```
R version 3.2.2 (2015-08-14) -- "Fire Safety"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

  Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> setwd("C:/Users/TonyLaptop/Desktop/rowan/DM1/lecture11/lecture11")
> houses=read.csv("Houses.csv",header=TRUE,sep=",")
> minn = min(houses$MedianHomeVal)
> maxx = max(houses$MedianHomeVal)
> houses$MedianHomeVal2 <- sapply(houses$MedianHomeVal,function(b) {
+     return((b-minn)/(maxx-minn))
+   }
+ )
> m<-mean(houses$MedianHomeVal)
> sigma<-sd(houses$MedianHomeVal)
> houses$zMedHomeVal<-(houses$MedianHomeVal-m)/sigma
> m<-mean(houses$MedianIncome)
> sigma<-sd(houses$MedianIncome)
> houses$zMedInc<-(houses$MedianIncome-m)/sigma
> m<-mean(houses$MedianHomeAge)
> sigma<-sd(houses$MedianHomeAge)
> houses$zMedHomeAge<-(houses$MedianHomeAge-m)/sigma
> m<-mean(houses$TotalBedRooms)
> sigma<-sd(houses$TotalBedRooms)
> houses$zBedrooms<-(houses$TotalBedRooms-m)/sigma
> m<-mean(houses$TotalRooms)
```

```
> sigma<-sd(houses$TotalRooms)
> houses$zTotRooms<-(houses$TotalRooms-m)/sigma
> m<-mean(houses$Population)
> sigma<-sd(houses$Population)
> houses$zPop<-(houses$Population-m)/sigma
> m<-mean(houses$Households)
> sigma<-sd(houses$Households)
> houses$zHouseholds<-(houses$Households-m)/sigma
> m<-mean(houses$Latitude)
> sigma<-sd(houses$Latitude)
> houses$zLatitude<-(houses$Latitude-m)/sigma
> m<-mean(houses$Longitude)
> sigma<-sd(houses$Longitude)
> houses$zLongitude<-(houses$Longitude-m)/sigma
> vars<-
c("MedianHomeVal2","zMedInc","zMedHomeAge","zBedrooms","zTotRooms","zPop","zHouseholds","zLatitude","
zLongitude")
> houses2<-houses[vars]
> head(houses2)
  MedianHomeVal2   zMedInc zMedHomeAge  zBedrooms  zTotRooms       zPop
1      0.9022664 2.34470896   0.9821189 -0.9706826 -0.8047996 -0.9744050
2      0.7082466 2.33218146  -0.6070042  2.0 1.3486168  2.0458405  0.8614180
3      0.6950507 1.78265622   1.8561366 -0.8258748 -0.5357329 -0.8207575
4      0.6727828 0.93294491   1.8561366 -0.7190493 -0.6241995 -0.7660095
5      0.6746385 -0.01288068   1.8561366 -0.6122238 -0.4623928 -0.7598283
6      0.5251545 0.08744452   1.8561366 -0.7712751 -0.7869229 -0.8940491
  zHouseholds  zLatitude zLongitude
1  -0.9770092 1.052523  -1.327803
2   1.6699206 1.043159  -1.322812
3  -0.8436165 1.038478  -1.332794
4  -0.7337637 1.038478  -1.337785
5  -0.6291419 1.038478  -1.337785
6  -0.8017678 1.038478  -1.337785
>
> dim(houses2)
[1] 20640     9
> library("neuralnet")
Loading required package: grid
Loading required package: MASS
> set.seed(2)
```

```
> train=sample(1:20640,16000)
> test=-train
> trainingData=houses2[train,]
> testingData=houses2[test,]
> dim(trainingData)
[1] 16000     9
> dim(testingData)
[1] 4640     9
> names(trainingData)
[1] "MedianHomeVal2"  "zMedInc"      "zMedHomeAge"    "zBedrooms"
[5] "zTotRooms"       "zPop"         "zHouseholds"    "zLatitude"
[9] "zLongitude"
> formula=
MedianHomeVal2~zMedInc+zMedHomeAge+zBedrooms+zTotRooms+zPop+zHouseholds+zLatitude+zLongitude
> names(testingData)
[1] "MedianHomeVal2"  "zMedInc"      "zMedHomeAge"    "zBedrooms"    "zTotRooms"       "zPop"
"zHouseholds"
[8] "zLatitude"       "zLongitude"
> dim(testingData)
[1] 4640     9
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=5, threshold=0.2)
> proc.time() - ptm
   user  system elapsed
 152.46    3.33  301.99
> results<-compute(nnet,testingData[,2:9])
> dim(results$net.result)
[1] 4640     1
>
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640    11
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2007115494
```
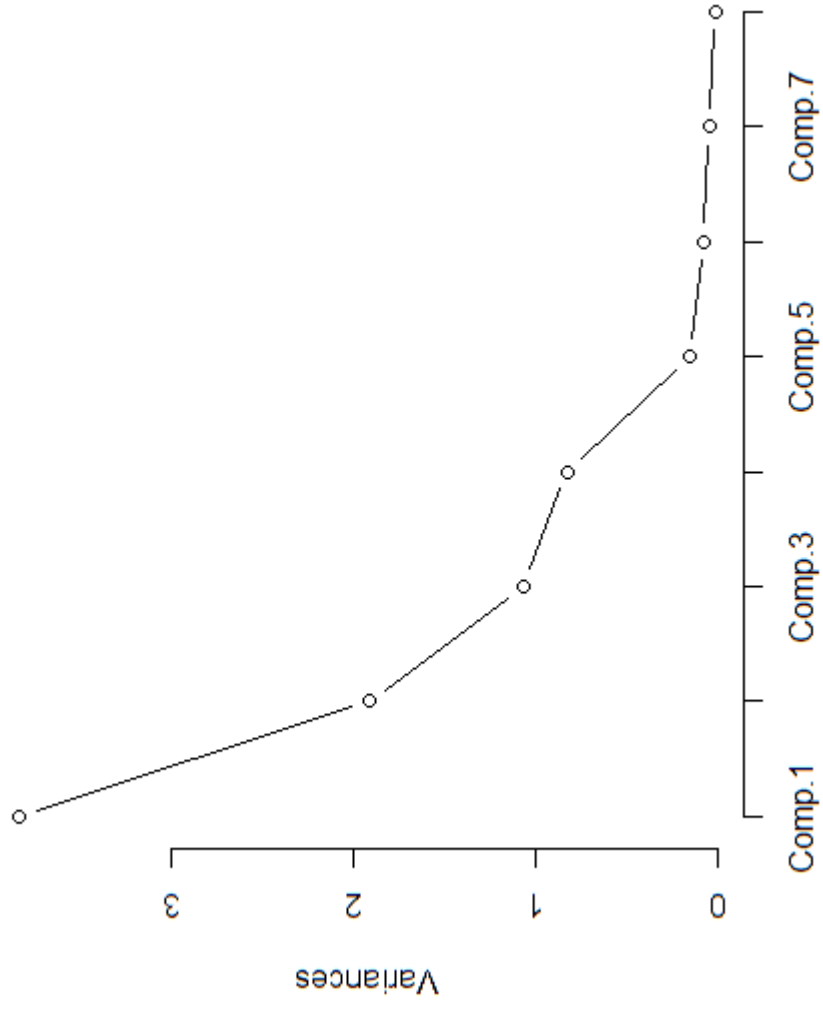
```
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.15)
> proc.time() - ptm
   user  system elapsed
 435.18   23.34  458.57
> results<-compute(nnet,testingData[,2:9])
> dim(results$net.result)
[1] 4640    1
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.1898130528
> names(houses2)
[1] "MedianHomeVal2" "zMedInc"     "zMedHomeAge"  "zBedrooms"    "zTotRooms"    "zPop"
"zHouseholds"
[8] "zLatitude"    "zLongitude"
> names(trainingData)
[1] "MedianHomeVal2" "zMedInc"     "zMedHomeAge"  "zBedrooms"    "zTotRooms"    "zPop"
"zHouseholds"
[8] "zLatitude"    "zLongitude"
> dim(trainingData)
[1] 16000     9
> PCA<-princomp(trainingData[,2:9])
> screeplot(PCA,type="lines")
```

# PCA



Variances

Comp.1  Comp.3  Comp.5  Comp.7

```
Loadings:
             Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7  Comp.8
zMedInc               0.219   0.897  -0.395                           0.171
zMedHomeAge          -0.490  -0.381  -0.891
zBedrooms            -0.479  -0.114           0.377  -0.216   0.230  -0.705
zTotRooms                            -0.111   0.307   0.546  -0.574   0.144
zPop         -0.475  -0.113          -0.850   0.117          -0.123
zHouseholds  -0.493  -0.106           0.160  -0.388   0.308   0.679
zLatitude     0.701                           0.482   0.505
zLongitude   -0.701                           0.497   0.486

               Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7  Comp.8
SS loadings     1.000   1.000   1.000   1.000   1.000   1.000   1.000   1.000
Proportion Var  0.125   0.125   0.125   0.125   0.125   0.125   0.125   0.125
Cumulative Var  0.125   0.250   0.375   0.500   0.625   0.750   0.875   1.000

> trainingData$C1<-0
> trainingData$C2<-0
> trainingData$C3<-0
> trainingData$C4<-0
> trainingData$C5<-0
> for (i in 1:16000){
+   trainingData$C1[i]<-trainingData$zMedHomeAge[i]*.219-trainingData$zBedrooms[i]*.49-
trainingData$zTotRooms[i]*.479-trainingData$zHouseholds[i]*.493-trainingData$zPop[i]*.475
+   trainingData$C2[i]<-trainingData$zLatitude[i]*.701-trainingData$zlongitude[i]*.701
+   trainingData$C3[i]<-trainingData$zMedInc[i]*.897-trainingData$zMedHomeAge[i]*.381-
trainingData$zBedrooms[i]*.114-trainingData$zPop[i]*.113-trainingData$zHouseholds[i]*.106
+   trainingData$C4[i]<-trainingData$zLatitude[i]*.109-trainingData$zMedInc[i]*.395-
trainingData$zMedHomeAge[i]*.891-trainingData$zTotRooms[i]*.111
+   trainingData$C5[i]<-trainingData$zBedrooms[i]*.377+trainingData$zTotRooms[i]*.307-
trainingData$zPop[i]*.85+trainingData$zHouseholds[i]*.16
+ }

> head(trainingData)
  MedianHomeVal2     zMedInc  zMedHomeAge    zBedrooms    zTotRooms        zPop
zHouseholds     zLatitude  zLongitude
```

```
3816   0.4329899671 -1.3233720718 -0.2891795935 -0.6897472321 -0.5999995330  -
0.2080394728 -0.6656803006  0.5389006047
14497  0.8299955468  0.6481287088 -1.0837411305 -0.09471385615 -0.2308922649  -
0.14778819784 -1.3023986103  1.1877548179
11833  0.3298976087 -0.4582908989 -1.0042849768 -0.8543615 4931 -1.2198878158  -
1.2516412251  1.7173316083 -0.3794776662
3469   0.3245368060 -0.4582908989  0.1875573287 -0.6691974 2410 -0.4710768987  -
0.6840683423 -0.6188627779  0.5438917910
19478  0.2195887852 -0.7741100068  0.1875573287 -0.5813631 5958 -0.5593800729  -
0.5375979209  0.9542059871 -0.6989135865
19469  0.2016507148 -0.3610186137 -0.1302672861  0.5199549 071  0.6280836722  -
0.5713924123  0.9588877394 -0.7039047728

              C1            C2            C3            C4            C5
3816   0.75802893700 -0.8444112147 -0.9629772123  0.7843937 7615  0.18540072636
14497 -0.06526847607 -1.7455975532  1.0468405178  0.5544511 5371  0.17337037261
11833  1.78333833910  1.4698633014  0.3394641495  1.3529752 7704  0.26578560911
3469   1.26916042014 -0.8150909528 -0.2805148384  0.0250524 1448 -0.17870425692
19478  1.07355963748  1.1588368212 -0.5793652902  0.2929254 4656  0.03128405822
19469 -1.03136156824  1.1656175511 -0.4333865370  0.2934719 8030  0.19919477694
> formula=MedianHomeVal2~C1+C2+C3+C4+C5
> dim(testingData)
[1] 4640   13
> testingData$C1<-0
> testingData$C2<-0
> testingData$C3<-0
> testingData$C4<-0
> testingData$C5<-0
> for (i in 1:4640){
+ testingData$C1[i]<-testingData$zMedHomeAge[i]*.49-testingData$zBedrooms[i]*.49-
testingData$zTotRooms[i]*.479-testingData$zHouseholds[i]*.493-testingData$zPop[i]*.475
+ testingData$C2[i]<-testingData$zLatitude[i]*.701-testingData$zLongitude[i]*.701
+ testingData$C3[i]<-testingData$zMedInc[i]*.897-testingData$zMedHomeAge[i]*.381-
testingData$zBedrooms[i]*.114-testingData$zPop[i]*.113-testingData$zHouseholds[i]*.106
+ testingData$C4[i]<-testingData$zLatitude[i]*.109-testingData$zMedInc[i]*.395-
testingData$zMedHomeAge[i]*.891-testingData$zTotRooms[i]*.111
+ testingData$C5[i]<-testingData$zBedrooms[i]*.377+testingData$zTotRooms[i]*.307-
testingData$zPop[i]*.85+testingData$zHouseholds[i]*.16
+ }
> names(testingData)
```

```
 [1] "MedianHomeVal2" "zMedInc"        "zMedHomeAge"    "zBedrooms"      "zTotRooms"      "zPop"
 "zHouseholds"
 [8] "zLatitude"      "zLongitude"     "guess"          "Actual"         "absDiff"        "dummy"
"C1"
[15] "C2"             "C3"             "C4"             "C5"
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.15)
> proc.time() - ptm
   user  system elapsed
 282.94    2.06  285.01
> results<-compute(nnet,testingData[,14:18])
> dim(results$net.result)
[1] 4640    1
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640   18
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2332437547
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=5, threshold=0.15)
> proc.time() - ptm
   user  system elapsed
 199.15    4.89  204.13
> results<-compute(nnet,testingData[,14:18])
> dim(results$net.result)
[1] 4640   18
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640   18
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2495409405
```

```
> #might be getting over trained; test on training data and see if results are better
> #if results are much better then the model is overtrained
> dim(trainingData)
[1] 16000    14
> names(trainingData)
 [1] "MedianHomeVal2" "zMedInc"       "zMedHomeAge"    "zBedrooms"      "zTotRooms"      "zPop"
 "zHouseholds"
 [8] "zLatitude"      "zLongitude"     "C1"             "C2"             "C3"             "C4"
 "C5"
> dim(testingData)
[1] 4640    18
> #save testing data
> t2<-testingData
> #use the first 4640 rows of training data as testing data
> testingData<-trainingData[1:4640,]
> dim(testingData)
[1] 4640    14
> results<-compute(nnet,testingData[,10:14])
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640    18
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2474397734
> #now turn the last 4640 rows of training data into testing data
> dim(trainingData)
[1] 16000    14
> 16000-4640
[1] 11360
> testingData<-trainingData[11360:16000,]
> results<-compute(nnet,testingData[,10:14])
> dim(results$net.result)
[1] 4641     1
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
```

```
> dim(testingData)
[1] 4641   16
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2475267321
> #apparently not overtrained
> #go back to original training data
> testingData<-t2
> dim(testingData)
[1] 4640   18
> PCA$loadings

Loadings:
            Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8
zMedInc             0.897 -0.395                      0.171
zMedHomeAge  0.219 -0.381 -0.891
zBedrooms   -0.490 -0.114         0.377 -0.216  0.230 -0.705
zTotRooms   -0.479        -0.111  0.307  0.546 -0.574  0.144
zPop        -0.475 -0.113        -0.850  0.117        -0.123
zHouseholds -0.493 -0.106         0.160 -0.388  0.308  0.679
zLatitude    0.701                0.109         0.482  0.505
zLongitude  -0.701                              0.497  0.497  0.486

                Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8
SS loadings      1.000  1.000  1.000  1.000  1.000  1.000  1.000  1.000
Proportion Var   0.125  0.125  0.125  0.125  0.125  0.125  0.125  0.125
Cumulative Var   0.125  0.250  0.375  0.500  0.625  0.750  0.875  1.000
```

```
> #perhaps the issue is that we only have 3 digits of accuracy from the loadings above
> #let's see if we can get a few more digits
> load <- with(PCA, unclass(loadings))
> load
```

|              | Comp.1        | Comp.2        | Comp.3        | Comp.4        | Comp.5        | Comp.6        | Comp.7        | Comp.8        |
|--------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| zMedInc      | -0.0425623597 | -0.0372126309 | 0.8969731005  | -0.3949814596 | -0.0529260728 | -0.0506998566 | 0.1713223917  | -0.0388530885 |
| zMedHomeAge  | 0.2190433696  | 0.0284699830  | -0.3809362597 | -0.8913116327 | 0.0359419449  | 0.0942858674  | 0.0382733138  | -0.0037557203 |
| zBedrooms    | -0.4898694489 | 0.0637953578  | -0.1135057640 | -0.0646287672 | 0.3767775509  | -0.2158356384 | 0.2299775707  | -0.7053030025 |
| zTotRooms    | -0.4792152063 | 0.0769738234  | 0.0961902348  | -0.1149814146 | 0.3068033969  | 0.5463959364  | -0.5738488422 | 0.1440973072  |
| zPop         | -0.4751637056 | 0.0323496332  | -0.1134462582 | -0.0875560865 | -0.8504245961 | 0.1171489851  | 0.0252002639  | -0.1227024879 |
| zHouseholds  | -0.4927034568 | 0.0671885065  | -0.1057168175 | -0.0979592998 | 0.1597922137  | -0.3878561593 | 0.3083810729  | 0.6788903447  |
| zLatitude    | 0.0766241402  | 0.7006673952  | 0.0158344959  | 0.1088144934  | 0.0461871623  | 0.4819969793  | 0.5050615798  | 0.0395104505  |
| zLongitude   | -0.0786171369 | -0.7009381357 | -0.0624009040 | 0.0619824172  | 0.0944910626  | 0.4969192936  | 0.4859641356  | 0.0524489881  |

```
> trainingData$C1<-0
> trainingData$C2<-0
> trainingData$C3<-0
> trainingData$C4<-0
> trainingData$C5<-0
> for (i in 1:16000){
+ trainingData$C1[i]<-trainingData$zMedHomeAge[i]*.21904336969-
trainingData$zBedrooms[i]*.4898694891-trainingData$zTotRooms[i]*.47921520638-
trainingData$zHouseholds[i]*.49270345689-trainingData$zPop[i]*.47516370568
+ trainingData$C2[i]<-trainingData$zLatitude[i]*.70066739522-trainingData$zLongitude[i]*.70093813570
+ trainingData$C3[i]<-trainingData$zMedInc[i]*.8969731005 -
trainingData$zMedHomeAge[i]*.3809625978-trainingData$zBedrooms[i]*.11350576403-
trainingData$zPop[i]*.11344625820-trainingData$zHouseholds[i]*.10571681754
+ trainingData$C4[i]<-trainingData$zLatitude[i]*.10881449345-trainingData$zMedInc[i]*.3949814561-
trainingData$zMedHomeAge[i]*.8913116277-trainingData$zTotRooms[i]*.1149814146
+ trainingData$C5[i]<-trainingData$zBedrooms[i]*.37677755094+trainingData$zTotRooms[i]*.30680339698-
trainingData$zPop[i]*.8504245961+trainingData$zHouseholds[i]*.15979221377
```

```
+ }
> testingData$C1<-0
> testingData$C2<-0
> testingData$C3<-0
> testingData$C4<-0
> testingData$C5<-0
> for (i in 1:4640){
+ testingData$C1[i]<-testingData$zMedHomeAge[i]*.2190436969-testingData$zBedrooms[i]*.4898694891-
testingData$zTotRooms[i]*.47921520638-testingData$zHouseholds[i]*.49270345689-
testingData$zPop[i]*.47516370568
+ testingData$C2[i]<-testingData$zLatitude[i]*.70066739522-testingData$zLongitude[i]*.70093813570
+ testingData$C3[i]<-testingData$zMedInc[i]*.89697310056 -testingData$zMedHomeAge[i]*.38093625978-
testingData$zBedrooms[i]*.11350576403-testingData$zPop[i]*.11344625820-
testingData$zHouseholds[i]*.10571681754
+ testingData$C4[i]<-testingData$zLatitude[i]*.10881449345-testingData$zMedInc[i]*.39498145961-
testingData$zMedHomeAge[i]*.89131163277-testingData$zTotRooms[i]*.11149814146
+ testingData$C5[i]<-testingData$zBedrooms[i]*.37677755094+testingData$zTotRooms[i]*.30680339698-
testingData$zPop[i]*.85042459610+testingData$zHouseholds[i]*.15979221377
+ }
>
> formula
MedianHomeVal2 ~ C1 + C2 + C3 + C4 + C5
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.15)
> proc.time() - ptm
   user  system elapsed
 287.31    8.14  295.53
> results<-compute(nnet,testingData[,14:18])
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640   18
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.234961167
> #a little bit better, but not much
> ptm <- proc.time()
```

```r
> nnet<-neuralnet(formula,trainingData, hidden=10, threshold=0.1)
> proc.time() - ptm
   user  system elapsed
 769.07   28.44  797.63
> results<-compute(nnet,testingData[,14:18])
> dim(results$net.result)
[1] 4640    1
> testingData$guess<-results$net.result
> testingData$Actual<-(testingData$MedianHomeVal2*(maxx-minn))+minn
> testingData$guess<-(testingData$guess*(maxx-minn))+minn
> dim(testingData)
[1] 4640   18
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2349090405
> #no improvement
```

```
> #A common trick is to change the target variable to log of that variable
> #this has the affect of making non-linear data mostly linear
> trainingData$LogHomeVal<-log(trainingData$Actual)
> head(trainingData)
      MedianHomeVal2       zMedInc   zMedHomeAge     zBedrooms     zTotRooms          zPop
zHouseholds    zLatitude    zLongitude
3816    0.4329899671 -1.3233720718 -0.2891795935 -0.21103490916 -0.6897472321 -0.5999995330 -
0.2080394728 -0.6656803006  0.5389006047
14497   0.8298955468  0.6481287088 -1.0837411305 -0.09471385615  0.1188279732 -0.2308922649 -
0.1478819784 -1.3023986103  1.1877548179
11833   0.3298976087 -0.4582908989 -1.0042849768 -0.85436154931 -0.8103001114 -1.2198878158 -
1.2516412251  1.7173316083 -0.3794776662
3469    0.3245368060 -0.4582908989  0.1875573287 -0.66919742410 -0.7080822708 -0.4710768987 -
0.6840683423 -0.6188627779  0.5438917910
19478   0.2195887852 -0.7741100068  0.1875573287 -0.58136315958 -0.4527668573 -0.5593800729 -
0.5375979209  0.9542059871 -0.6989135865
19469   0.2016507148 -0.3610186137 -0.1302672861  0.53199549071  0.6280836722  0.3360141134
0.5713924123  0.9588877394 -0.7039047728
                 C1            C2            C3           C4            C5 Actual LogHomeVal
3816   0.75817381320 -0.8441564675 -0.9628555061 0.7849264379 1 0.18588126247 225000 12.32385568
14497 -0.06535947012 -1.7450908895  1.0467683546 0.5549833098 5 0.17349684347 417500 12.94203982
11833  1.78318616447  1.4692686326  0.3391801498 1.3533648160 8 0.26691300191 175000 12.07254125
3469   1.26910783603 -0.8148514686 -0.2808047913 0.0254529968 1 -0.17807402608 172400 12.05757264
19478  1.07352146524  1.1584762099 -0.5795224531 0.2929011757 3 0.03185161328 121500 11.70766954
19469 -1.03131849785  1.1652550738 -0.4331103385 0.2930151279 6 0.19869155414 112800 11.63337162
> max2<-max(trainingData$LogHomeVal)
> min2<-min(trainingData$LogHomeVal)
> max2
[1] 13.12236538
> min2
[1] 9.615738811
> trainingData$MedianHomeVal3<-(trainingData$LogHomeVal-min2)/(max2-min2)
> head(trainingData)
```

```
        MedianHomeVal2       zMedInc    zMedHomeAge      zBedrooms      zTotRooms          zPop
zHouseholds      zLatitude     zLongitude
3816       0.432989671 -1.3233720718 -0.2891795935 -0.2110349916 -0.6897472321 -0.599995330 -
0.2080394728 -0.6656803006  0.5389006047
14497      0.8298955468  0.6481287088 -1.0837411305 -0.0947138615  0.1188279732 -0.2308922649 -
0.1478819784 -1.3023986103  1.1877548179
11833      0.3298976087 -0.4582908989 -1.0042849768 -0.8543615493 -0.8103001114 -1.2198878158 -
1.2516412251  1.7173316083 -0.3794776662
3469       0.3245368060 -0.4582908989  0.1875573287 -0.6691974241 -0.7080822708 -0.4710768987 -
0.6840683423 -0.6188627779  0.5438917910
19478      0.2195887852 -0.7741100068  0.1875573287 -0.5813631595 -0.4527668573 -0.5593800729 -
0.5375979209  0.9542059871 -0.6989135865
19469      0.2016507148 -0.3610186137 -0.1302672861  0.5319954907  0.6280836722  0.336014134
0.5713924123  0.9588877394 -0.7039047728

                  C1            C2            C3            C4            C5  Actual  LogHomeVal
MedianHomeVal3
3816     0.7581738132 -0.8441564675 -0.9628555061  0.7849264379  0.18588126247  225000 12.32385568
0.7722855054
14497   -0.0653594701 -1.7450908895  1.0467683546  0.5549833098  0.17349684347  417500 12.94203982
0.9485757748
11833    1.7831861644  1.4692686326  0.3391801498  1.3533648160  0.26691300191  175000 12.07254125
0.7006170732
3469     1.2691078360 -0.8148514686 -0.2808047913  0.0254529968 -0.17807402608  172400 12.05757264
0.6963484078
19478    1.0735214652  1.1584762099 -0.5795224531  0.2929011757  0.03185161328  121500 11.70766954
0.5965650151
19469   -1.0313184978  1.1652550738 -0.4331103385  0.2930151279  0.19869155414  112800 11.63337162
0.5753771520
> exp(12.32385568)
[1] 224775.1122
> #rounding costs us $25 which is about .01% error
> formula
MedianHomeVal2 ~ C1 + C2 + C3 + C4 + C5
> #rerun model with MedianHomeVal3 which is normalized Log of median home val
> formula<-MedianHomeVal3~C1+C2+C3+C4+C5
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.1)
> proc.time() - ptm
   user  system elapsed
 970.09   19.77  990.00
```

```
> results<-compute(nnet,testingData[,14:18])
> testingData$guess<-results$net.result
> testingData$guess<-testingData$guess*(max2-min2)
> testingData$guess<-testingData$guess+min2
> testingData$guess<-exp(testingData$guess)
> testingData$absDiff<-abs(testingData$guess-testingData$Actual)
> avgErr<-mean(testingData$absDiff)
> avgHomeVal<-mean(testingData$Actual)
> avgErr/avgHomeVal
[1] 0.2341468598
> #no improvement with log home values
> #Let's try a regression model to see if we can get a better result
> regmodel=lm(trainingData$MedianHomeVal3~trainingData$C1+trainingData$C2+trainingData$C3+
trainingData$C4+trainingData$C5)
> summary(regmodel)

Call:
lm(formula = trainingData$MedianHomeVal3 ~ trainingData$C1 +
    trainingData$C2 + trainingData$C3 + trainingData$C4 + trainingData$C5)

Residuals:
        Min          1Q      Median          3Q         Max
-0.77864590 -0.07110792  0.00272059  0.07356449  0.68574262

Coefficients:
                      Estimate    Std. Error   t value               Pr(>|t|)
(Intercept)       0.7045675482  0.0009144070  770.51854 < 0.00000000000000222 ***
trainingData$C1  -0.0195710714  0.0004945693  -39.57195 < 0.00000000000000222 ***
trainingData$C2  -0.0041434612  0.0006735235   -6.15192   0.0000000000078374 ***
trainingData$C3   0.0808126417  0.0009068346   89.11508 < 0.00000000000000222 ***
trainingData$C4  -0.0779133549  0.0010346019  -75.30757 < 0.00000000000000222 ***
trainingData$C5   0.0517195795  0.0023857551   21.67850 < 0.00000000000000222 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1156565 on 15994 degrees of freedom
Multiple R-squared:  0.493714,  Adjusted R-squared:  0.493213
F-statistic: 3115.095 on 5 and 15994 DF,  p-value: < 0.00000000000000022204
> #terrible R squared. This is not going to work too well
> dim(testingData)
```

```
[1] 4640   18
> for (i in 1:4640)
+ {
+   testingData$guess[i]<-testingData$C5[i]*.0517195795-
trainingData$C4[i]*.0779133549+trainingData$C3[i]*.0808126417-trainingData$C2[i]*
+ .0041434612-trainingData$C1[i]*.0195710714+.7045675482
+ }
> testingData$guess<-testingData$guess*(max2-min2)
> testingData$guess<-testingData$guess+min2
> testingData$guess<-exp(testingData$guess)
> mean(testingData$absDiff)
[1] 110002.2071
> mean(testingData$Actual)
[1] 204041.6976
> mean(testingData$absDiff)/mean(testingData$Actual)
[1] 0.53911631
> #awful
>
> #next thing to try is the use of user defined composite variable
> #go back to original houses data
> head(houses)
  MedianHomeVal MedianIncome MedianHomeAge TotalRooms TotalBedRooms Population Households Latitude
Longitude MedianHomeVal2
1        452600       8.3252            41        880           129        322        126    37.88
-122.23     0.9022663824
2        358500       8.3014            21       7099          1106       2401       1138    37.86
-122.22     0.7082465639
3        352100       7.2574            52       1467           190        496        177    37.85
-122.24     0.6950507421
4        341300       5.6431            52       1274           235        558        219    37.85
-122.25     0.6727827926
5        342200       3.8462            52       1627           280        565        259    37.85
-122.25     0.6746384551
6        269700       4.0368            52        919           213        413        193    37.85
-122.25     0.5251545354
     zMedHomeVal       zMedInc    zMedHomeAge      zBedrooms      zTotRooms           zPop    zHouseholds
zLatitude      zLongitude
1   2.1295798911    2.3447089561    0.9821188657  -0.9706826023  -0.8047995998  -0.9744049915  -0.9770091850
1.052522785   -1.327803055
```

```
2  1.314124014   2.3321814648  -0.6070042083   1.348167609   2.0458405374   0.8614179998   1.6699205726
   1.043159280  -1.322811868
3  1.258662922   1.782662721    1.856136556   -0.825874760  -0.535732907  -0.820757468  -0.843616479
   1.038477528  -1.332794241
4  1.165071845   0.932944907    1.856136556   -0.719049304  -0.624199468  -0.766009500  -0.733763663
   1.038477528  -1.337785427
5  1.172871101  -0.012880683    1.856136556   -0.612223847  -0.462392752  -0.759828278  -0.629141934
   1.038477528  -1.337785427
6  0.544597666   0.087444519    1.856136556   -0.771275082  -0.786922937  -0.894049103  -0.801767788
   1.038477528  -1.337785427
> min3<-min(houses$MedianIncome)
> max3<-max(houses$MedianIncome)
> houses$medInc<-(houses$MedianIncome-min3)/(max3-min3)
> names(houses)
 [1] "MedianHomeVal"   "MedianIncome"    "MedianHomeAge"   "TotalRooms"    "TotalBedRooms"
     "Population"      "Households"
 [8] "Latitude"        "Longitude"       "MedianHomeVal2"  "zMedHomeVal"   "zMedInc"
     "zBedrooms"       "zMedHomeAge"
[15] "zTotRooms"       "zPop"            "zHouseholds"     "zLatitude"     "zLongitude"
     "medInc"
> min3<-min(houses$MedianHomeAge)
> max3<-max(houses$MedianHomeAge)
> houses$medAge<-(houses$MedianHomeAge-min3)/(max3-min3)
> min3<-min(houses$TotalRooms)
> max3<-max(houses$TotalRooms)
> houses$tRooms<-(houses$TotalRooms-min3)/(max3-min3)
> min3<-min(houses$Population)
> max3<-max(houses$Population)
> houses$Pop<-(houses$Population-min3)/(max3-min3)
> min3<-min(houses$Households)
> max3<-max(houses$Households)
> houses$Hholds<-(houses$Households-min3)/(max3-min3)
> min3<-min(houses$Latitude)
> max3<-max(houses$Latitude)
> houses$lat<-(houses$Latitude-min3)/(max3-min3)
> min3<-min(houses$Longitude)
> max3<-max(houses$Longitude)
> houses$long<-(houses$Longitude-min3)/(max3-min3)
> min3<-min(houses$TotalBedRooms)
> max3<-max(houses$TotalBedRooms)
> houses$tBedRooms<-(houses$TotalBedRooms-min3)/(max3-min3)
```

```
> names(houses)
 [1] "MedianHomeVal"  "MedianIncome"   "MedianHomeAge"  "TotalRooms"     "TotalBedRooms"
 [8] "Population"     "Households"     "Latitude"       "Longitude"      "zMedHomeVal"    "zMedInc"
[15] "zTotRooms"      "zBedrooms"      "zPop"           "zHouseholds"    "zLatitude"      "zLongitude"     "medInc"
"medAge"
[22] "tRooms"         "tBedRooms"      "Pop"            "Hholds"         "lat"            "long"           "tBedRooms"
> vars<-c("MedianHomeVal2","medInc","medAge","tRooms","Pop","Hholds","lat","long","tBedRooms")
> houses3<-houses[vars]
> dim(houses3)
[1] 20640      9
> head(houses3)
  MedianHomeVal2    medInc    medAge     tRooms       Pop       Hholds         lat        long
          tBedRooms
1    0.902663824 0.539668184 0.7843137255 0.022330739l0 0.0089408833544 0.0205558296  0.5674814028
0.2115537850 0.019863343886
2    0.708246563  0.538027061l7 0.3921568627 0.1805025688 0.067210403879 0.18697582634 0.565356 0043
0.2121513944 0.17147734327
3    0.6950507421 0.4660280548 1.000000000 0.03726028791 0.013817651840 0.02894260812 0.5642933050
0.2101593625 0.029329608 94
4    0.6727827926 0.3546985559 1.000000000 0.03235159469 0.01555536870 4 0.03584936688 0.5642933050
0.2091633466 0.036312849l6
5    0.6746384551 0.2307761272 1.000000000 0.04l32967089 0.015751562544 0.04242723236 0.5642933050
0.2091633466 0.043296089 39
6    0.5251545354 0.2439207735 1.000000000 0.023322651l20 0.011491353457 0.03157375432 0.5642933050
0.2091633466 0.0328 982061
> max(houses3$long)
[1] 1
> #make a composite variable
> houses3$comp1<- (houses3$tRooms+houses3$Pop+houses3$Hholds+houses3$tBedRooms)/4
> head(houses3)
  MedianHomeVal2    medInc    medAge     tRooms       Pop       Hholds         lat        long
          tBedRooms
1    0.902663824 0.539668184 0.7843137255 0.022330739l0 0.0089408833544 0.0205558296  0.5674814028
0.2115537850 0.019863343886
2    0.708246563  0.538027061l7 0.3921568627 0.1805025688 0.067210403879 0.18697582634 0.565356 0043
0.2121513944 0.17147734327
3    0.6950507421 0.4660280548 1.000000000 0.03726028791 0.013817651840 0.02894260812 0.5642933050
0.2101593625 0.029329608 94
```

```
4  0.6727827926 0.3546985559 1.0000000000 0.0323515946  0.0155536870  0.5642933050
   0.2091633466 0.0363128491 6
5  0.6746384551 0.2307761272 1.0000000000 0.0413296708  0.0157516254  0.5642933050
   0.2091633466 0.0432960893 9
6  0.5251545354 0.2439207735 1.0000000000 0.0233226512  0.0114935345  0.5642933050
   0.2091633466 0.0328982061

         comp1
1  0.0179227102 8
2  0.1515415355 7
3  0.0273375392 0
4  0.0300172948 6
5  0.0357011388 0
6  0.0248216449 0
> max(houses3$comp1)
[1] 0.9075818929
> min(houses3$comp1)
[1] 0.0001009282068
> houses3$comp2<-(houses3$lat-houses3$long)/2
> max(houses3$comp2)
[1] 0.491665114 3
> min(houses3$comp2)
[1] -0.4742327184
> #add .5 so min and max go back to 0 to 1
> houses3$comp2<-houses3$comp2+.5
> formula<-MedianHomeVal2~comp1+comp2+medInc+medAge
> trainingData<-houses3[train,]
> testingData<-houses3[test,]
> dim(trainingData)
[1] 16000    11
> dim(testingData)
[1] 4640    11
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=5, threshold=0.1)
> proc.time() - ptm
   user  system elapsed
 399.67   23.70  423.45
> formula
MedianHomeVal2 ~ comp1 + comp2 + medInc + medAge
> vars<-c("comp1","comp2","medInc","medAge")
> test2<-testingData[vars]
```

```
> results<-compute(nnet,test2)
> testingData$guess<-results$net.result
> min3<-min(houses$MedianHomeVal)
> max3<-max(houses$MedianHomeVal)
> testingData$Actual<-(testingData$MedianHomeVal2*(max3-min3))+min3
> head(testingData)
   MedianHomeVal2    medInc       medAge    tRooms           Pop         Hholds        lat
          tBedRooms
8      0.4668042606 0.1806940593 1.00000000000 0.0788951625 0.0323439558 0.1062325275 0.5632306057
   0.2091633466 0.1064556176
10     0.5074226498 0.2200866194 1.00000000000 0.0902131339 0.0433868662 0.1172504522 0.5632306057
   0.2091633466 0.1095592799
19     0.2962894998 0.1028399608 0.9607843137 0.0568950607 0.0276633313 0.0687386942 0.5632306057
   0.2081673307 0.0704531347
21     0.2731968116 0.0591647011 0.7647058824 0.0190497990 0.0113792426 0.0271336951 0.5642933050
   0.2071713147 0.0283985102
24     0.1746405169 0.1159087460 1.00000000000 0.0428811231 0.0238235376 0.0532807104 0.5632306057
   0.2071713147 0.0521415270
33     0.1967022816 0.0894470421 0.9215686275 0.0488325957 0.0286723282 0.0549251767 0.5632306057
   0.2071713147 0.0633147113
           comp1         comp2         guess        absDiff    Actual
8      0.0809818158 0.677033629  0.508663765  0.0418595046  241400
10     0.0901024330 0.677033629  0.618262097  0.1108394476  261100
19     0.0559375552 0.677531637  0.265669248  0.0306202515  158700
21     0.0214903117 0.678560995  0.129320229  0.1438765824  147500
24     0.0430317245 0.678029645  0.292556804  0.1179162871   99700
33     0.0489362030 0.678029645  0.217713746  0.0210114644  110400
> testingData$guess<-(testingData$guess*(max3-min3))+min3
> testingData$absDiff<-abs(testingData$Actual-testingData$guess)
> mean(testingData$absDiff)
[1] 55019.87943
> mean(testingData$absDiff)/mean(testingData$Actual)
[1] 0.269501748
> #worse than PCA
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.05)
> proc.time() - ptm
   user  system elapsed
 729.37   45.05  774.56
> vars<-c("comp1","comp2","medInc","medAge")
```

```
> test2<-testingData[vars]
> results<-compute(nnet,test2)
> testingData$guess<-results$net.result
> testingData$guess<-(testingData$guess*(max3-min3))+min3
> testingData$absDiff<-abs(testingData$Actual-testingData$guess)
> testingData$percErr<-testingData$absDiff/testingData$Actual
> mean(testingData$absDiff)/mean(testingData$Actual)
[1] 0.2687390437
> #let's see if we improve if we undo the lat and longitude composite
> formula<-MedianHomeVal2~comp1+medInc+medAge+lat+long
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=8, threshold=0.05)
Warning message:
algorithm did not converge in 1 of 1 repetition(s) within the stepmax
> proc.time() - ptm
   user  system elapsed
2394.73  143.58 2541.89
> #took forever and didn't work
> #try again with more layers and higher threshold
> ptm <- proc.time()
> nnet<-neuralnet(formula,trainingData, hidden=10, threshold=0.09)
Warning message:
algorithm did not converge in 1 of 1 repetition(s) within the stepmax
> proc.time() - ptm
   user  system elapsed
2835.84  152.79 2992.83
> #still doesn't work.  Give up
```