

Chapter 01: Transitioning from Data Developer to Data Scientist

	New Data	New Data Source
Data Developer	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Data Scientist		<input checked="" type="checkbox"/>

	Data Quality
Data Developer	<input checked="" type="checkbox"/>
Data Scientist	<input checked="" type="checkbox"/>

	Query	Mining
Data Developer	<input checked="" type="checkbox"/>	
Data Scientist	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

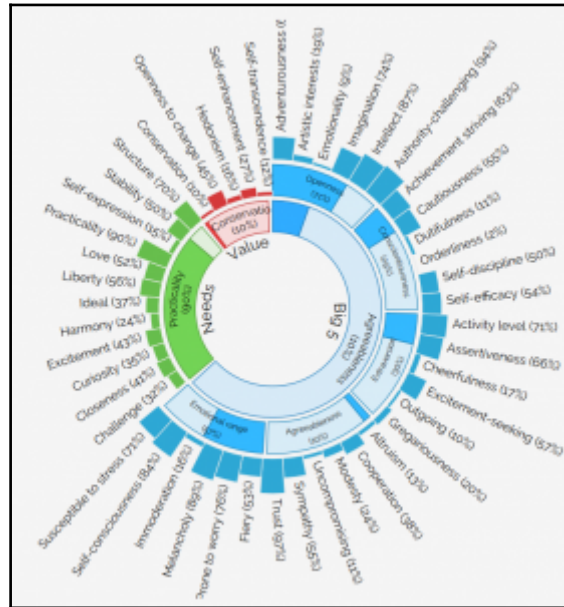
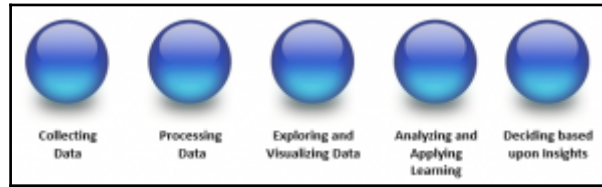
	Performance
Data Developer	<input checked="" type="checkbox"/>
Data Scientist	<input checked="" type="checkbox"/>

Financial Reporting	
Data Developer	<input checked="" type="checkbox"/>
Data Scientist	

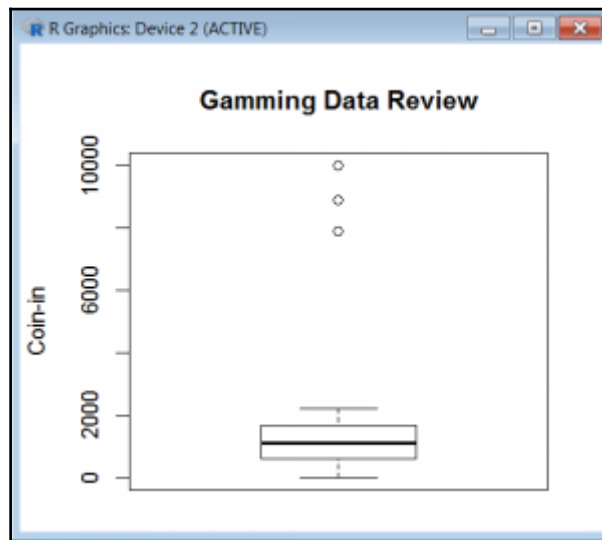
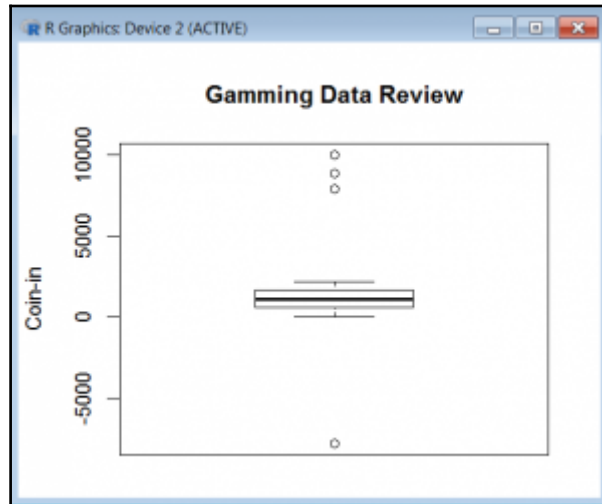
Visualization	
Data Developer	<input checked="" type="checkbox"/>
Data Scientist	<input checked="" type="checkbox"/>

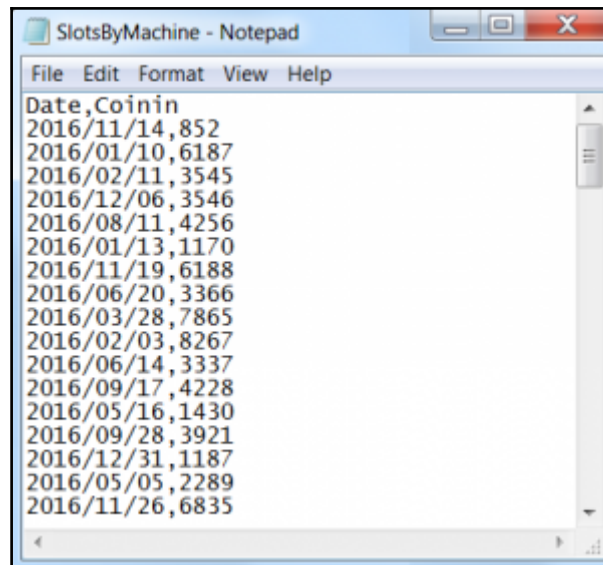
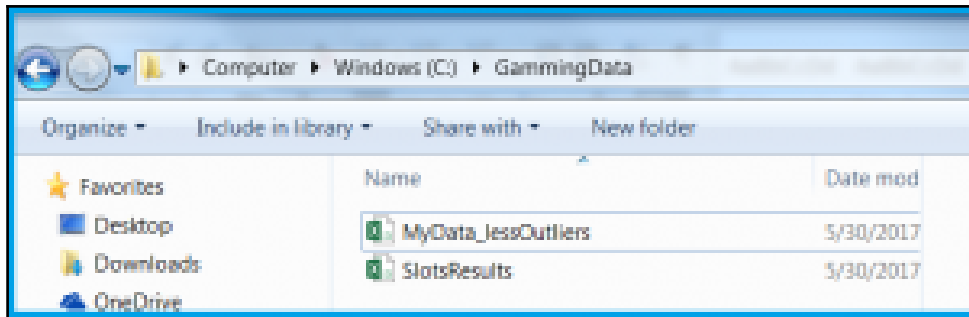
Tools of the Trade	
Data Developer	<input checked="" type="checkbox"/>
Data Scientist	<input checked="" type="checkbox"/>

Chapter 02: Declaring the Objectives



Chapter 03: A Developer's Approach to Data Cleaning

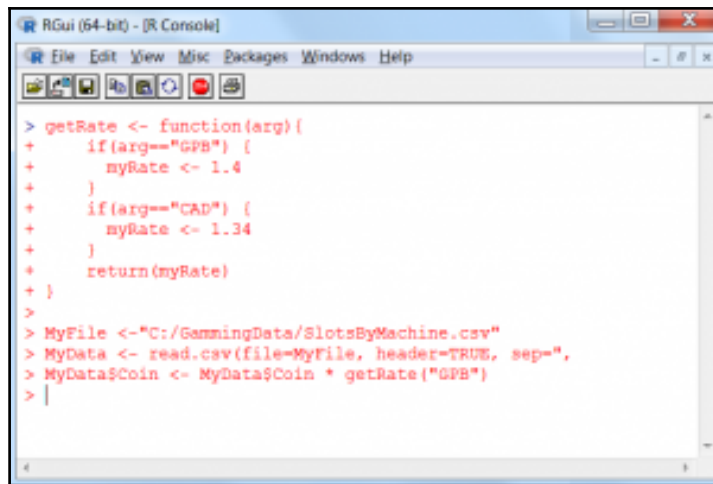




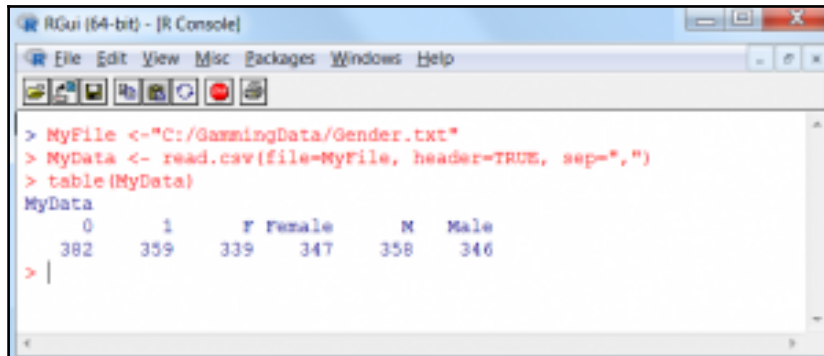
```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help
> MyFile <- "C:/GamingData/SlotsByMachine.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, se
> class(MyData)
[1] "data.frame"
> class(MyData$Date)
[1] "factor"
> class(MyData$Coinin)
[1] "integer"
> |
```

```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help
> MyFile <- "C:/GamingData/SlotsByMachine.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> MyData$Date <- parseDate(MyData$Date, "%m/%d/%Y", return=MyData$Date, %Y, sep="")
> class(MyData$Date)
[1] "Date"
> |
```

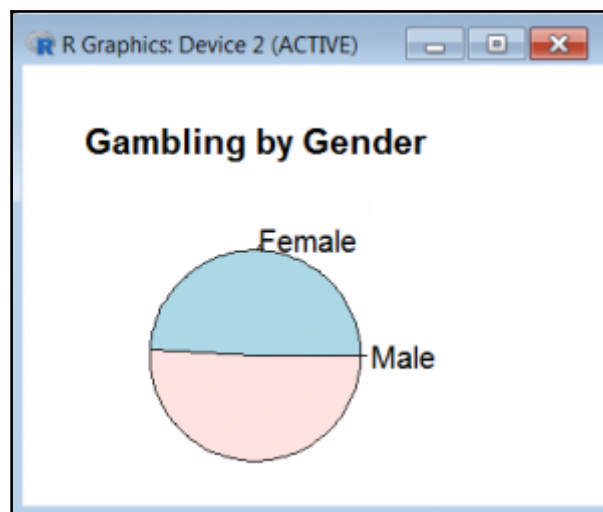
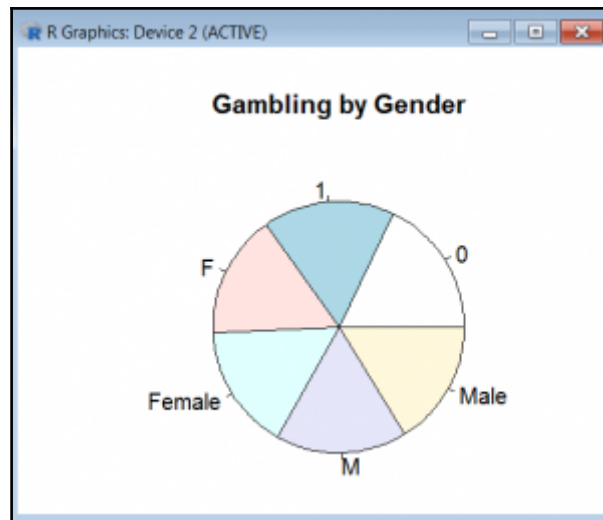
```
RGui (64-bit) - [R Console]
File Edit View Misc Packages Windows Help
> MyFile <- "C:/GamingData/SlotsByMachine.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> MyData$Date <- parseDate(MyData$Date, "%m/%d/%Y", return=MyData$Date, %Y, sep="")
> class(MyData$Date)
[1] "Date"
> class(MyData$Date)
[1] "Date"
> |
```



```
> getRate <- function(arg){
+   if(arg=="GBP") {
+     myRate <- 1.4
+   }
+   if(arg=="CAD") {
+     myRate <- 1.34
+   }
+   return(myRate)
+ }
>
> MyFile <- "C:/GamingData/SlotsByMachine.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> MyData$Coin <- MyData$Coin * getRate("GBP")
> |
```



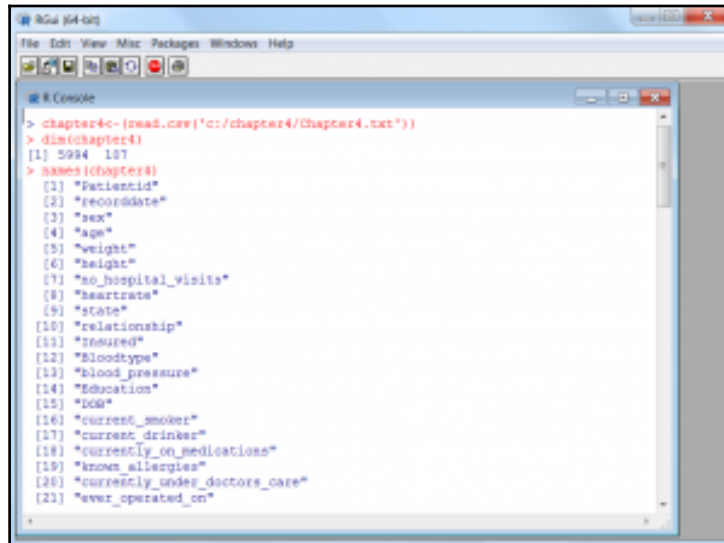
```
> MyFile <- "C:/GamingData/Gender.txt"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> table(MyData)
MyData
  0      1      F Female      M  Male
382  359  339  347  358  346
> |
```




```
R Console
> MyFile <- "C:/GamingData/SlotsResults.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> MyData[[1]]
  Coin.in
1      9999
2      1505
3       673
4       872
5      7900
6       671
7       475
8      1750
9      8888
10       88
11       790
12     -7792
13      1030
14       353
15       823
16      1789
17      1648
```

```
R Console
> MyFile <- "C:/GamingData/SlotsResults.csv"
> MyData <- read.csv(file=MyFile, header=TRUE, sep=",")
> scale(MyData[[1]], center = TRUE, scale = TRUE)
  Coin.in
[1,] 1.200796e+01
[2,] 5.190288e-01
[3,] -6.063296e-01
[4,] -3.371633e-01
[5,] 9.168869e+00
[6,] -6.090348e-01
[7,] -8.741432e-01
[8,] 8.504144e-01
[9,] 1.050523e+01
[10,] -1.397597e+00
[11,] -4.480760e-01
[12,] -1.205604e+01
[13,] -1.234534e-01
[14,] -1.039160e+00
[15,] -4.034404e-01
[16,] 9.031655e-01
[17,] 7.124498e-01
```

Chapter 04: Data Mining and the Database Developer



```
RGui [64-bit]
File Edit View Misc Packages Windows Help

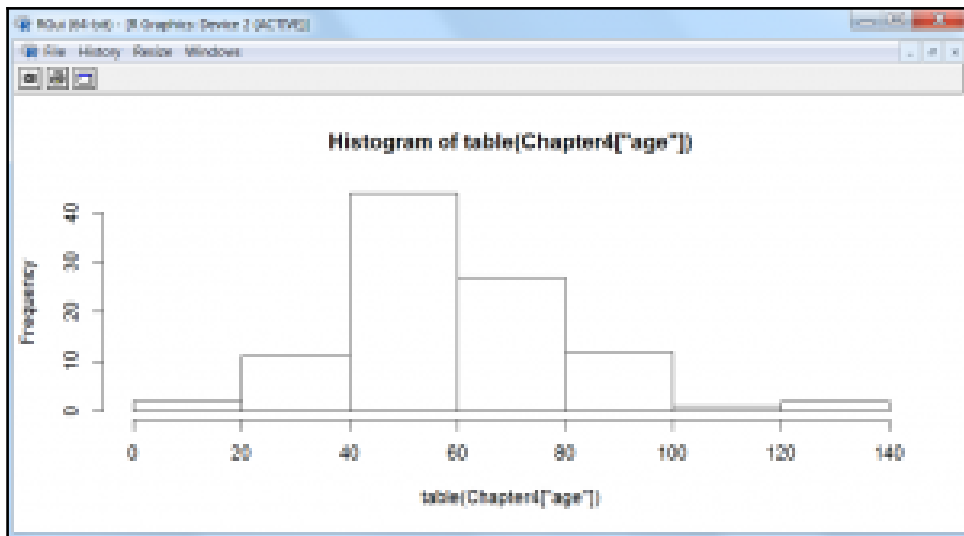
Console
> chapter4<-read.csv('c:/chapter4/Chapter4.txt')
> dim(chapter4)
[1] 5994 187
> names(chapter4)
 [1] "patientid"
 [2] "recorddate"
 [3] "sex"
 [4] "age"
 [5] "weight"
 [6] "height"
 [7] "no_hospital_visits"
 [8] "heartrate"
 [9] "state"
[10] "relationship"
[11] "insured"
[12] "bloodtype"
[13] "blood_pressure"
[14] "Education"
[15] "DOB"
[16] "current_smoker"
[17] "current_drinker"
[18] "currently_on_medications"
[19] "known_allergies"
[20] "currently_under_doctors_care"
[21] "ever_operated_on"
```

```
> table(chapter4["current_smoker"])

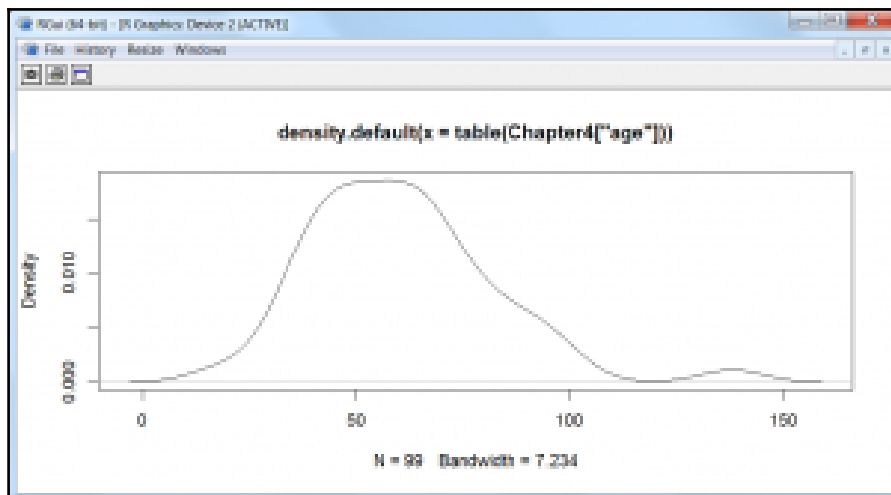
No Yes
5466 528
> |
```

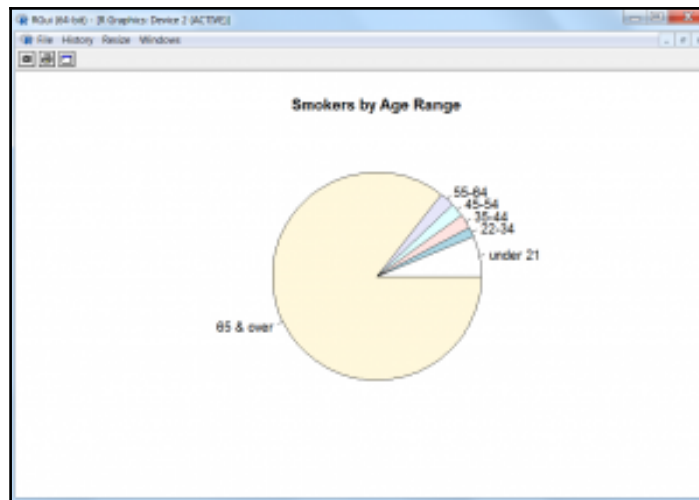
```
> range(Chapter4["age"])
[1] 1 99
> |
```

```
> Chapter4<-read.csv('c:/chapter4/Chapter4.txt')
> hist(table(Chapter4["age"]))
> |
```



```
> plot(density(table(Chapter4["age"])))  
> |
```





```
> Chapter4<-read.csv('c:/chapter4/Chapter4_NA.csv')
> nrow(Chapter4)
[1] 5994
> Chapter4<-na.omit(Chapter4)
> nrow(Chapter4)
[1] 5988
> |
```

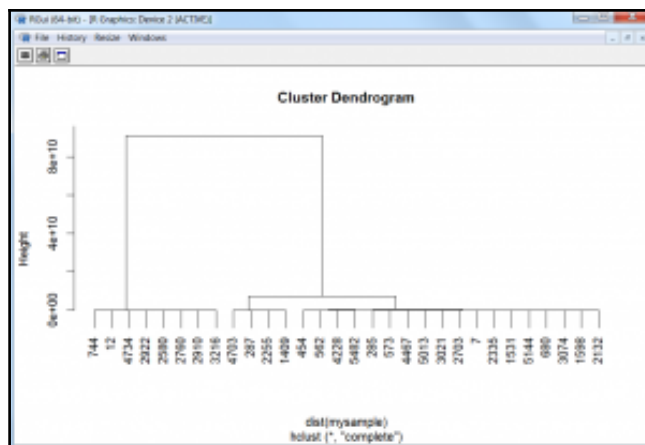
```
> # --- create a subset of smokers only
> mysub <- subset(Chapter4,Chapter4["current_smoker"]=="Yes")
> # --- confirm row count
> nrow(mysub)
[1] 528
> |
```

```
> # --- create a random sample of 30 smokers
> mysample <- mysub[sample(1:nrow(mysub), 30,
+   replace=FALSE),]
> # --- confirm the row count in our sample
> nrow(mysample)
[1] 30
> |
```

```

> # -- perform hierarchical cluster
> smokerclust<-hclust(dist(mysample))
> # -- create results in a dendrogram
> plot(smokerclust)
> |

```



```

> Chapter4<-read.csv('c:/chapter4/varainces.csv')
> var(Chapter4["No_servings_per_week_skin_milk"])
No_servings_per_week_skin_milk      No_servings_per_week_skin_milk
No_servings_per_week_skin_milk      0.003160316
> |

```

```

> Chapter4<-read.csv('c:/chapter4/varainces.csv')
> var(Chapter4["No_servings_per_week_skin_milk"])
No_servings_per_week_skin_milk      No_servings_per_week_skin_milk
No_servings_per_week_skin_milk      0.003160316
> var(Chapter4["No_servings_per_week_regular_or_diet_soda"])
No_servings_per_week_regular_or_diet_soda      No_servings_per_week_regular_or_diet_soda
No_servings_per_week_regular_or_diet_soda      8.505655
> |

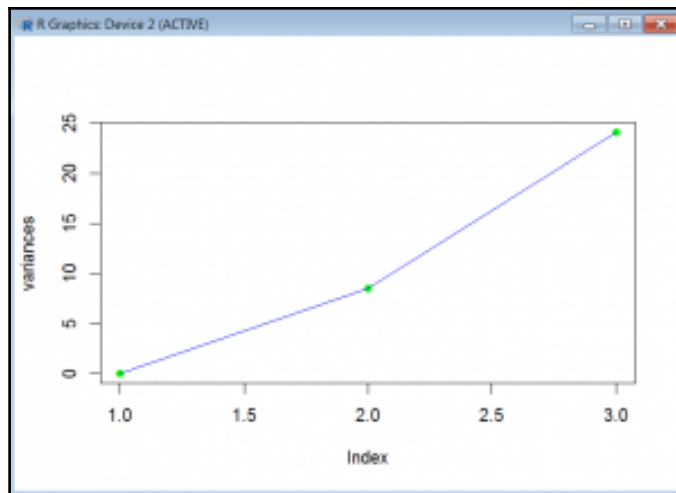
```

```

> var(Chapter4["No_servings_per_week_water"])
No_servings_per_week_water      No_servings_per_week_water
No_servings_per_week_water      24.10477
> |

```

```
> Chapter4<-read.csv('c:/chapter4/varainces.csv')
> c1<-var(Chapter4["No_servings_per_week_skin_milk"])
> c2<-var(Chapter4["No_servings_per_week_regular_or_diet_soda"])
> c3<-var(Chapter4["No_servings_per_week_water"])
> variances<-c(c1, c2, c3)
> plot(variances, pch=16, col="green")
> lines(variances, col="blue")
> |
```

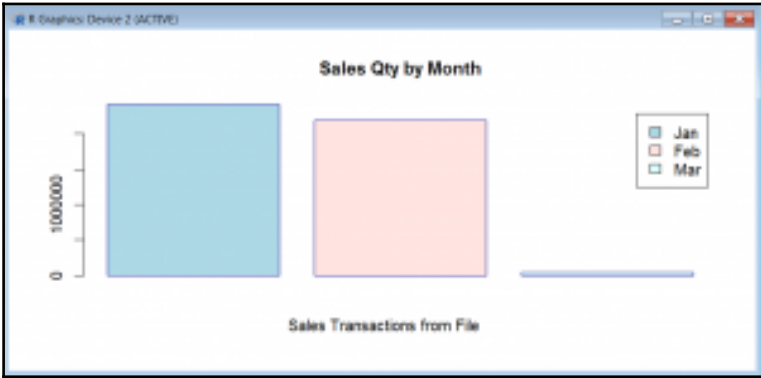


Chapter 05: Statistical Analysis for the Database Developer

```
> myData <- read.csv("01/Worksp/ExampleSalesTrans.csv")
> colnames(myData)
[1] "product_code" "product_name" "quantity" "sales_date" "return_date" "sales_region"
> summary(myData)
  product_code  product_name  quantity      sales_date    return_date  sales_region
min.   : 1.0    Bell       : 39   min.   : 3          2129          2911   min.   : 1.00
1st Qu.:25.0    Back       : 39   1st Qu.:2449    2/14/2013: 25   1/22/2013: 9   1st Qu.:2.00
Median :51.0    Head badge: 38   Median :5840    1/17/2013: 22   1/9/2013 : 8   Median :3.00
Mean   :59.4    Locknut    : 37   Mean   :5873    1/18/2013: 22   2/22/2013: 5   Mean   :2.99
3rd Qu.:75.0    Cable     : 16   3rd Qu.:7853    1/9/2013 : 22   2/28/2013: 5   3rd Qu.:4.00
Max.   :99.0    Chainguard: 16   Max.   :8883    2/2/2013 : 20   1/2/2013 : 4   Max.   :5.00
      (Other) :835      (Other) :860      (Other) :101
> |
```

```
> nrow(myData)
[1] 1040
> list(unique(myData$product_name))
[[1]]
[1] Freehub
[2] Seat lug
[3] Handlebar plug
[4] Handlebar tape
[5] Locknut
[6] Kickstand
```

```
> # --- use sort and unique functions to list our year(s) and month(s)
>
> sort(unique(YearsInData))
[1] "" "2013"
> sort(unique(MonthsInData))
[1] "" "1" "2" "3"
> |
```



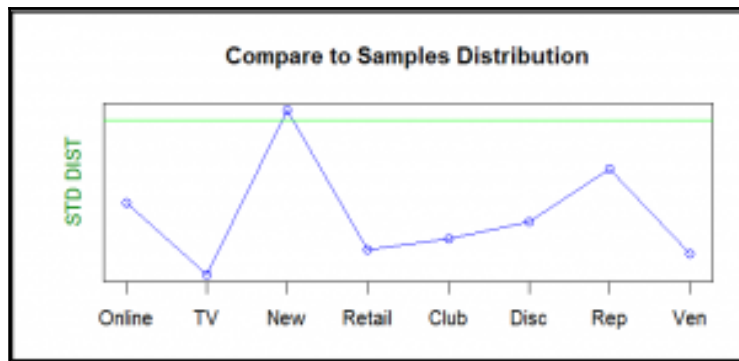
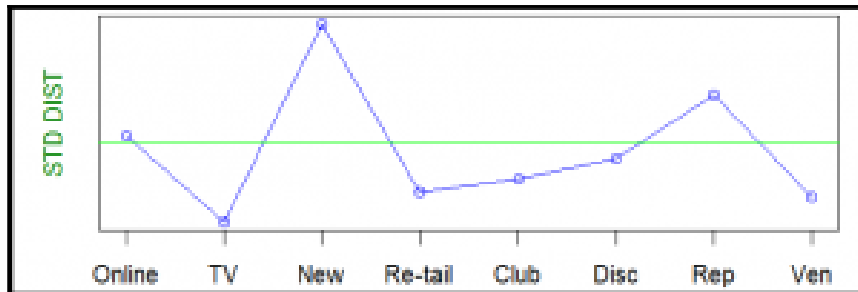
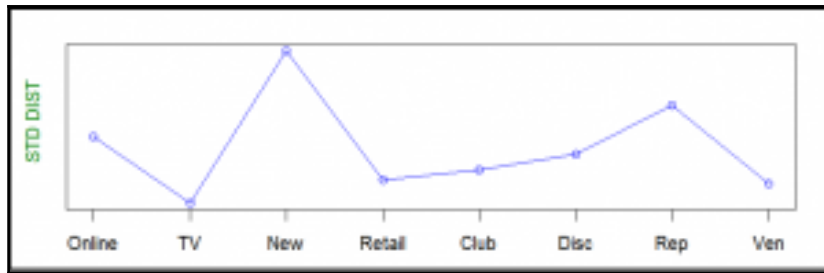
```
> # read.csv("c:/Mocker/Samples/SalesTrans_2.csv", sep=";", header = TRUE, skip = 0)
> colnames(x)
[1] "product_code" "product_name" "quantity" "sales_date" "return_date" "sales_region" "sale_type"
> |
```

```
> summary(x)
product_code
Min.   : 1.00
1st Qu.:29.00
Median :58.00
Mean   :49.50
3rd Qu.:74.00
Max.   :99.00

Subset
Ricycle bike case
Qty
Description
1) the disc component of a disc brake.[12] another name for a detangler - a device that allows the handlebars and
bar ends
RObject
  quantity    sales_date    return_date    sales_region    sale_type
Min.   : 18             : 279             :1960   Min.   :1.000   Vendor   :294
1st Qu.:2590   1/22/2014 : 27   2/2/2014 : 7   1st Qu.:12.000   Club     :288
Median :3222   2/27/2012 : 27   1/7/2014 : 6   Median :1.000   Repeat   :287
Mean   :3044   2/12/2014 : 26   2/14/2014 : 6   Mean   :2.989   Online   :285
3rd Qu.:7480   1/5/2014 : 24   2/15/2014 : 6   3rd Qu.:4.000   Retailer :284
Max.   :9932   2/2/2013 : 24   2/17/2014 : 6   Max.   :9.000   New Customer:274
      (Other) :1832 (Other) : 248 (Other) :523
```

```
> # --- mean for all sale types:
> MeanAll <-mean(data.df[["quantity"]])
>
> # --- standard deviation for all sales types:
> StdDAll<-sd(data.df[["quantity"]])
>
> # --- median for all sales types:
> MeanAll <-mean(data.df[["quantity"]])
> MeanAll
[1] 5043.971
> StdDAll
[1] 2876.796
> MeanAll
[1] 5043.971
> |
```

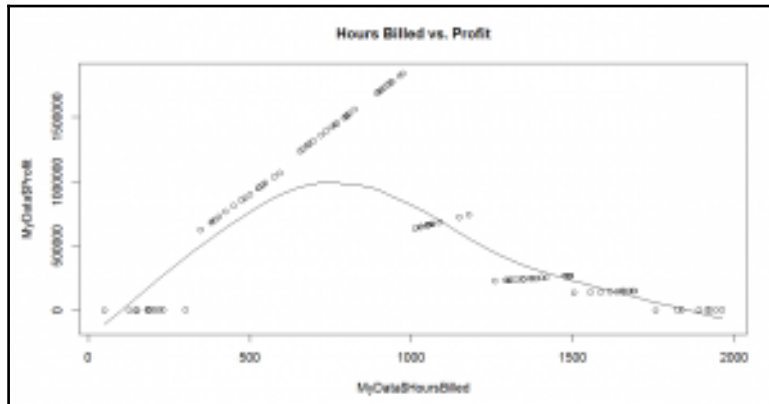
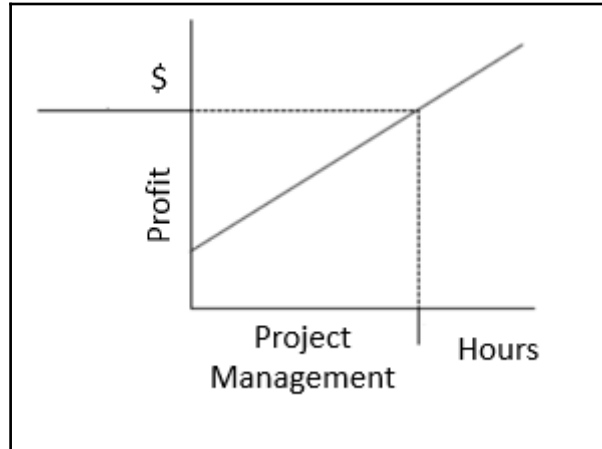


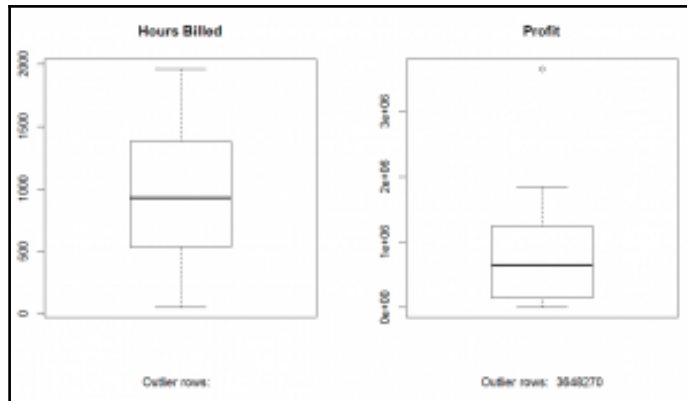


```
> df
  group max    mean    sd median min    sum
1  Online 9982 5160.888 2883.941 5100.0  33 1470853
2    TV 9970 4516.475 2803.657 4305.0  30 1151701
3 New Customer 9954 5052.507 2986.567 4966.5 122 1394492
4  Retailer 9986 4972.775 2832.253 4862.0  10 1412268
5    Club 9980 5282.920 2843.501 5560.5  30 1521481
6 Discounted 9918 5408.142 2862.614 5622.0  10 1449382
7  Repeat 9950 4937.014 2920.933 5068.0  90 1416923
8  Vendor 9848 4987.676 2827.729 5137.0  40 1476352
> |
```



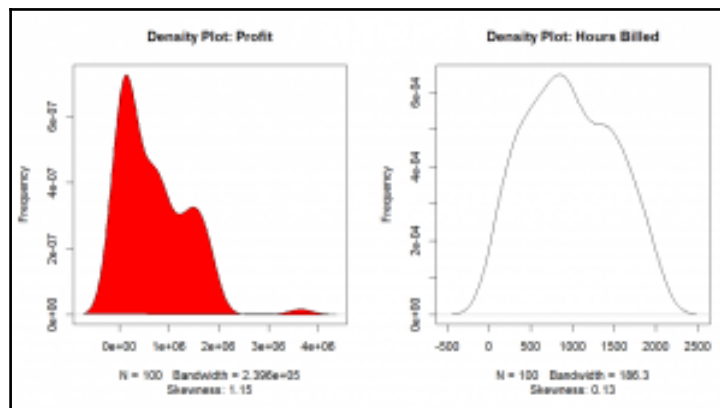
Chapter 06: Database Progression to Database Regression





```

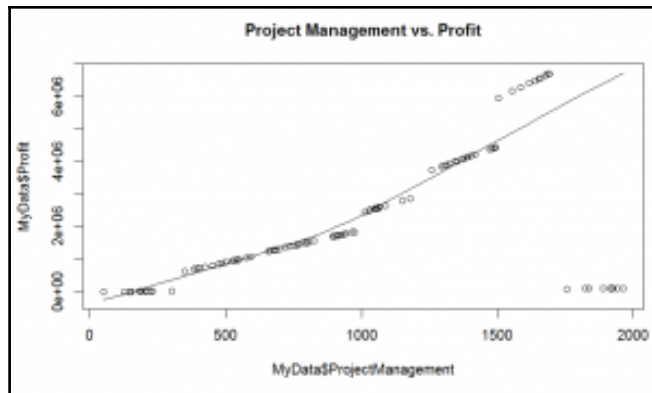
HoursBilledProfit - Notepad
File Edit Format View Help
ProjectID, hoursBilled, Profit
1,1029.00,9648270.00
2,528.00,948800.00
3,1824.00,820.80
4,807.00,1525230.00
5,1737.00,790.65
6,184.00,1821.60
7,1616.00,145440.00
8,1915.00,861.75
9,696.00,1315440.00
10,1837.00,147330.00
11,1368.00,246240.00
12,1303.00,234540.00
13,663.00,1253070.00
14,1963.00,883.35
15,404.00,727200.00
16,1312.00,236160.00
17,737.00,1392930.00
18,764.00,1443960.00
19,1943.00,874.35
    
```



```

> cor(MyData$HoursBilled, MyData$Profit)
[1] -0.2578628
> |
    
```

```
> MyData <- read.csv(file="c:/Worker/ProjectManagementProfit0.csv", header=TRUE, sep=",")
> cor(MyData$ProjectManagement, MyData$Profit)
[1] 0.6129142
> |
```



```
> # --- build linear regression model on full data
> aLinearMod <- lm(ProjectManagement ~ Profit, data=MyData)
> print(aLinearMod)

Call:
lm(formula = ProjectManagement ~ Profit, data = MyData)

Coefficients:
(Intercept)      Profit
      6.180         1.629

> |
```

```
> # --- predict project profitability
> ProfitPred <- predict(linMod, testData)
> ProfitPred
      1      3      4      5      7      14      23      25      33      41
1098.0782 561.3375 833.4988 980.7987 1490.0242 282.5688 787.6671 707.8238 945.5875 1297.7181
  44  59  80  80  80  78  82  88  98  100
1342.8028 889.3103 849.7438 783.3471 1029.6436 1883.9952 792.6667 1000.2492 894.8933 732.8004

> |
```

```
> summary(lmMod)

Call:
lm(formula = ProjectManagement ~ Profit, data = trainingData)

Residuals:
    Min       1Q   Median       3Q      Max
-516.48 -210.30  -43.53   59.69 1360.63

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.655e+02  6.448e+01  8.763 3.03e-13 ***
Profit      1.757e-04  2.136e-05  8.229 3.39e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 383.4 on 78 degrees of freedom
Multiple R-squared:  0.4647,    Adjusted R-squared:  0.4578
F-statistic: 67.71 on 1 and 78 DF, p-value: 3.393e-12

> |
```

```
> actuals_preds <- data.frame(ohio$actuals=testData$ProjectManagement, predicted=distPred)
> correlation_accuracy <- cor(actuals_preds)
> head(actuals_preds)
  actuals predicted
1    1029    1004.8763
2    1024     981.3375
4     807     833.4984
5    1737     980.7547
7    1816    1690.8242
14   1963     883.5466

> |
```

Chapter 07: Regularization for Database Improvement

```
> ols <- lm(y~ x1 + x2 + x3)
> summary(ols)

Call:
lm(formula = y ~ x1 + x2 + x3)

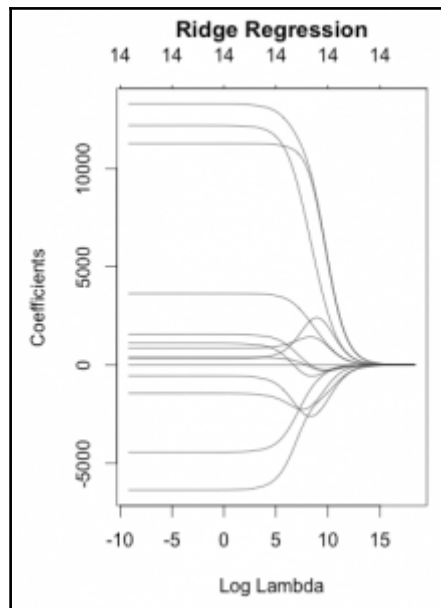
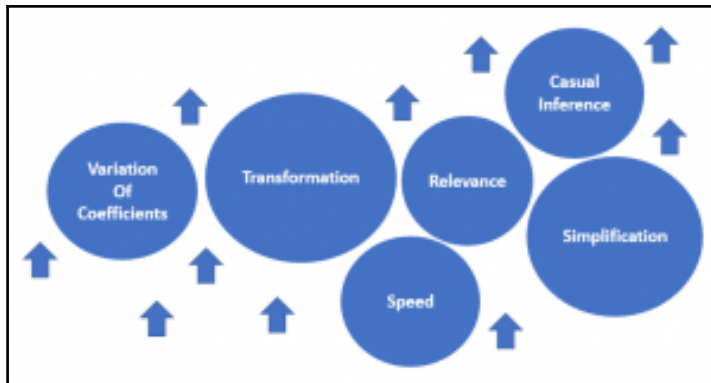
Residuals:
    Min       1Q   Median       3Q      Max
-1.19698 -0.28592  0.04026  0.24016  1.20322

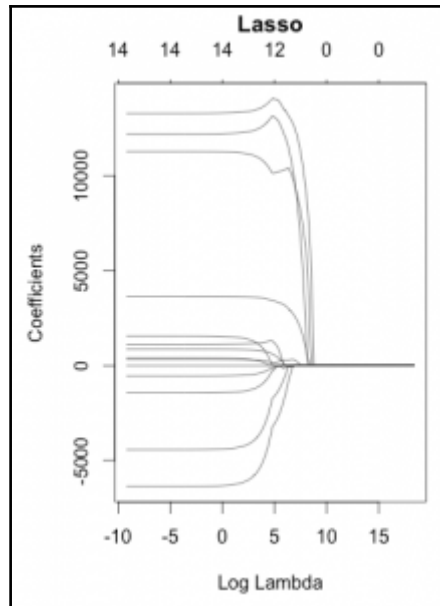
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -0.4293    0.4916  -0.873   0.3954
x1             1.7851    0.4812   3.710  0.0019 **
x2             0.7119    0.4622   1.540  0.1430
x3             0.2839    0.5122   0.554  0.5870
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6306 on 16 degrees of freedom
Multiple R-squared:  0.4831,    Adjusted R-squared:  0.3862
F-statistic: 4.984 on 3 and 16 DF,  p-value: 0.0125

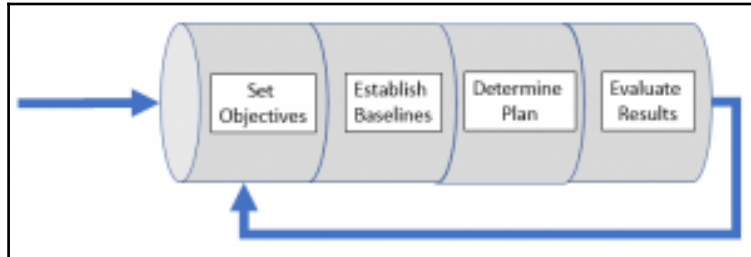
> |
```

```
> # --- Fit model using ridge regression using independent variables
>
> ridge <- lm.ridge (y ~ x1 + x2 + x3, lambda = seq(0, .1, .001))
> summary(ridge)
      Length Class Mode
coef    303  -none- numeric
scales    3  -none- numeric
Inter     1  -none- numeric
lambda  101  -none- numeric
ym         1  -none- numeric
xm         3  -none- numeric
GCV      101  -none- numeric
kHKB      1  -none- numeric
kLW       1  -none- numeric
> |
```





Chapter 08: Database Development and Assessment



```
> ols <- lm(y~ x1 + x2 + x3)
> summary(ols)

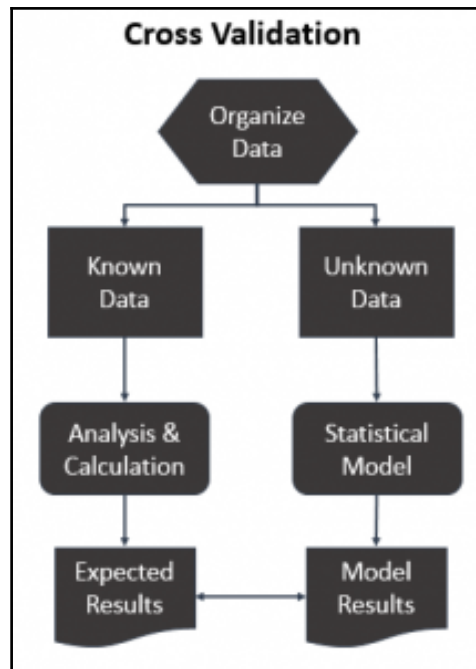
Call:
lm(formula = y ~ x1 + x2 + x3)

Residuals:
    Min       1Q   Median       3Q      Max
-1.19698 -0.28592  0.04026  0.24016  1.20322

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -0.4293    0.4916   -0.873   0.3954
x1             1.7851    0.4812    3.710   0.0019 **
x2             0.7119    0.4622    1.540   0.1430
x3             0.2839    0.5122    0.554   0.5870
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

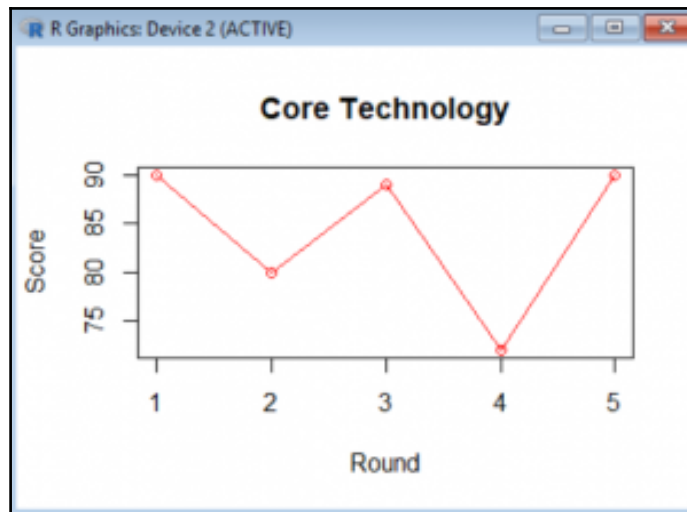
Residual standard error: 0.6306 on 16 degrees of freedom
Multiple R-squared:  0.4831,    Adjusted R-squared:  0.3862
F-statistic: 4.984 on 3 and 16 DF,  p-value: 0.0125

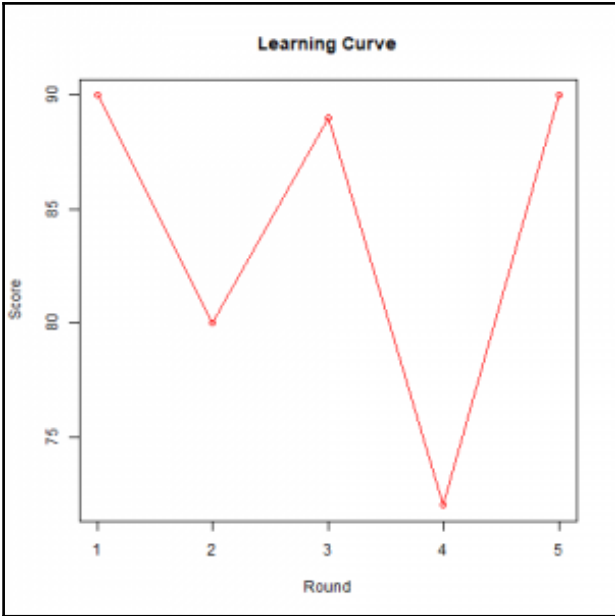
> |
```



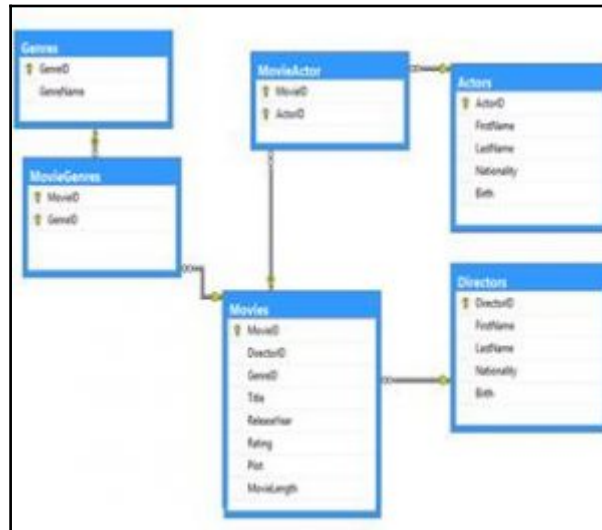
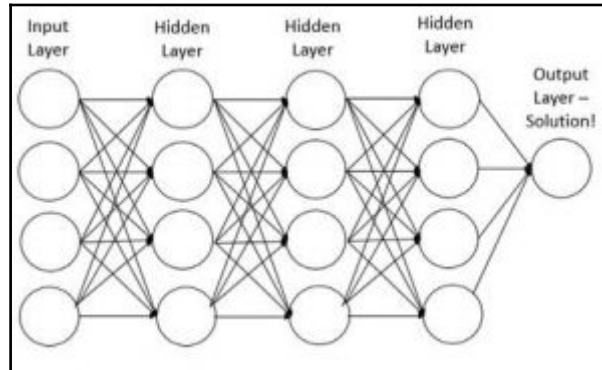
Characteristic	Validation Percent	Round 1 Percent	Round 2 Percent	Round 3 Percent	Round 4 Percent	Round 5 Percent	Average
Core Technology	90%	90%	80%	89%	72%	90%	85%
FT Project Manager	17%	30%	17%	30%	15%	17%	15%
FT Client Resource	79%	78%	77%	79%	80%	78%	79%
Sub-Contracted	51%	9%	5%	15%	44%	79%	41%
Time & Materials	65%	99%	99%	89%	75%	99%	89%
Not to Exceed	90%	69%	69%	69%	59%	70%	71%
Formal QA	95%	89%	85%	91%	92%	99%	93%
On-site	75%	78%	75%	81%	84%	88%	81%
Remote	89%	99%	99%	89%	99%	19%	82%

Characteristic	Validation Percent	Round 1 Percent	Round 2 Percent	Round 3 Percent	Round 4 Percent	Round 5 Percent	Average
Core Technology	90%	90%	80%	89%	72%	90%	85%
FT Project Manage	17%	30%	17%	30%	15%	17%	25%
FT Client Resource	79%	78%	77%	79%	80%	78%	79%
Sub-Contracted	92%	9%	5%	15%	44%	79%	41%
Time & Materials	69%	99%	99%	89%	79%	99%	89%
Not to Exceed	90%	69%	69%	69%	59%	70%	71%
Formal QA	99%	89%	89%	91%	92%	99%	93%
On-site	75%	78%	79%	81%	84%	88%	81%
Remote	89%	99%	99%	89%	99%	19%	82%





Chapter 09: Databases and Neural Networks

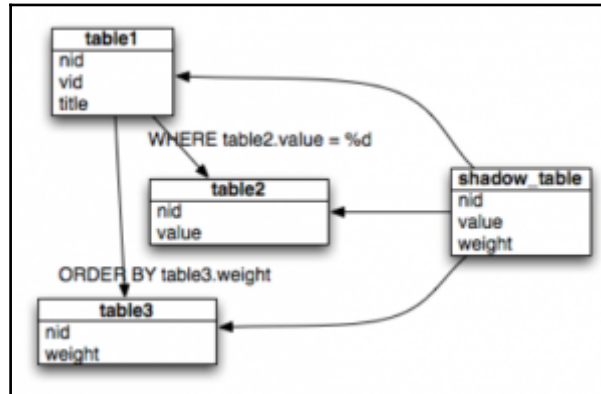




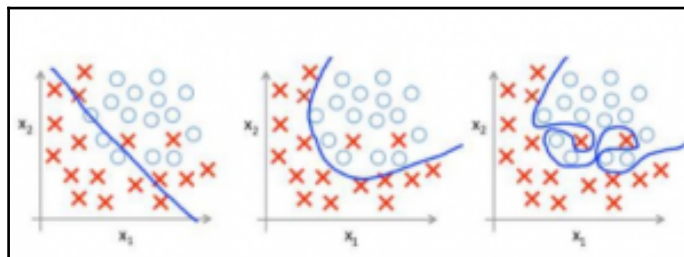
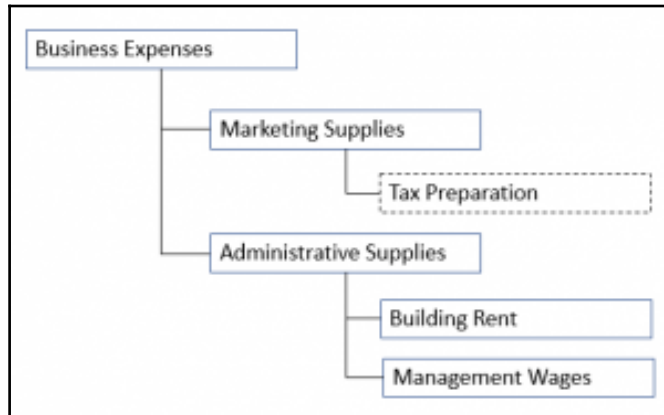
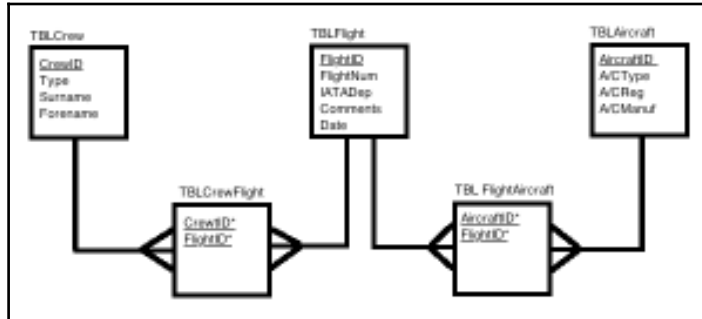
Input	Expected Output	Neural Net Output
1	1	0.9623402772
4	2	2.0083461217
9	3	2.9958221776
16	4	4.0009548085
25	5	5.0028838579
36	6	5.9975810435
49	7	6.9968278722
64	8	8.0070028670
81	9	9.0019220736
100	10	9.9222007864

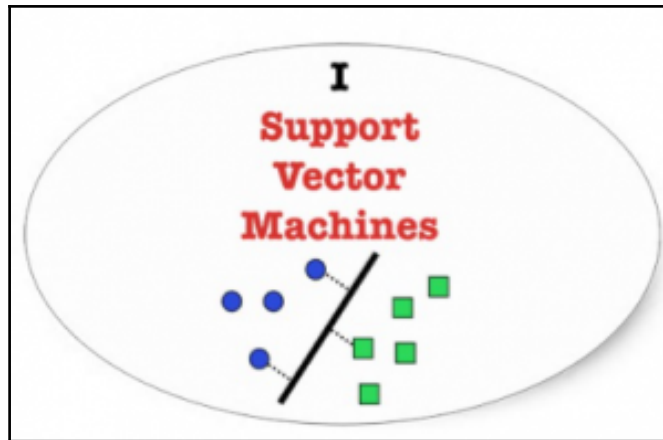
```
R Console
> set.seed(1)
> x <- runif(7)
>
> # Manually scaling
> (x - mean(x)) / sd(x)
[1] -1.01951259013 -0.48940037389 -0.04788275305  0.97047345797 -1.21713897716
[6]  0.94007370915  1.08338752711
>
> scale(x)
      [,1] [,2]
[1,] -1.01951259013
[2,] -0.48940037389
[3,] -0.04788275305
[4,]  0.97047345797
[5,] -1.21713897716
[6,]  0.94007370915
[7,]  1.08338752711
attr(,"scaled:center")
[1] 0.3847772287
attr(,"scaled:scale")
[1] 0.3229666497
> |
```

Chapter 10: Boosting your Database



Chapter 11: Database Classification using Support Vector Machines





UCI Machine Learning Repository

Statlog (German Credit Data) Data Set

Download: [Data Folder](#) [Data Set Description](#)

Abstract: This dataset classifies people according to a set of attributes as good or bad credit risks. Comes in two formats: one all numeric. Also comes with a post matrix.

Data Set Characteristics	Multi-variant	Number of Instances	1000	Area	Predefined
Attribute Characteristics	Categorical, Integer	Number of Attributes	20	Date Created	1994-11-17
Associated Tasks:	Classification	Missing Values?	Yes	Number of Attributes	10000

Source:
 Professor Dr. Hans-Hermann
 Institut für Statistik und Ökonometrie
 Universität Heidelberg
 F3 Wirtschaftswissenschaften
 Von-Miller-Platz 5
 69118 Heidelberg, FRG

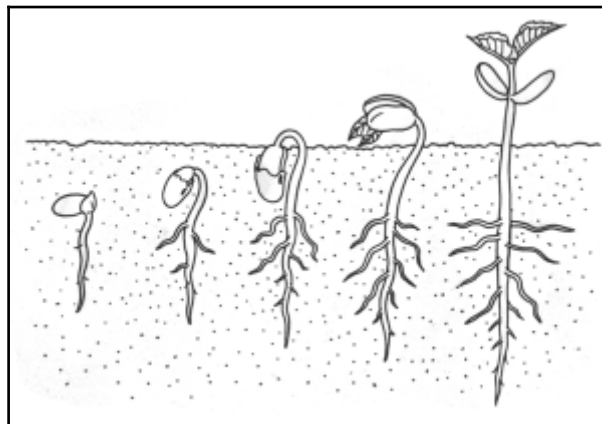
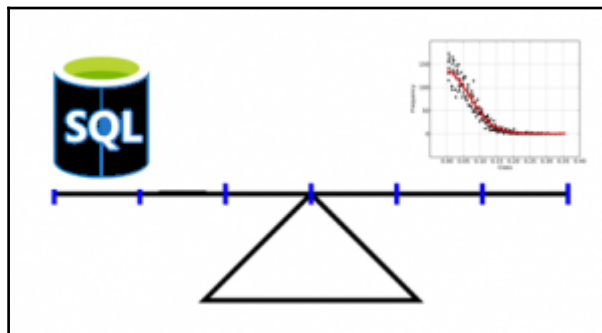
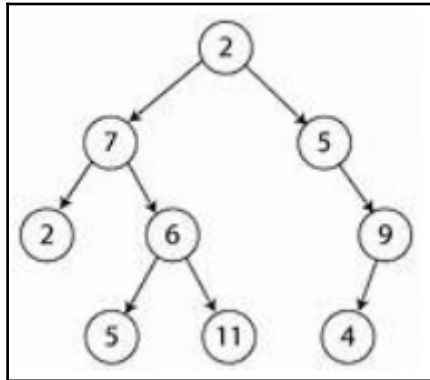
Index of /ml/machine-learning-databases/statlog/german

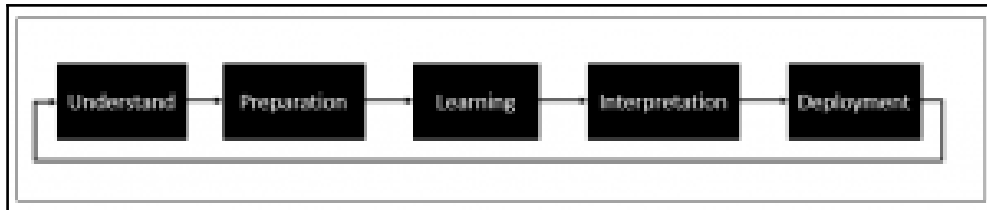
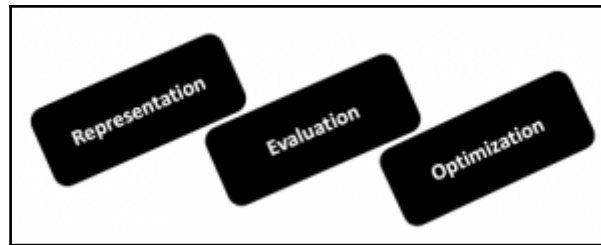
Name	Last modified	Size	Description
Parent Directory	-	-	-
index	05-Dec-1996 04:07	150	
german.data	17-Nov-1994 00:51	788K	
german.data.gz	17-Nov-1994 00:51	1008K	
german.doc	17-Nov-1994 00:51	4.68K	

Apache/2.2.15 (CentOS) Server at archive.ics.uci.edu Port 443

```
A11 6 A34 A43 1169 A65 A75 4 A93 A101 4 A121 67 A143 A152 2 A173 1 A192 A201 1
A12 48 A32 A43 5951 A61 A73 2 A92 A101 2 A121 22 A143 A152 1 A173 1 A191 A201 2
A14 12 A34 A46 2096 A61 A74 2 A93 A101 3 A121 49 A143 A152 1 A172 2 A191 A201 1
A11 42 A32 A42 7882 A61 A74 2 A93 A103 4 A122 45 A143 A153 1 A173 2 A191 A201 1
A11 24 A33 A40 4870 A61 A73 3 A93 A101 4 A124 53 A143 A153 2 A173 2 A191 A201 2
A14 36 A32 A46 9055 A65 A73 2 A93 A101 4 A124 35 A143 A153 1 A172 2 A192 A201 1
A14 24 A32 A42 2835 A63 A75 3 A93 A101 4 A122 53 A143 A152 1 A173 1 A191 A201 1
A12 36 A32 A41 6948 A61 A73 2 A93 A101 2 A123 35 A143 A151 1 A174 1 A192 A201 1
A14 12 A32 A43 3059 A64 A74 2 A91 A101 4 A121 61 A143 A152 1 A172 1 A191 A201 1
A12 30 A34 A40 5234 A61 A71 4 A94 A101 2 A123 28 A143 A152 2 A174 1 A191 A201 2
A12 12 A32 A40 1295 A61 A72 3 A92 A101 1 A123 25 A143 A151 1 A173 1 A191 A201 2
A11 48 A32 A49 4308 A61 A72 3 A92 A101 4 A122 24 A143 A151 1 A173 1 A191 A201 2
A12 12 A32 A43 1567 A61 A73 1 A92 A101 1 A123 22 A143 A152 1 A173 1 A192 A201 1
A11 24 A34 A40 1199 A61 A75 4 A93 A101 4 A123 60 A143 A152 2 A172 1 A191 A201 2
A11 15 A32 A40 1403 A61 A73 2 A92 A101 4 A123 28 A143 A151 1 A173 1 A191 A201 1
A11 24 A32 A43 1282 A62 A73 4 A92 A101 2 A123 32 A143 A152 1 A172 1 A191 A201 2
A14 24 A34 A43 2424 A65 A75 4 A93 A101 4 A122 53 A143 A152 2 A173 1 A191 A201 1
A11 30 A30 A40 8072 A65 A72 2 A93 A101 3 A123 25 A141 A152 3 A173 1 A191 A201 1
A12 24 A32 A41 12579 A61 A75 4 A92 A101 2 A124 44 A143 A153 1 A174 1 A192 A201 2
A14 24 A32 A43 3430 A63 A75 3 A93 A101 2 A123 31 A143 A152 1 A173 2 A192 A201 1
A14 9 A34 A40 2134 A61 A73 4 A93 A101 4 A123 48 A143 A152 3 A173 1 A192 A201 1
A11 6 A32 A43 2647 A63 A73 2 A93 A101 3 A121 44 A143 A151 1 A173 2 A191 A201 1
A11 10 A34 A40 2241 A61 A72 1 A93 A101 3 A121 48 A143 A151 2 A172 2 A191 A202 1
A12 12 A34 A41 1804 A62 A72 3 A93 A101 4 A122 44 A143 A152 1 A173 1 A191 A201 1
A14 10 A34 A42 2069 A65 A73 2 A94 A101 1 A123 26 A143 A152 2 A173 1 A191 A202 1
A11 6 A32 A42 1374 A61 A73 1 A93 A101 2 A121 36 A141 A152 1 A172 1 A192 A201 1
A14 6 A30 A43 426 A61 A75 4 A94 A101 4 A123 39 A143 A152 1 A172 1 A191 A201 1
```

Chapter 12: Database Structures and Machine Learning





id	A	B	C	D	E	F	G	H	I	J	K	L
1	Passenger	Survived	Sex	Name	Age	Age	Height	Weight	Class	Fare	Cabin	Embarked
2	1	0	0	0 Braund, Mr. Owen Louis	male	22	7	0	A/5 11131	7.25		S
3	2	1	1	0 Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	7	0	PC 17509	71.000	C85	C
4	3	1	1	0 Heikkinen, Mrs. Laina	female	26	0	0	0 P28162	5.000		S
5	4	1	1	0 Fatou, Mrs. Jacques-Francois (Jily May Belle)	female	35	1	0	21280	01.1 1128		S
6	5	0	0	0 Allen, Mr. William Henry	male	35	0	0	0 24034	8.000		S
7	6	0	0	0 Brown, Mr. James	male	0	0	0	0 00000	0.000		S
8	7	0	1	0 McCulloch, Mr. Thomas I	male	34	0	0	0 2460	51.000	F40	S
9	8	0	0	0 Paton, Mr. Charles (James)	male	2	1	1	0 00000	31.000		S
10	9	1	1	0 Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27	0	1	0 34700	51.000		S
11	10	1	1	0 Nassau, Mrs. John (John Nassau)	female	33	1	0	0 23700	50.000		C
12	11	1	1	0 Sandstrom, Miss. Margareta Est	female	4	1	1	0 000	01.1 00		S
13	12	1	1	0 Gibson, Miss. Elizabeth	female	28	0	0	0 21070	00.000	F100	S
14	13	0	0	0 Anderson, Mr. William Henry	male	20	0	0	0 0/5 1111	0.000		S
15	14	0	0	0 Anderson, Mr. Anders Johan	male	39	1	1	0 24000	11.000		S
16	15	0	0	0 Virtanen, Mrs. Sofia (Anneli Maria Paronensmarja)	female	39	0	0	0 00000	1.000		S
17	16	1	1	0 Stewart, Mrs. Mary B (Margaret)	female	55	0	0	0 24000	00		S
18	17	0	0	0 Cook, Martin Joseph	male	0	4	1	0 00000	00.000		S
19	18	1	1	0 Williams, Mr. Charles Eugene	male	0	0	0	0 24400	11		S
20	19	0	0	0 Vander Planke, Mrs. Julia (Emelia Maria Vandermonten)	female	30	1	0	0 24000	00		S
21	20	1	1	0 Gibson, Mrs. Fatma	female	0	0	0	0 2000	7.000		C