Andrew Hautau
Intro to Data Science
Week 3 Assignment

1. K-Means Clustering
    a. Data normalization is important for K-means clustering because it puts the data on equal terms.
    b. K-means methods necessarily rely on numerical data, since the method relies on finding the means of data points.
    c. Clustering is un-supervised learning because the task focuses on learning the structure of a given dataset. The task does not involve inference of a desired function from training data, as in supervised learning
2. R Code and Mammographic Masses Dataset:

```r
1  # Andrew Hautau
2  # Intro To Data Science
3  # Class 3 Assignment
4
5  # Fetch csv file.
6  df <- read.csv("http://archive.ics.uci.edu/ml/machine-learning-databases/mammographic-masses/mammographic_masses.data")
7
8  # Give Column Names
9  names(df) <- c("BI-RADS assessment","Age","Shape","Margin","Density","Severity")
10 names(df)
11
12 # Coerce all values of data frame into numeric values.
13 df <- data.frame(lapply(df, function(x) as.numeric(as.character(x))))
14
15 # How many NAs are present?
16 stuff <- is.na(df)
17 # 162, apparently.
18 sum(stuff)
19 # How about by column?
20 colSums(stuff)
21
22 # Let's fetch the max and mins for each column, and then histogram those vectors!
23 maxs <- c()
24 mins <- c()
25
26 for(i in seq_along(df)) {
27    maxs <- cbind(maxs, max(df[,i], na.rm=T))
28    mins <- cbind(mins, min(df[,i], na.rm=T))
29 }
30
31 # 96 years is an outlier in this dataset.
32 maxs
33 hist(df$Age)
34
35 # Find the index of 96 in column 1.
36 which(df$Age == 96)
37 # Let's get rid of the observation.
38 df[726,] = NULL
39 df[726,2]
40
41 # Changed Severity to Diseased.
42 colnames(df)[6] <- "Diseased"
43 colnames(df)[6]
44
45 # Write to csv file.
46 write.csv(df,"mammographic_mass.csv")
47
48
```

3.  Predixion Insight: After attempting into the wee-hours, I'm putting up the white flag on this one. I setup a Windows VM with Excel (also with PowerPivot), but the Predixion installation wizard would always blow up, complaining about not having an install of SQL Server. This all seems to clash with the comments on the LinkedIn group. My one theory is that SQL Server is a dependency for Insight with Excel 2010, but not other versions.
4.  Playing with Octave (following tutorials from (http://en.wikibooks.org/wiki/Octave_Programming_Tutorial).

```
Octave was configured for "x86_64-apple-darwin10.7.3".

Additional information about Octave is available at http://www.octave.org.

Please contribute if you find this software useful.
For more information, visit http://www.octave.org/help-wanted.html

Read http://www.octave.org/bugs.html to learn how to submit bug reports.

For information about changes from previous versions, type `news'.

octave-3.4.0:1> x = [1, 3, 2]
x =

   1   3   2

octave-3.4.0:2>  x = [1; 3; 2]
x =

   1
   3
   2

octave-3.4.0:3>  A = [1, 1, 2; 3, 5, 8; 13, 21, 34]
A =

    1    1    2
    3    5    8
   13   21   34

octave-3.4.0:4> A = [1, 6, 3; 2, 7, 4]
A =

   1   6   3
   2   7   4

octave-3.4.0:5> B = [2, 7, 2; 7, 3, 9]
B =

   2   7   2
   7   3   9

octave-3.4.0:6> A ./ B
ans =

   0.50000   0.85714   1.50000
   0.28571   2.33333   0.44444

octave-3.4.0:7>
```

```
octave-3.4.0:14> A = [1, 2, 3; 4, 5, 6; 7, 8, 9]
A =

   1   2   3
   4   5   6
   7   8   9

octave-3.4.0:15> diag(x)
ans =

Diagonal Matrix

   1   0   0
   0   5   0
   0   0   9

octave-3.4.0:16> fliplr(A)
ans =

   3   2   1
   6   5   4
   9   8   7

octave-3.4.0:17> flipud(A)
ans =

   7   8   9
   4   5   6
   1   2   3

octave-3.4.0:18> 
```