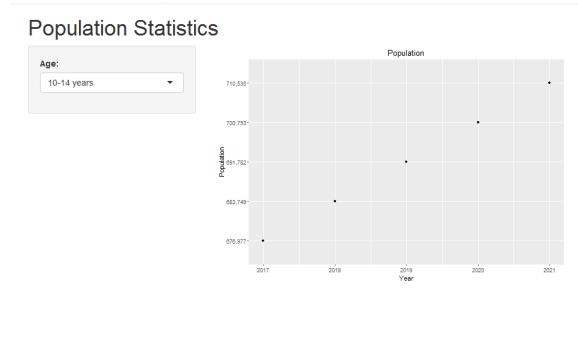


Figure 3.4: Final Model

## 3.4 Deliverable

You are expected to deliver properly formatted code written in the R language, using the Shiny package. Additionally, you will need to provide screenshots of various input combinations to demonstrate that your application is fully functional. This model will display Olympic swimming times over the modern Olympics. The variables for this model include time, year, and medal. The inputs are the event and gender of the swimmer. The graph should display the medaling times for each event and gender for each year of the modern Olympics. The data for this project is found within the `SwimmingExampleM` dataset provided.

### 3.4.1 Programming Hints

In the deliverable you will need to be able to add color to the model. It is quite simple. In the previous section we explained that the `ggplot` command makes the

actual plot. In order to create color we just add `color = variablename`, inside the `aes()` function.

---

```
1 aes(x=Time, y=Year, color=Medal)
```

---

The final shiny app screen should look like figure 3.5

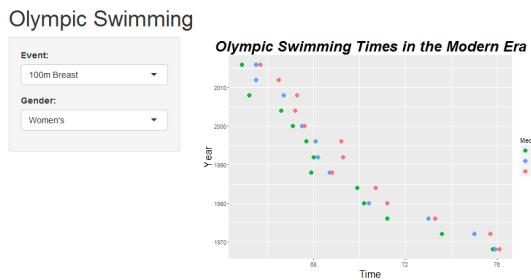


Figure 3.5: Deliverable

## 3.5 Teaching Code

---

```

1 #Ishaan Rao
2 #September 1, 2016
3 #pop.R
4
5 #clean up and set the directory
6 rm(list=ls())
7 setwd("C:/Users/Ishaan/Documents/NCSSM 2015-2016")
8
9 #load libraries
10 library(ggplot2)
11 library(scales)
12 library(ggmap)
13 library(dplyr)
14 library(car)
15 library(gcookbook)
16 library(shiny)
17
18 #read data file
19 popData <- read.csv("Population.csv")
20
21 #setting the server, code instructing R what to do with inputs
22 #output$my plot is what will be created when the app is run
23 #renderPlot allows Shiny to adjust when the user changes an option
24 #By defining the subset we indicate which inputs we want to be selected by
25 #the user
26 #select are the values in the dataset that are going to be plotted
27 #ggplot give direction to create plot, aes() defines the axes, ggtitle
#creates
28 title and geom_point creates the scatter plot
29 server <- function(input, output) {

```

```

29   output$myPlot <- renderPlot({
30     dataInput <- subset(popData, (Age==input$age), select=c(Year,
31                           Population, State))
32 
33     ggplot(data = dataInput, aes(x=Year, y=Population))
34     + geom_point() + ggtitle("Population")
35   })
36 
37 }
38 # Define UI for application that draws a scatterplot
39 #headerPanel creates the title
40 #sidebarPanel creates the user entered section
41 #selectInputs is where you insert the options that the user can select
42 #mainPanel outputs the app using all of the previous code we inputed
43 ui <- shinyUI(fluidPage(
44   headerPanel("Population Statistics"),
45   sidebarPanel(
46     selectInput(
47       inputId = "age",
48       label = "Age:",
49       choices = c(
50         "Total", "0-4 years", "5-9 years",
51         "10-14 years", "15-19 years", "20-24 years",
52         "25-29 years", "30-34 years", "35-39 years",
53         "40-44 years", "45-49 years", "50-54 years",
54         "55-59 years", "60-64 years", "65-69 years",
55         "70-74 years", "75-79 years", "80-84 years",
56         "85+ years")
57       )
58     ),
59     mainPanel(
60       plotOutput(outputId = "myPlot")
61     )
62   )))
63 #this runs the Shiny App
64 shinyApp(ui = ui, server = server)

```

---

### 3.6 Example Student Code

---

```

1 #Ishaan Rao
2 #September 1, 2016
3 #olympic.R
4
5 #load libraries
6 library(ggplot2)
7 library(scales)
8 library(ggmap)
9 library(dplyr)
10 library(car)
11 library(gcookbook)
12 library(shiny)

```

```
13
14 rm(list=ls())
15 setwd("C:/Users/Ishaan/Documents/NCSSM 2015-2016")
16
17 swimData <- read.csv("SwimExampleM.csv")
18 #setting the server, code instructing R what to do with inputs
19 #output\$my plot is what will be created when the app is run
20 #renderPlot allows Shiny to adjust when the user changes an option
21 #By defining the subset we indicate which inputs we want to be selected by
22     the user
23 #select are the values in the dataset that are going to be plotted
24 #ggplot give direction to create plot, aes() defines the axes, ggtitle
25     creates
26 title and geom_point creates the scatter plot
27 #scale_color_hue creates different colors for each of the three models\\\
28 server <- function(input, output) {
29
30     output\$myPlot <- renderPlot({
31
32         dataInput <- subset(swimData,
33             (Event==input\$event & Gender==input\$gender),
34             select=c(Time, Year, Medal, Name))
35
36         ggplot(data = dataInput, aes(x=Time, y=Year, color=Medal)) +
37             geom_point(size=3) +
38             ggtitle("Olympic Swimming Times in the Modern Era") +
39             theme(
40                 plot.title = element_text(color="black", size=24,
41                     family="Times", face="bold.italic"), axis.title = element_text(
42                         size=16)
43             ) +
44             scale_color_hue(breaks=c("Gold", "Silver", "Bronze"))
45     })
46 }
47 # Define UI for application that draws a scatterplot
48 #headerPanel creates the title
49 #sidebarPanel creates the user entered section
50 #selectInputs is where you insert the options that the user can select
51 #In the model we have two different inputs one for event and one for gender,
52     we
53 create two different selectInput commands to account for the options
54 #mainPanel outputs the app using all of the previous code we inputted
55 ui <- shinyUI(fluidPage(
56     headerPanel("Olympic Swimming"),
57     sidebarPanel(
58         selectInput(
59             inputId = "event",
60             label = "Event:",
61             choices = c(
62                 "50m Free", "100m Free", "100m Back", "100m Breast",
63                 "100m Butterfly", "200m Free", "200m Back", "200m Breast",
64                 "200m Butterfly", "200m IM", "400m IM"
65             ),
66             selected = "50m Free"
67         ),
68         selectInput(
```

```
65     inputId = "gender",
66     label = "Gender:",
67     choices = c("Men's", "Women's"),
68     selected = "Women's"
69   )
70 ),
71 mainPanel(
72   plotOutput("myPlot")
73 )
74 )))
75 #runs application
76 shinyApp(ui = ui, server = server)
```

---

### 3.7 Further Readings

- For a full tutorial on creating a Shiny app, visit this tutorial on the Shiny website: <http://shiny.rstudio.com/tutorial/>.
- To learn more advanced Shiny topics, functions, and capabilities visit: <http://shiny.rstudio.com/articles/>.
- To see sample Shiny apps in action visit Shiny's gallery: <http://shiny.rstudio.com/gallery/>



## 4. deSolve

### 4.1 Introduction to Computational Differential Solutions

Many of the systems studied in life sciences, environmental sciences, and physical sciences can be modeled with differential equations. Differential equations serve to model rates of change, which are often easier to measure than every interaction in the system. A differential equation is an equation for a function that relates the values of the function to the values of its derivatives. An ordinary differential equation (ODE) is a differential equation for a function of a single variable (such as  $x(t)$ ), while a partial differential equation (PDE) is a differential equation for a function of several variables (such as  $v(x, y, z, t)$ ). An ODE contains ordinary derivatives and a PDE contains partial derivatives. Typically, PDE's are much harder to solve than ODE's (**chasnov**).

Initial Value Problems (IVPs) involve differential equations where certain initial values in the model are known. The `deSolve` package uses computational means to analyze and compute the values of states in such models. In a differential equation, the state is the variable that is varied with respect to the independent variable in the system. Parameters, on the other hand, are values that signify some characteristic or attribute of the system. For example, in the given system of equations:

$$\frac{dY}{dt} = aY$$

The state of the differential equation is  $Y$ . The independent variable in this system is  $t$ . Finally,  $a$  is a parameter that can be varied to vary the solution. In this lab, we will analyze systems of ordinary first-order differential equations, which relate the first derivative to states and parameters.

### 4.1.1 Examples of Differential Equations

The *Malthusian law of population growth* is used to model the populations of certain kinds of organisms living in ideal environments for limited lengths of time  $t$ . It gives the rate at which a population  $p$  changes with respect to  $t$ . The value of the constant  $r$  depends on the organism.

$$\frac{dp}{dt} = rp$$

This *second-order reaction rate* law gives the rate at which a single chemical species combines to produce a new species, such as methyl radicals combining in a gas to form ethane molecules.

$$\frac{dx}{dt} = k(A - x)^2$$

This equation models the motion of a damped mass-spring system subjected to a time-dependent force  $F(t)$ .

$$m \frac{d^2x}{dt^2} + b \frac{dx}{dt} + kx = F(t)$$

## 4.2 Installation

To install the `deSolve` package in R, enter the following code into the command line. Installation may produce warnings that are generally inconsequential to the overall usage of the package. However, if the package does not install correctly, contact the package developers.

---

```

1 > install.packages("deSolve")
2 > library(deSolve)

```

---

## 4.3 Objectives of Case Study

There are certain objectives that students should focus to accomplish while completing this activity.

- Know what ordinary differential equations and partial differential equations are. For the case studies that follow, it is inconsequential for students to understand the structure and function of ordinary differential equations.
- Understand how initial value problems arise in ordinary differential equation systems in the natural sciences. That is, use initial values found in nature and differential equation models to produce valid solutions.

- Know how to set up a differential equation given a system and the relationship between rates. This is important for modeling a given system.
- Know how to use `deSolve` to computationally solve systems of ordinary first-order differential equations. This R package allows for computational solutions that effectively model the changes in the dependent variables of the system.

#### 4.4 Case Study - Damped Mass-Spring System

As discussed in the examples section, a well-known second-order differential equation is used to analyze the motion of an oscillating mass-spring.

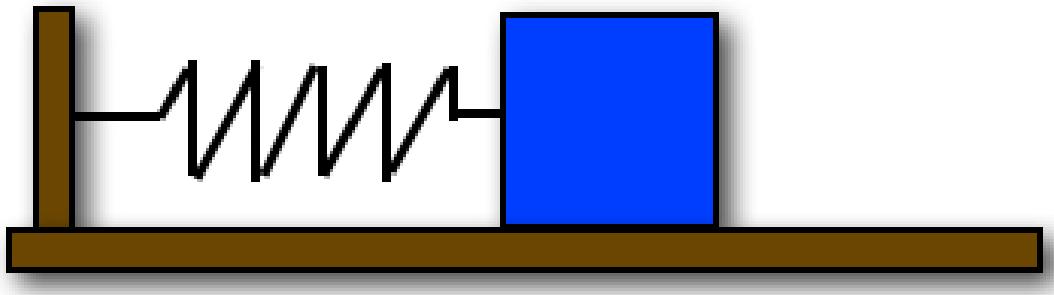


Figure 4.1: A mass-spring system

A mass-spring system is often called a harmonic oscillator. In an idealized system where no heat loss occurs, the mass-spring oscillator will continue oscillating forever.

The nature of a spring can be described using a single parameter, often noted as  $k$ . The spring constant  $k$  is used primarily in Hooke's Law, which relates the amount a spring stretches from rest,  $x$ , to the amount of force  $F$  that is applied on the mass  $m$ . Hooke's Law states that:

$$F = -kx$$

By Newton's Second Law, we also know that:

$$F = m \frac{d^2x}{dt^2}$$

Therefore, we can create a second-order differential equation for the idealized mass-spring oscillator:

$$m \frac{d^2x}{dt^2} = -kx$$

$$m \frac{d^2x}{dt^2} + kx = 0$$

Unfortunately, due to the complexity of behaviors and solutions in second-order differential equations, it is computationally easier to create a system of differential equations  $x'_1, x'_2, \dots, x'_n$  that model the system.

We begin with a few basic states that will represent the differential equations in our system.

$$x_1 = x$$

$$x_2 = \frac{dx}{dt}$$

We compute the derivatives of these states:

$$x'_1 = x' = x_2$$

$$x'_2 = x''$$

By solving the second-order differential equation in terms of the second derivative  $x''$ :

$$m \frac{d^2x}{dt^2} + kx = 0$$

$$mx'' + kx = 0$$

$$x'' = -\frac{k}{m}x = -\frac{k}{m}x_1$$

Therefore, the system of ordinary first-order differential equations that we are trying to solve is:

$$\frac{dx_1}{dt} = x_2$$

$$\frac{dx_2}{dt} = -\frac{k}{m}x_1$$

#### 4.4.1 Computational Solutions

We begin by setting up the environment to begin using the `deSolve` package. We import the library and clear the environment variables.

---

```

2 # Date:           Month Date Year
3 # Name:          DiffEq.R
4
5 # Import the library into R
6 library(deSolve)
7
8 # Clear variables from environment
9 rm(list = ls())

```

---

As stated in the introduction, we are solving an initial value problem. In this case, we use arbitrary values for states and parameters. We initialize  $x_1 = 100$  to suggest that the spring is stretched 100 meters past its rest phase. Furthermore,  $x_2 = 0$  to suggest that the motion begins from rest, when the rate of change is 0.

---

```

1 # m is the mass of the object (1 kg)
2 # k is the spring constant (50 kg / s^2)
3 parameters <- c(m = 1, k = 50)
4
5 # X is the first differential equation (100 m)
6 # Y is the second differential equation (0 m/s)
7 state <- c(X = 100, Y = 0)

```

---

We can define a function that represents the system of differential equations we have defined. Such a function will effectively relate the state and parameters to the derivatives in question.

---

```

1 # Function that defines the system
2 system <- function(t, state, parameters) {
3   with(as.list(c(state, parameters)), {
4     dX <- Y
5     dY <- -1*(k/m)X
6
7     list(c(dX, dY))
8   })
9 }

```

---

We run the simulation over a certain time interval with discrete differences between every time value. Finally, we run the `ode()` function to produce the discrete values.

---

```

1 # The sequence of time run in the simulation
2 times <- seq(0, 10, by=0.01)
3
4 # Generate curves that are solutions to the system of differential equations
5 out <- ode(y=state, times=times, func=system, parms=parameters)

```

---

Running the command above will generate data that is stored in `out`. We can plot this data out. From our derivation, we know that the value of  $x_1 = x$ . Therefore,

the solution to  $x_1$  is equal to the solution of the second-order differential equation. We also plot how the two variables,  $x_1$  and  $x_2$ , vary in comparison to each other. We plot the following:

---

```

1 plot(out[, "time"], out[, "X"], xlab="time", ylab="displacement", pch=19)
2 plot(out[, "X"], out[, "Y"], xlab="displacement", ylab="velocity", pch=19)

```

---

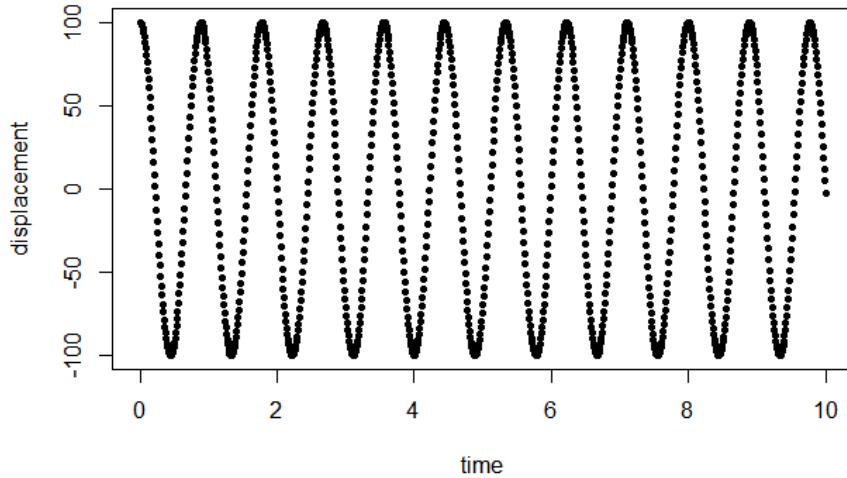


Figure 4.2: The displacement of the spring over time

#### 4.4.2 Improving the System

However, this idealized system does not capture the motion of a spring that is dampened by other factors, such as heat loss by friction. Such a system is known as a **Damped Mass-Spring Oscillator**. The spring loses momentum with relation to its velocity, which is the derivative of the position  $x$ . Therefore, we can add a new parameter to the second-order differential equation,  $b$ , that accounts for the dampening of the motion of the spring.

$$m \frac{d^2x}{dt^2} + b \frac{dx}{dt} + kx = 0$$

Using the method described above to convert the second-order differential equation into a system of first-order differential equations, we get:

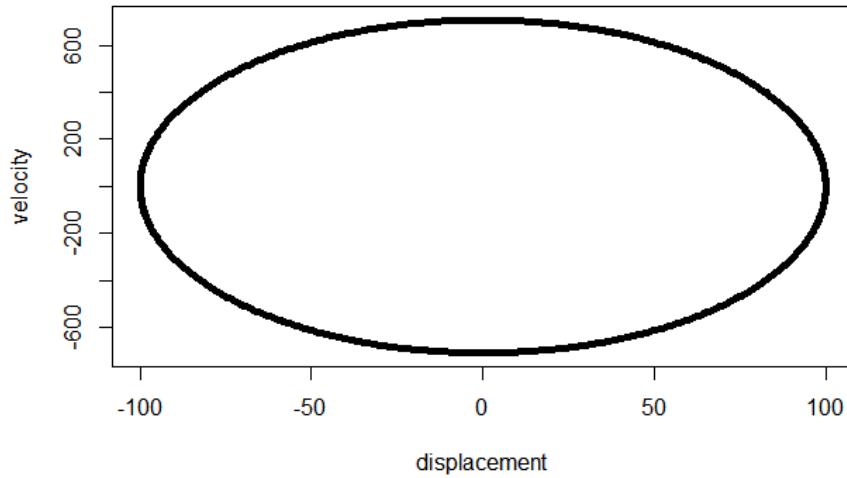


Figure 4.3: The relationship between  $x_1$  and  $x_2$

$$\begin{aligned}\frac{dx_1}{dt} &= x_2 \\ \frac{dx_2}{dt} &= -\frac{k}{m}x_1 - \frac{b}{m}x_2\end{aligned}$$

We modify our systems function in R to account for the velocity-dependent dampening of the spring as follows:

---

```

1 # m is the mass of the object (1 kg)
2 # b is the dampening of the object (0.25)
3 # k is the spring constant (50 kg / s^2)
4 parameters <- c(m = 1, b = 0.25, k = 50)
5
6 # Function that defines the system
7 system <- function(t, state, parameters) {
8   with(as.list(c(state, parameters)), {
9     dX <- Y
10    dY <- -1*(k/m)X - (b/m)Y
11  })
12}

```

---

Running the simulation again, we get the graphs shown in figures 4.4 and 4.5. Notice that the range of the motion of the spring-mass system continually decreases as time passes. This accounts for much of the heat loss that a damped spring experiences.

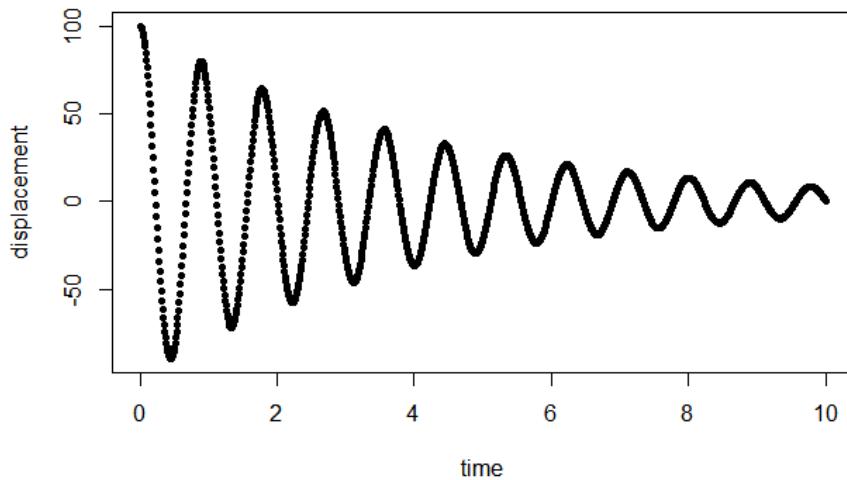


Figure 4.4: The displacement of the spring over time

## 4.5 Student Assignment - Modelling of Genetic Processes

Many of the genetic processes that are studied in the field of molecular biology can be modeled using systems of differential equations. The student is tasked with understanding many of these systems and modeling them. It is important for us to go over the *central dogma of molecular biology* before proceeding any further.

The three main components of this dogma are DNA, mRNA, and proteins. Proteins are generally the effectors in our cells, responsible for life-critical tasks, enzymatic processes, and regulation. Every protein has its own function and specificity. Ribosomal subunits, responsible for synthesizing proteins, get instructions for synthesis from polymer strands made of nucleotides known as mRNA. mRNA strands, in turn, are synthesized using instructions from another polymer known as DNA, which is the only heritable molecule present in the cell (as long as we ignore mitochondrial DNA that is passed down from a maternal lineage). We will also quickly review the processes of transcription and translation in depth. However, a strong foundation in biology is not necessary for modelling the problems that will follow.

- **Transcription:** Transcription is the process by which the information contained in a section of DNA (most often a gene) is transferred to a newly assembled piece of messenger RNA (mRNA). Transcription is triggered by the binding of RNA Polymerase to specific regions on the DNA called promoters, and is enabled/facilitated or impeded by a range of promoter-specific proteins called transcription factors (**stan**).
- **Translation:** During translation, proteins are produced by ribosomes that bind to a specific site of the mRNA (the ribosome binding site - RBS) and

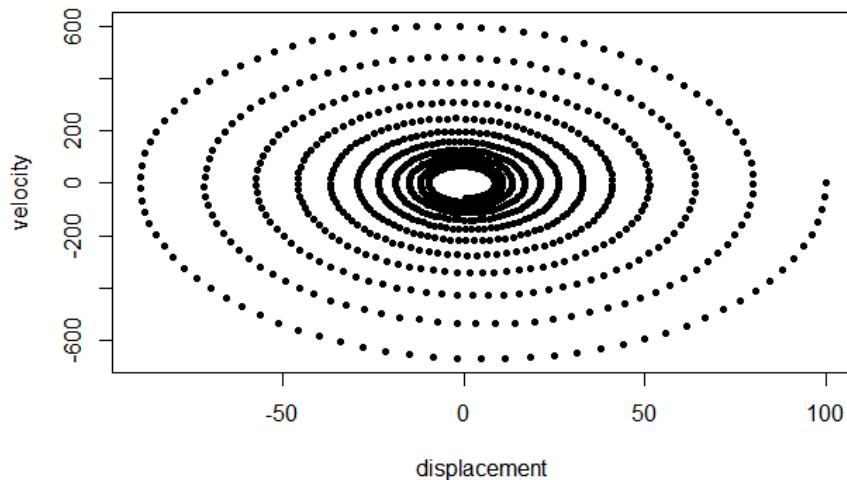


Figure 4.5: The relationship between  $x_1$  and  $x_2$

then move along the mRNA chain converting triplets of bases or “codons” on the mRNA into the appropriate peptide chain of amino acids defining the desired protein. The mRNA is read by the ribosome as triplet of bases, usually beginning with an AUG, or an initiator methionine codon downstream of the ribosome binding site. Complexes of initiation factors and elongation factors bring aminoacylated transfer RNAs (tRNAs) into the ribosome-mRNA complex, matching the codon in the mRNA to the anti-codon in the tRNA, thereby adding the correct amino acid in the growing peptide sequence encoding the emerging protein. As the amino acids are linked together into the growing peptide chain, they begin folding into the correct conformation. This folding continues until the nascent polypeptide chain is released from the ribosome as a mature protein (**stan**).

#### 4.5.1 Constitutive Gene Expression

**Problem Statement:** Model the system of constitutive gene expression in terms of the central dogma of molecular biology, mRNA degradation, and protein degradation.

Constitutive gene expression is the process that occurs when a gene is always on. That is, the gene is not regulated in any manner and simply follows the central dogma of biology. In such a case, we have to consider a few rates that occur in the system.

- $\frac{dm}{dt}$  is the rate at which mRNA is synthesized from DNA. It will be one of the differential equations students will use in their system.

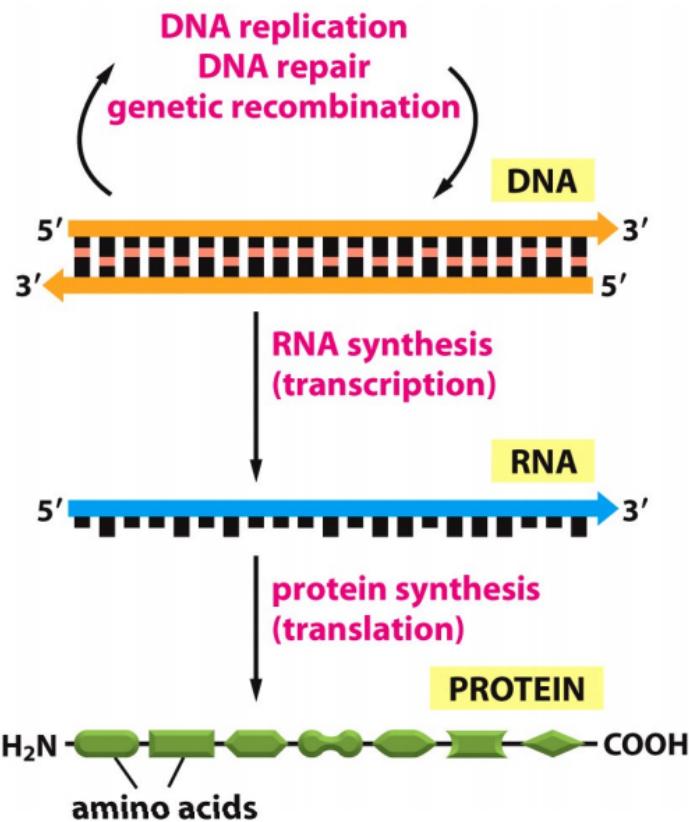


Figure 7-1 Essential Cell Biology 3/e (© Garland Science 2010)

Figure 4.6: The Central Dogma of Molecular Biology

- $\frac{dp}{dt}$  is the rate at which the protein in question is synthesized from DNA. It will be the other differential equation students will consider in their system.
- $k_m$  is the constitutive transcription rate. It is considered to be constant, and it represents the number of mRNA molecules produced per gene, per unit of time. We assume that there is only one copy of the gene in the cell. In the case of bacterial cells, where plasmid-located genes occur in multiple copies in different plasmids, we would multiply the rate  $k_1$  by the number of copies  $N$  in the cell to obtain a proper rate.
- $k_p$  is the translation rate. It is considered to be constant, and it represents the number of protein molecules produced per mRNA molecule, per unit of time.
- $d_m$  is the mRNA degradation rate, measured as the frequency at which mRNA molecules degrade.
- $d_p$  is the protein degradation rate, measured as the frequency at which protein molecules degrade.

### Sample Student Solution

The student should be capable of modeling the system using the rates provided above. The system we are considering should account for transcription, translation, and the degradation of components in the process.

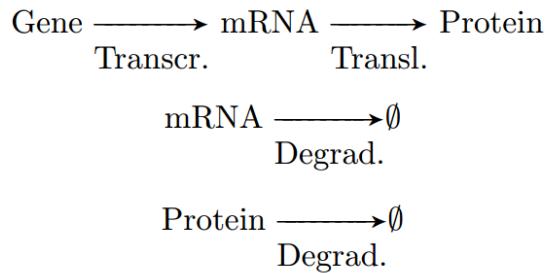


Figure 4.7: The constitutive gene expression system

The rate at which mRNA is transcribed is equal to the rate provided to us,  $k_m$ . Furthermore, we account for the rate of degradation by multiplying the frequency rate  $d_m$  by the number of mRNA molecules in the system. The rate at which protein molecules are translated is equal to the rate  $k_p$ . Similar to the mRNA system, the degradation occurs as a product of  $d_p$  and the amount of protein produced.

$$\begin{aligned}
 \frac{dm}{dt} &= k_m - d_m m \\
 \frac{dp}{dt} &= k_p m - d_p p
 \end{aligned}$$

The student simply has to change certain aspects of the program written by the instructor. In this case, the student must define a new set of states, parameters, and a new system function.

---

```

1 library(deSolve)
2
3 parameters <- c(km = 10, kp = 5, dm = 1, dp = 3)
4 state <- c(M = 0, P = 0)
5
6 system <- function(t, state, parameters) {
7   with(as.list(c(state, parameters)), {
8     dM <- km - dm*M
9     dP <- dp*M - dp*P
10
11     list(c(dM, dP))
12   })
13 }
14
15 times <- seq(0, 10, by=0.01)
16
17 out <- ode(y = state, times = times, func = system, parms = parameters)
18
19 # Plotting Functions

```

```

20 xaxis.range <- range(times)
21 yaxis.range <- range(min(range(out[,"M"]), range(out[,"P"])), max(range(out
22   [, "M"]), range(out[, "P"])))
23 plot(out[, "time"], out[, "M"], xlim=xaxis.range, ylim=yaxis.range, col="red",
24   xlab="Time", ylab="Concentration", main="System", pch=19)
25 points(out[, "time"], out[, "P"], col="blue", pch=19)
26 legend("bottomright", c("mRNA", "Protein"), pch=c(19, 19), col=c("red", "blue"))

```

---

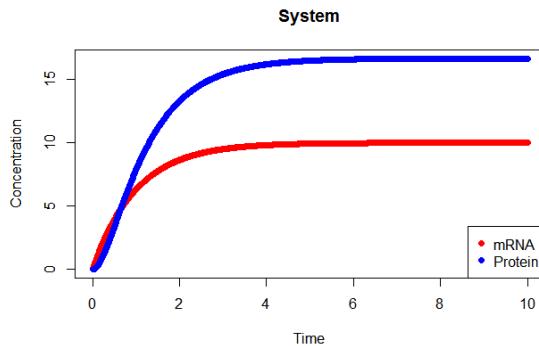


Figure 4.8: Constitutive gene expression concentrations in the system

We see that the concentration of proteins and mRNA in the system reaches a steady state with the initial values supplied. It is important for the student to understand that varying the parameters can change the solution to the system. For example, if the rate of degradation of a protein is much larger than the rate at which it is produced, we expect to see no overall protein production. Therefore, it is just as important to choose appropriate initial values as it is to produce a sound differential model. Although arbitrary values are chosen for these samples, students are encouraged to read literature in their field of research or study to find common values for such problems. For instance, the average half life of mRNA in the model system *E.coli* is around 5 minutes. Such information is inconsequential to producing an appropriate model.

#### 4.5.2 Gene Transcription Regulation

Interestingly, most genes that are known today do not follow the constitutive gene expression model. The constant production of proteins can be unnecessary, if not harmful, to the proper function of the organism. Therefore, multiple layers of regulation occur to ensure that the gene produces proteins at the correct time and the proper rate. There are two major arguments for the necessity of gene regulation in cellular and molecular biology. Firstly, all cells carry the same DNA (and therefore, the same heritable information) but have different developmental stages and functions. Secondly, cells are subject to many environmental perturbations, which can only be dealt with by regulation.

Proteins known as transcription factors are generally responsible for the regulation of genes. Transcription factors attach to a region upstream of the gene known as the promoter. The transcription factor can either be an activator (promotes the transcription of the gene) or a repressor (inhibits the transcription of the gene).

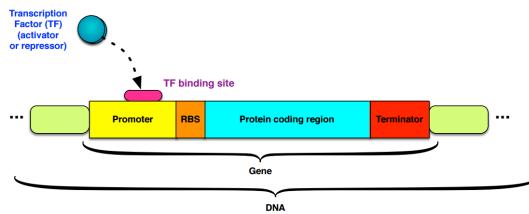


Figure 4.9: Action of regulative transcription factors on protein-encoding DNA regions

### Activators and Repressors

**Problem Statement:** Model the case of a gene whose transcription is activated by the binding of an activator/repressor to its transcription factor binding site.

A gene begins transcription once an activator binds to its corresponding transcription factor binding site. Once activated, the gene goes through the same process that it would during constitutive gene expression. These processes are inhibited once a repressor binds to the transcription factor binding site. This process as a whole is the basis of gene regulation. Because of the fact that, save for the introduction of activators and repressors, the process that the gene undergoes is identical to constitutive gene expression, many of the rates from the constitutive gene expression section are the same in this case. These include  $\frac{dm}{dt}$ ,  $\frac{dp}{dt}$ ,  $k_m$ ,  $k_p$ ,  $d_m$ , and  $d_p$ . With the introduction of a new parameter (an activator or a repressor), there are subsequently new variables added to this model. The repressor and activator cannot occur at the same time, as the gene can only be activated or inhibited at a certain point in time. These are as follows.

- $A$  represents the activator.
- $R$  represents the repressor.
- $K$  represents the activation or repression coefficient.
- $n$  represents the Hill coefficient, or the number of activators/repressors that need to cooperatively bind to the promoter (transcription factor binding site) in order for gene expression to be activated/inhibited.

### Sample Student Solution

The student should be capable of modeling the system using the rates provided in the constitutive gene expression section, as well as using the new variables

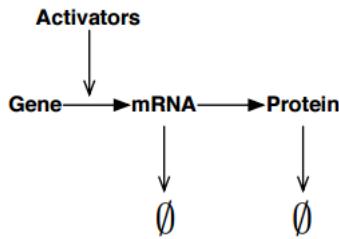


Figure 4.10: The activation and subsequent activities of a gene

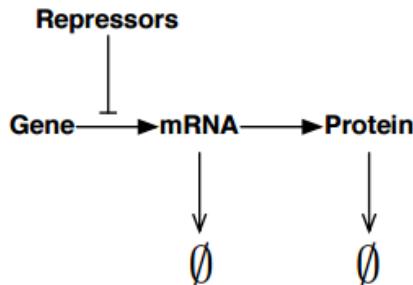


Figure 4.11: The inhibition and subsequent activities of a gene

provided above. The system we are considering should account for transcription, translation, and the degradation of components in the process.

The differential equations used in this model are almost identical to the differential equations used in previous models (the constitutive gene expression model). The only difference is the introduction of the activator/repressor and its coefficients, which is then multiplied by  $k_m$ .

These are the new differential equations with the introduction of the activator:

$$\frac{dm}{dt} = k_m \frac{A^n}{K^n + A^n} - d_m m$$

$$\frac{dp}{dt} = k_p m - d_p p$$

These are the new differential equations with the introduction of the repressor:

$$\frac{dm}{dt} = k_m \frac{K^n}{K^n + R^n} - d_m m$$

$$\frac{dp}{dt} = k_p m - d_p p$$

As with other cases, the student must define a new set of parameters and system functions in the code. Most of the code remains the same. The following shows the code when using the activator differential equations.

---

```
library(deSolve)
```

```

2
3 parameters <- c(km = 10, kp = 5, dm = 4, dp = 4, A = 0.25, K = 0.01, n = 4)
4 state <- c(M = 0, P = 0)
5
6 system <- function(t, state, parameters) {
7   with(as.list(c(state, parameters)), {
8     dM <- km*(A^n/(K^n+A^n)) - dm*M
9     dP <- dp*M - dp*P
10
11   list(c(dM, dP))
12 })
13 }

```

---

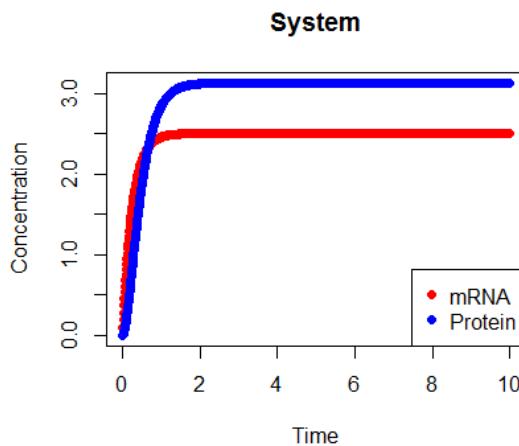


Figure 4.12: Gene expression in a system containing activators

### 4.5.3 Regulation

Unfortunately, by simply solving each system individually, it is hard to tell how activation and inhibition affect gene expression models. With the ability to model constitutive gene expression along with activation and regulation, we can compare models where activation and regulation does occur. Using the following values, we can create three concurrent systems and compare them:

---

```

1 # New parameters in the code
2 parameters <- c(km = 0.346, kp = 6.931, dm = 0.138, dp = 0.017, A = 200, R =
200, K = 100, n = 2)
3 state <- c(M = 0, P = 0)
4
5 # Constitutive Gene Expression Model
6 system.const <- function(t, state, parameters) {
7   with(as.list(c(state, parameters)), {
8     dM <- km - dm*M
9     dP <- kp*M - kp*P

```

```

10
11     list(c(dM, dP))
12 }
13 }
14
15 # Activation in Gene Expression
16 system.act <- function(t, state, parameters) {
17   with(as.list(c(state, parameters)), {
18     dM <- dm*(A^n / (K^n + A^n)) - dm*M
19     dP <- kp*M - dp*P
20
21     list(c(dM, dP))
22   })
23 }
24
25 # Inhibition in Gene Expression
26 system.rep <- function(t, state, parameters) {
27   with(as.list(c(state, parameters)), {
28     dM <- kM*(K^n / (K^n + R^n)) - dm*M
29     dP <- kp*M - dp*P
30
31     list(c(dM, dP))
32   })
33 }
34
35 # Independent variable
36 times <- seq(0, 400, by=1)
37
38 # Computing the solution to all three models
39 out.const <- ode(y = state, times = times, func = system.const, parms =
40   parameters)
40 out.act <- ode(y = state, times = times, func = system.act, parms =
41   parameters)
41 out.rep <- ode(y = state, times = times, func = system.rep, parms =
42   parameters)
42
43 # Plotting functions
44 plot(out.const[, "time"], out.const[, "P"], col="red", xlab="Time", ylab="Concentration", main="System", pch=19)
45 points(out.act[, "time"], out.act[, "P"], col="green", pch=19)
46 points(out.rep[, "time"], out.rep[, "P"], col="yellow", pch=19)
47 legend("bottomright", c("Protein CONST", "Protein ACT", "Protein REP"), pch=
48   c(19, 19, 19), col=c("red", "green", "yellow"))

```

---

The code will produce the graphics shown in figure 4.13. Notice that the constitutive protein expression runs unaffected. The activation gene expression curve takes longer to achieve expression because it takes time for activators to act as transcription factors before transcription can take place. Furthermore, the repressed curve is much lower than both, as is expected from such systems.

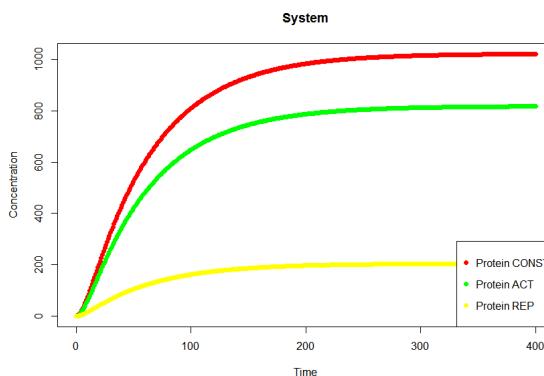


Figure 4.13: Comparing constitutive, activation, and inhibition in gene expression.

## 4.6 Conclusion

Differential equations are used in a variety of different areas to model a wide range of functions, such as the rate that a protein is transcribed or the rate that a population grows. The main function of a differential equation is to model a rate of change. Some of these equations are easily solved by hand, whereas others are more complicated; the `deSolve` package in R serves to solve and graph a system of differential equation once certain parameters and initial values are entered. Through learning and using `deSolve`, it is possible to easily solve a system of complex differential equations, as well as view the results graphically and be able to change parameters and values as needed. This streamlines a process of solving equations that occur within all types of science, math, and other processes.

## 4.7 Acknowledgements

We would like to thank the creators of the `deSolve` package, Karline Soetaert, Thomas Petzoldt, and R. Woodrow Setzer of the Royal Netherlands Institute of Sea Research. We would also like to thank Mr. Robert Gotwals and the North Carolina School of Science and Mathematics (NCSSM).



## 5. Clustering

Often, we are not only interested in specific data points, but also in the relationships between those data points. Cluster analysis refers to the grouping of data points into *clusters* based on similar characteristics. R has variety of libraries that support cluster analysis. We will use `cluster.` datasets for this lesson, in which we will analyze clustering in data about crime rates in US cities. (**Novomestky**)

Here are some typical commands used in cluster analysis:

1. `scale`
2. `k.means.fit`
3. `attributes`
4. `wssplot`
5. `clusplot`
6. `dist`
7. `hclust`
8. `groups`
9. `rect.hclust`
10. `table`

For this study, you will learn a basic framework for cluster analysis using two different methods, K-means and hierarchical clustering. Each method produces a different type of cluster graph of the same data.

- Your first task is to graph a set of data using the K-means clustering method.
- Your second task is to graph a set of data using the hierarchical clustering method.

## 5.1 Building the Model

You should begin by copying and pasting the teaching code into RStudio. This will serve as a foundation upon which to build your code. We have included comments to guide you as you complete the code by adding commands. When you successfully complete the code, it should output two graphics, shown in figures 5.1 and 5.2.

### 5.1.1 Programming Hints

1. Different datasets will often require different clustering methods. Therefore, not all clustering methods will work (or at least, work well) for every dataset! Be prepared to use a variety of clustering methods in the future.
2. When installing packages, install them in the "console" window of RStudio.
3. If a piece of code is indented in the teaching or student code, it must also be indented in the code that you write. This is important!

## 5.2 Deliverable

The final product that you create should consist of one K-means cluster graph composed of groups of data points embedded in ovals and one hierarchical tree plot.

## 5.3 Teaching Code

---

```
1 # Your Name
2 # Date
3 # Script name: dataset.R
4 # Objective of script
5 #
6 # set working directory where results will be saved
7 setwd("file location")
8 #
9 #
10 # install the necessary package and datasets, then read in library
11 install.packages("data package")
```

```
12 # this particular step should be completed in the console window
13 library(data package)
14 data(dataset, package = 'data package')
15 # "dataset" is the name that you choose to call your specific dataset
16 head(dataset)
17 #
18 # now we will begin to complete the first objective
19 #
20 # begin by standardizing the variables
21 dataset <- scale(original.dataset.name[-1])
22 #
23 # this next piece of code is the K-means function, where the number
24 # # symbol denotes the number of cluster groups
25 k.means.fit <- kmeans(dataset, #)
26 #
27 # this code segment includes all of the outputs of the k.means.fit function
28 attributes((k.means.fit))
29 #
30 # this next piece of code denotes the centroids
31 k.means.fit$centers # a nice table should result
32 #
33 # this piece of code shows the clusters
34 k.means.fit$cluster
35 #
36 # this code string denotes the size of the clusters
37 k.means.fit$size
38 #
39 # now we will begin to complete the first objective
40 #
41 # load the cluster library
42 library(cluster)
43 #
44 # this next section of code creates the a cluster graph using the K-means
        method
45 clusplot(dataset, k.means.fit$cluster, main = "Name of Cluster Graph", color
        = TRUE,
46 shade = TRUE, labels = 2, lines = 0) # look at those nice colorful ovals!
47 #
48 # this next line creates a confusion matrix which evaluates the performance
        of the clustering
49 table(dataset[,1],k.means.fit$cluster)
50 #
51 # Okay, now that you've mastered k-means clustering,
52 # let's move on to hierarchical clustering
53 #
54 # this code segments defines the distance "d", which in
55 # this case is Euclidean, for the clustering algorithm
56 d <- dist(us.crime, method = "euclidean")
57 #
58 H.fit <- hclust(d, method = "ward.D")
59 #
60 # this displays the dendrogram
61 plot(H.fit)
62 #
63 # this separates the diagram into 4 groups...
64 groups <- cutree(H.fit, k=4)
```

---

```

65 #
66 # ... while this piece of code creates a red border around those same groups
67 rect.hclust(H.fit, k=4, border = "blue")
68 #
69 # lastly, another confusion matrix is used to evaluate the clustering
70 table(us.crime[,1],groups)
71 #
72 # end code

```

---

## 5.4 Example Student Code

---

```

1
2 # Katherine L. Bennett
3 # September 1, 2016
4 # Script name: crime.R
5 # The purpose of this script is to analyse a sample of U.S. city crime
6 # using two cluster analysis methods: K-means and Hierarchical
7 #
8 setwd("/Users/Katherine/Desktop/RFiles")
9 #
10 library(cluster.datasets)
11 data(us.crime, package = 'cluster.datasets')
12 head(us.crime)
13 #
14 us.crime <- scale(sample.us.city.crime.1970[-1])
15 #
16 k.means.fit <- kmeans(us.crime, 4)
17 #
18 attributes((k.means.fit))
19 k.means.fit$centers
20 k.means.fit$cluster
21 k.means.fit$size
22 #
23 library(cluster)
24 #
25 # the following two lines should be one line in your code
26 # we have separated them in this guide for the reader's convenience
27 #
28 clusplot(us.crime, k.means.fit$cluster, main = "U.S. City Crime Cluster
   Representation", color = TRUE,
29 shade = TRUE, labels = 2, lines = 0)
30 #
31 table(us.crime[,1],k.means.fit$cluster)
32 #
33 # hierarchical cluster method
34 #
35 d <- dist(us.crime, method = "euclidean")
36 H.fit <- hclust(d, method = "ward.D")
37 plot(H.fit)
38 groups <- cutree(H.fit, k=4)
39 rect.hclust(H.fit, k=4, border = "blue")
40 table(us.crime[,1],groups)

```

```
41 #  
42 # end code
```

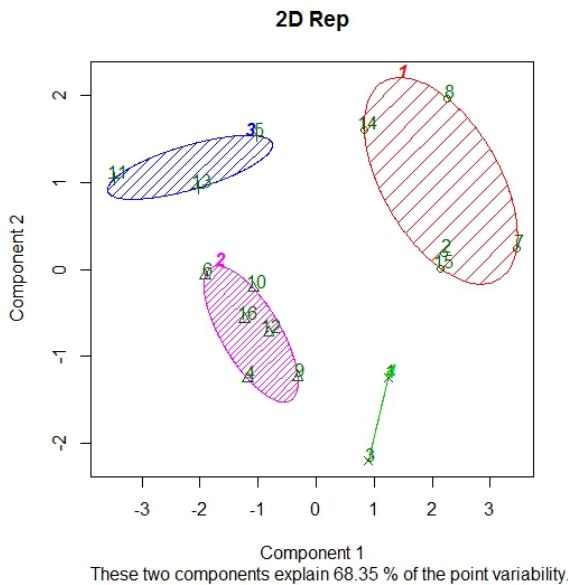


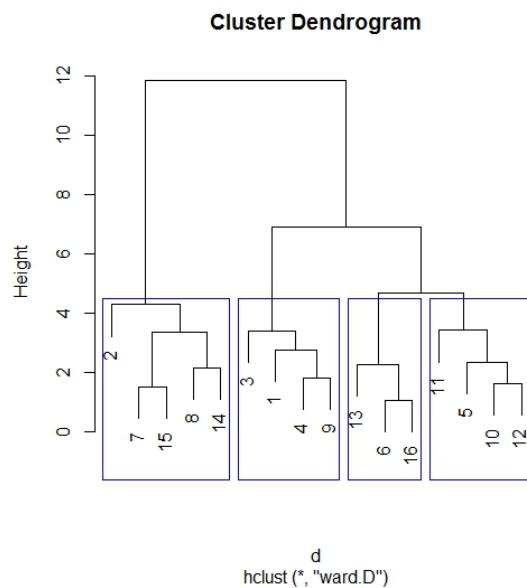
Figure 5.1: K-means Clustering Method Result

## 5.5 Further Readings

There are a wide variety of sources available that cover cluster analysis. We recommend chapter 8 of *Introduction to Data Mining* by Tan, Steinbach, and Kumar (2006). A link to the chapter is provided below:

<https://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>

Figure 5.2: Hierarchical Clustering Method Result





## 6. Microseq

### 6.1 Introductory Reading

The conversion of working DNA to proteins is one of the most important processes in biology. After DNA undergoes transcription and becomes an mRNA sequence, it can be turned into a string of amino acids via sequences known as codons. Each codon consists of a combination of three nucleotides, and there are twenty amino acids that match with these combinations.

DNA and all DNA processes are integral to scientists' understanding of biology. Studying this molecule is a fascinating and essential task, but because of the large quantities of data that genetics studies can obtain, manually analyzing a sequence can prove to be insignificant. Bioinformatics is a relatively recent field that combines computational science, statistics, engineering, and biology in order to interpret this type of data. Being able to use technology to aid in the analysis of different sequences can greatly lessen the tediousness of translating sequences and manually visualizing data. It also eliminates the possibility of human error, and allows genetics research to move much faster.

For example, existing algorithms, such as JASPAR and STARR-seq, are databases or functional tools that allow bioinformaticians to analyze sequence data. These programs can do things such as identify even remotely homologous sequences from a DNA multiple sequence alignment, or identify *cis*-regulatory regions throughout the genome. ([jaspar](#)) The focus of this case study are R packages known as Biostrings and microseq, which, in conjunction, allow for sequence data analysis.

## 6.2 Objectives of the Case Study

The objective of this study is to use the Biostrings and microseq packages in conjunction to translate a plain sequence format of DNA.

### 6.2.1 Objective 1

Translation of DNA to RNA is a vital biological process, and using the Biostrings and microseq functions will allow for a significantly streamlined process that will aid projects or studies in need of translated DNA. The objective is to obtain an amino acid sequence translated automatically from the given DNA sequence.

## 6.3 Building the Model

Tips for building the model are located in the comment sections of the code.

If the code is written properly, the output should be a list with the frequencies of each amino acid in the DNA sequence. The output should look similar to figure 6.1.

```
> alphabetFrequency(dn, as.prob=TRUE)
      A          R          N          D          C          Q          E          G          H
[1,] 0.03918031 0.04961676 0.03601664 0.023308 0.03593311 0.03796178 0.03690582 0.05007206 0.03404047
      I          L          K          M          F          P          S          T          W
[1,] 0.05424414 0.1079028 0.05324396 0.01835742 0.05273559 0.04981224 0.08967853 0.05192683 0.01759073
      Y          V          U          O          B          J          Z          X          *          -          +          .          other
[1,] 0.03112627 0.0520845 0 0 0 0 0 0.02775455 0.05050743 0 0 0 0
```

Figure 6.1: Frequency of amino acids

## 6.4 Deliverable

The student should submit a list containing the frequencies of each amino acid in the sequence analyzed. This list can be compared to the instructor's frequencies list to check for accuracy.

## 6.5 Teaching Code

You should also put any notes, teaching suggestions, etc. in this space!

---

```

1 # Name
2 # Date
3 # Name of script
4 # Objective of script
5
6 # clean up and set working directory
7 rm(list=ls())
8 setwd()
9
10 ### download data set from
11 ftp://ftp.ensembl.org/pub/release-85/fasta/mus_musculus/dna/
12
13 ### install Biostrings from
```

```
14 https://bioconductor.org/packages/release/bioc/html/Biostrings.html
15 ### then install microseq
16
17 #load microseq and Biostrings library
18 library()
19 library()
20
21 # view accepted IUPAC DNA letters
22 # these are the values that the functions will accept and put out
23 DNA_BASES
24 DNA_ALPHABET
25
26 # obtain plain sequence format from .txt file on desktop or source
27 # turn fasta into string
28 d <- readDNAStringSet("data set")
29 # translate object to amino acid sequence
30 dn <- translate(data set, genetic.code=GENETIC_CODE, if.fuzzy.codon = "solve"
  ")
31 ##"solve" allows unidentifiable codons to be represented as *
32
33 #characteristics of DNA strings
34 d
35 summary(data set)
36 dn
37 summary(translated data set)
38
39 #frequency of each letter or present symbol
40 alphabetFrequency(translated data set)
41
42 #frequency values
43 alphabetFrequency(translated data set, as.prob=TRUE)
44
45 ### EXTRA CREDIT!
46 letterFrequency(translated data set, "A", as.prob=TRUE)
47 letterFrequency(data set[[1]], letters="ACGT", OR=0)
48
49 # What do these two functions show?
50 #show frequency of base "A" in translated data set, and quantity of "ACGT"
      in
51 original data set
```

---

## 6.6 Example Student Code

---

```
1 # Christa Parrish
2 # September 8, 2016
3 # translate.R
4 # this script translates DNA sequence to Amino Acid sequence and observes
  frequencies
5 of letters
6
7 # clean up and set working directory
8 rm(list=ls())
```

```
9 setwd("~/Desktop/R")
10
11 ### download data set from ftp://ftp.ensembl.org/pub/release-85/fasta/mus_
12 musculus/dna/
13
14 #load microseq and Biostrings library
15 library(microseq)
16 library(Biostrings)
17
18 #view accepted IUPAC DNA letters
19 DNA_BASES
20 DNA_ALPHABET
21
22 # obtain plain sequence format from .txt file on desktop or source,
23 # whichever works
24 # turn fasta into object
25 d <- readDNAStringSet("MusMusculus.fa")
26 # translate object
27 dn <- translate(d, genetic.code=GENETIC_CODE, if.fuzzy.codon = "solve")
28 # "solve" allows unidentifiable codons to be represented as *
29
30 #characteristics of DNA strings
31 d
32 summary(d)
33 dn
34 summary(dn)
35
36 #frequency of each letter or present symbol
37 alphabetFrequency(dn)
38
39 #frequency values
40 alphabetFrequency(dn, as.prob=TRUE)
41
42 ### EXTRA CREDIT!
43 letterFrequency(dn, "A", as.prob=TRUE)
44 letterFrequency(d[[1]], letters="ACGT", OR=0)
45
46 # What do these two functions show?
```

---

## 6.7 Further Readings

By looking at the package manual online, students can read about the other functions of the microseq and Biostrings packages.

## 6.8 Introduction to Biodiversity

The 21st century marks the era of big data. This increase of data has effected all spheres of science, especially the study of biodiversity. Biodiversity data covers an increasingly wide range of areas including bio-geography, ecology, invasive species biology, and climate change. The information collected in biodiversity informatics projects has been used for studies relating to food security, control of disease vectors, and marine productivity (**Barve**).

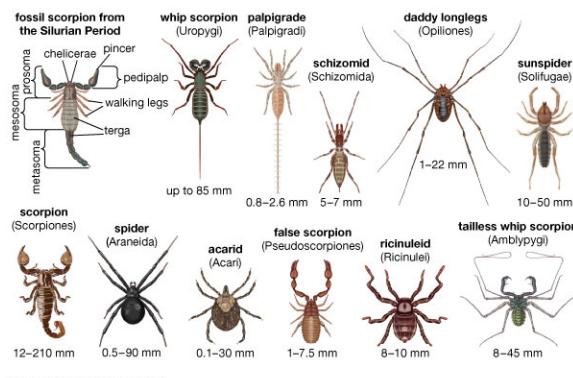


Figure 6.2: Example of biodiversity in arachnids

Primary biodiversity data is one of the most important biodiversity statistics; it documents a species' occurrences in time and space. Associated metadata, or data that describes statistics given in another set, may include information about who identified and recorded the species, who verified the identification, climatic parameters, micro-habitat information, number of individuals, sex, size, etc (**Barve**). When these two types of data are formatted in universally accepted formats, published openly and integrated with other such data streams, they are known "Digitally Accessible Knowledge" or DAK.

There are many sources of biodiversity data including: directed surveys, broad-scale surveys, and biological collections.

### 6.8.1 Direct Surveys

Directed surveys are often used where prior knowledge of a given system or biological mechanism exists. Surveys for this type of experiment are designed to control for known sources of variation(**GBIF**). Direct surveys aim to establish a causal relationship between some experimental treatment and its effect and are widely accepted in the scientific community. Regardless, the data collected in this type of survey is extremely expensive. Additionally, the data is aggregated by many researchers working independently in small areas over small time scales, which creates a network of heterogeneous data repositories with little opportunity for integration and eventual data degradation.

### 6.8.2 Broad Surveys

Broad-scale surveys involve over ten million observations annually. This data is used to generate probabilistic estimates of species occurrence, which can then be used to determine general patterns over broad geographic regions and time scales (**GBIF**). Data for broad surveys can be collected by anyone, ranging from trained observers to interested citizens. As a result, these surveys are relatively inexpensive to operate; however, the regulations behind these surveys are more relaxed, resulting in more potential for bias. Nonetheless, these studies provide the bulk of non-specimen observational data available via the Internet (**SP**).

### 6.8.3 Biological Collections

Biological collections are data sets of specimens. This type of data is primarily assembled to document the phenotypes and genotypes of the biotas worldwide (**Barve**). Biological collections are typically considered as "high-quality, low-volume" sources of data on biological diversity (**Barve**).

Given all the different ways to collect data, each with their own pros and cons, it is important for researchers utilizing biodiversity data to determine the completeness and accuracy of their data sets. Data visualization is one technique to observe patterns within very large data sets that are extremely difficult to analyze manually. The goal of this study is to learn how analyze biodiversity data sets and visualize these sets in order to paint a comprehensive representation of the data.

## 6.9 Objectives of the Case Study

The bdvis packages allows users to create models that visualize the biodiversity of the data collected. Using data that catalogs the occurrences of the species in question, bdvis creates a heat-map with geography superimposed for every point on the globe where the species was recognized, create a temporal visualization of the data by day, week, or month, and create a chronohorogram which plots the number of records on each day with colors indicating the value, and concentric circles for each year. The models that are created not only show the user if there is a lack of biodiversity, or lack of data, but also exactly where this lack occurs. The user can also see variations in the data on different months, and days. This variation in data may be due to differences in temperature.

By the end of this lesson, the reader should be able to:

1. Create a map of completeness of biodiversity
2. Visualize geographical and temporal coverage
3. Visualize gaps in data

## 6.10 Building the Model

Before loading you data into RStudio, it is necessary to install the packages: `maps`, `sqldf`, `plotrix`, `treemap`, `plyr`, `ggplot2`, `grid`, `lattice`, `chron`, `devtools`, `bdvis`, and

rinat. After installing these packages, load them into your library.

```
1 #make sure to include quotation marks when installing packages
2 install.packages("maps")
```

---

```
1 #leave off quotation marks when loading packages into library
2 library(maps)
```

---

Download the "Historical bird ringing records" from: <http://www.gbif.org/dataset/b4ae1720-1431-49ee-bfeb-8146fc42b1a3>. After data is configured to fit the bdvis package, set the working directory to the location of the data on your laptop. Next, read the data from the file, making sure to set stringsAsFactors as false.

```
1 #set working directory to where data is located
2 setwd("/Users/margauxwinter/Desktop/Rfolder/dwca-safring-v1.0")
3
4 #read your data into RStudio
5 birds <- read.delim('occurrence.txt', quote='', stringsAsFactors=FALSE)
```

---

You may begin using functions that create visual models of the data. The first model used in this example is a graph of completeness vs. number of species. To create this graph use the function bdcomplete, which divides the extent of the dataset in cells, via the getcellid function. Next, the function calculates the Chao2 estimator of species richness, which uses occurrence data from many samples to estimate the entirety of the species diversity. This function requires a minimum number of records that must be present on each cell to properly compute the index. If there are too few records in the cells, the function is unable to finish, an error will result. Use the following code to generate the plot shown in figure 6.3.

```
1 comp=bdcomplete(birds)
```

---

Next use the function mapgrid, which builds a grid map colored according to the density of records in each cell. Mapgrid allows the user to visually see gap and blank spaces in data. Grids are 1-degree cells, build with the getcellid function. Record-density maps apply a color gradient according to the number of records in the cell. Species-density maps apply a color gradient according to the number of different species in the cell, disregarding the number of records. Completeness maps apply a color gradient according to the completeness index, from 0 (incomplete) to 1 (complete). Use the following function to generate a map grid as shown in figure ??

```
1 mapgrid(comp, ptype='complete')
```

---

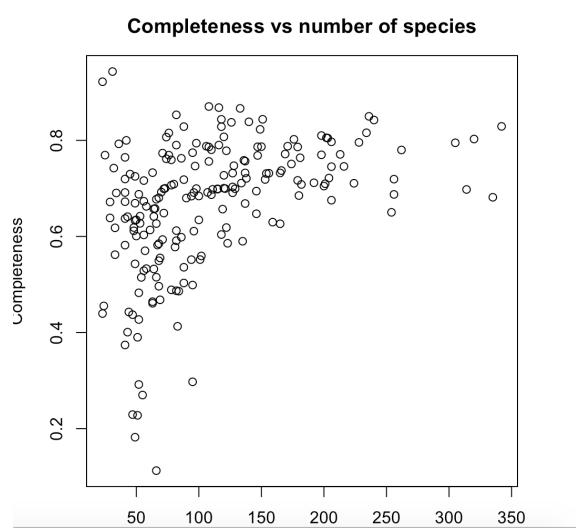


Figure 6.3: Completeness



Figure 6.4: Mapgrid

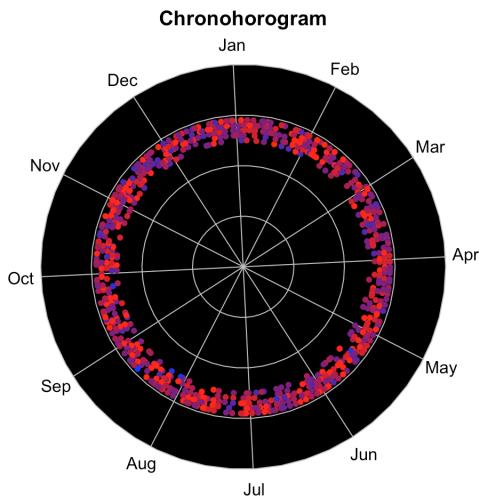
In order to create a chronohorogram of the dates in the provided data set, use the function `chronohorogram`. A chronohorogram plots the number of records on each day with colors indicating the value and concentric circles for each year. The following code is used to create a chronohorogram like the one shown in figure 6.5.

---

```
1 chronohorogram(birds)
```

---

Figure 6.5: Chronohorogram



Use the function `distrigraph` to build plots displaying distribution of biodiversity records among user-defined features. The next four distrographs all use different p-types (feature to represent) allowing each graphic to depict different aspects of the data. The following code is used to create all four distrograph respectively as shown in figures ?? and 6.7.

---

```
1 distrigraph(birds, ptype = "cell", col = "tomato")
2 distrigraph(birds, ptype = "species", ylab = "Species")
3 distrigraph(birds, ptype = "efforts", col = "red")
4 distrigraph(birds, ptype = "efforts", col = "red", type = "s")
```

---

`Bdcalendarheat` create a calendar heat map in which a calendar is structured as a matrix-like plot where each cell represents a unique date. The cells are colored as to show the density of records for that date. Use the following function to create the heat map shown in figure 6.8.

---

```
1 bdcalendarheat(birds)
```

---

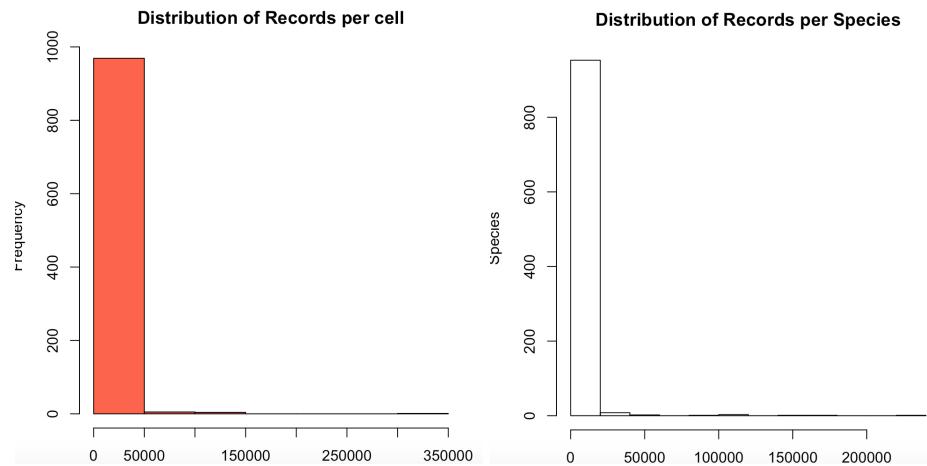


Figure 6.6: Left: Distrigraph 1, Right: Distrigraph 2

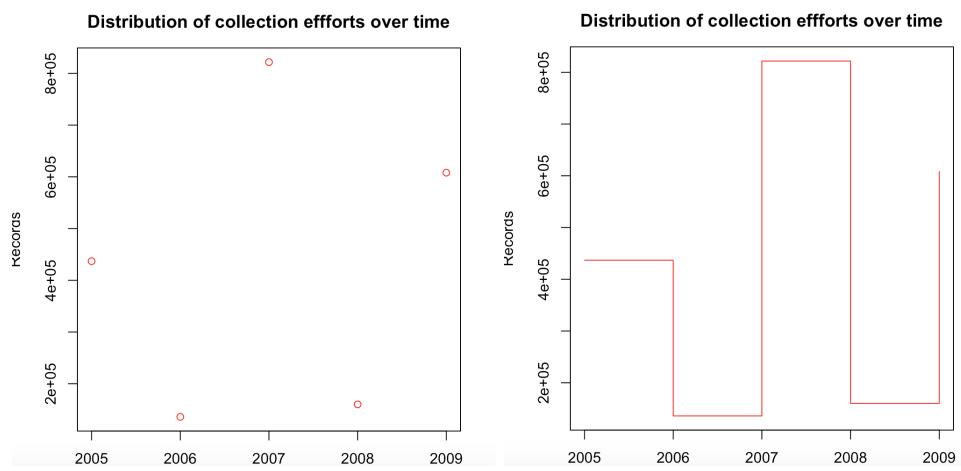


Figure 6.7: Left: Distrigraph 3, Right: Distrigraph 4

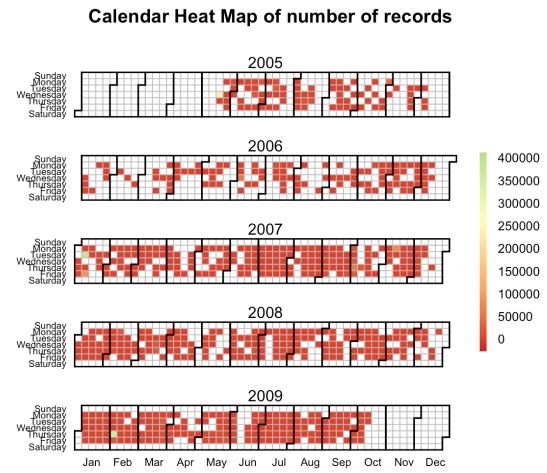


Figure 6.8: Heat Calendar

### 6.10.1 Programming Hints

Some functions, such as `mapgrid`, are customizable. Arguments within the function, including `indf`, `ptype`, `title`, `bbox`, `legscale`, `collow`, `colhigh`, `mapdatabase`, `region`, and `customize`, allow the map to be configured to fit the data as desired. Using the function `help()` for `mapgrid`, and other functions, within the console of RStudio gives a more comprehensive overview of these functions.

Other functions can be used within the package `bdvis`. These functions include `tempolar`, `gettaxo`, and `taxotree`. Each function provides a different view of the data by compiling data that are not used in the other models, and by allowing users to view the same data through different means.

## 6.11 Teaching Code

```

1
2 This is the code shown in the above examples:
3
4 #set working directory to data location
5 setwd("/Users/margauxwinter/Desktop/Rfolder/dwca-safring-v1.0")
6
7 #read data into RStudio
8 birds <- read.delim('occurrence.txt', quote='', stringsAsFactors=FALSE)
9
10 #configure data
11 conf <- list(Latitude='decimalLatitude', Longitude='decimalLongitude',
12 Date_collected='eventDate', Scientific_name='specificEpithet')
13 birds <- format_bdvis(birds, config=conf)
14 comp=bdcomplete(birds)
15
16 #visualize biodiversity data on the world map
17 mapgrid(comp, ptype='complete')
18
19 #create chronohorogram

```

---

```

20 chronohorogram(birds)
21 comp=bdcomplete(birds, recs=5)
22
23 #create various distigraphs
24 distigraph(birds, ptype="cell", col="tomato")
25 distigraph(birds, ptype="species", ylab="Species")
26 distigraph(birds, ptype="efforts", col="red")
27 distigraph(birds, ptype="efforts", col="red", type="s")
28
29 #create heat calendar
30 bdcalendarheat(birds)

```

---

## 6.12 Example Student Code

The final deliverable should be submitted as visual models of the data in the form of distrographs, maps, chronohorograms, and plots. Use the arachnid data set downloaded from: <http://www.gbif.org/occurrence/download/0012867-160910150852091>

---

```

1 #set working directory
2 setwd("/Users/margauxwinter/Desktop/Rfolder/
3 dwca-arc_national_survey_of_arachnida-v1.1")
4
5 #read in data set
6 spider <- read.delim('occurrence.txt', quote='', stringsAsFactors=FALSE)
7
8 #configure data to make it readable for bdvis package
9 conf <- list(Latitude='decimalLatitude', Longitude='decimalLongitude',
10 Date_collected='eventDate', Scientific_name='specificEpithet')
11 spider <- format_bdvis(spider, config=conf)
12 comp=bdcomplete(spider)
13
14 #begin functions for data visualization
15 mapgrid(comp, ptype='complete')
16 tempolar(spider, color="green", title="Arachnida daily", plottype="r",
17 timescale="d")
17 tempolar(spider, color="blue", title="Arachnida weekly", plottype="p",
18 timescale="w")
18 tempolar(spider, color="red", title="Arachnida monthly", plottype="r",
19 timescale="m")
19 chronohorogram(spider)
20 comp=bdcomplete(spider, recs=5)
21 distigraph(spider, ptype="cell", col="tomato")
22 distigraph(spider, ptype="species", ylab="Species")
23 distigraph(spider, ptype="efforts", col="red")
24 distigraph(spider, ptype="efforts", col="red", type="s")

```

---

- (1) Kelling, S. 2008. Significance of organism observations: Data discovery and access in biodiversity research. Report for the Global Biodiversity Information Facility, Copenhagen. 20 pp. Available online at [http://www.gbif.org/orc/?doc\\_id=1642](http://www.gbif.org/orc/?doc_id=1642).
- (2) Shanmughavel, P. (2007). An overview on biodiversity information in databases. *Bioinformation*, 1(9), 367–9. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/17597923>
- (3) Barve, V. (2014). Discovering and developing primary biodiversity data from social networking sites: A novel approach. *Ecological Informatics*, 24, 194–199. <http://doi.org/10.1016/j.ecoinf.2014.08.008>
- (4) RDocumentation. Web. 09 Oct. 2016. "Search All 11,558 CRAN, BioConductor and Github Packages." <https://www.rdocumentation.org>
- (5) Vijay Barve. Web. 09 Oct. 2016. "Vijay Barve." <https://vijaybarve.wordpress.com/tag/bdvis/>
- (6) Rbloggers. 08 May 2014. Web. 09 Oct. 2016. "Package Bdvis Is on CRAN." <https://www.r-bloggers.com/package-bdvis-is-on-cran/>
- (7) Accessing INaturalist Data. Web. 09 Oct. 2016. By Checking the Dimensions, We Can See Only 47 Records Were Found. We Could Try the Samething for a given Day, Month or Year. Let's Try Searhing for Cumulative Totals of Observations of Ephemeroptera and See If We Can Detect Seasonality. "Accessing INaturalist Data." <https://ropensci.org/blog/2014/03/26/rinat/>

- (8) GitHub. Web. 09 Oct. 2016. Vijaybarve. "Vijaybarve/bdvis." <https://github.com/vijaybarve/bdvis>
- (9) SAFRING: Historical Bird Ringing Records (2005-2009). Web. 09 Oct. 2016. Gbif. "SAFRING: Historical Bird Ringing Rec..." <http://www.gbif.org/dataset/b4ae1720-1431-49ee-bfeb-8146fc42b1a3>