

557_Project_2BS

Ben Straub, Hillary Koch, Jiawei Huang, Arif Masrur

3/15/2017

Boring Stuff on the dataset

Names of Variables

```
[1] "seismic"      "seismoacoustic" "shift"          "genergy"
[5] "gpuls"        "gdenergy"       "gdpuls"         "ghazard"
[9] "nbumps"       "nbumps2"        "nbumps3"        "nbumps4"
[13] "nbumps5"      "nbumps6"        "nbumps7"        "nbumps89"
[17] "energy"       "maxenergy"      "class"
```

Summary Statistics

```
seismic  seismoacoustic shift      genergy      gpuls
a:1682   a:1580          N: 921   Min.    :    100   Min.    :    2.0
b: 902   b: 956          W:1663   1st Qu.: 11660   1st Qu.: 190.0
        c: 48           Median : 25485   Median : 379.0
                        Mean    : 90242   Mean    : 538.6
                        3rd Qu.: 52832   3rd Qu.: 669.0
                        Max.    :2595650   Max.    :4518.0

      gdenergy      gdpuls      ghazard      nbumps
Min.    : -96.00   Min.    : -96.000   a:2342   Min.    :0.0000
1st Qu.: -37.00   1st Qu.: -36.000   b: 212   1st Qu.:0.0000
Median :  -6.00   Median :  -6.000   c: 30    Median :0.0000
Mean    : 12.38   Mean    :  4.509           Mean    :0.8595
3rd Qu.: 38.00   3rd Qu.: 30.250           3rd Qu.:1.0000
Max.    :1245.00   Max.    :838.000           Max.    :9.0000

      nbumps2      nbumps3      nbumps4      nbumps5
Min.    :0.0000   Min.    :0.0000   Min.    :0.00000   Min.    :0.000000
1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.00000   1st Qu.:0.000000
Median :0.0000   Median :0.0000   Median :0.00000   Median :0.000000
Mean    :0.3936   Mean    :0.3928   Mean    :0.06772   Mean    :0.004644
3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:0.00000   3rd Qu.:0.000000
Max.    :8.0000   Max.    :7.0000   Max.    :3.00000   Max.    :1.000000

      nbumps6      nbumps7      nbumps89      energy      maxenergy
Min.    :0   Min.    :0   Min.    :0   Min.    :    0   Min.    :    0
1st Qu.:0   1st Qu.:0   1st Qu.:0   1st Qu.:    0   1st Qu.:    0
Median :0   Median :0   Median :0   Median :    0   Median :    0
Mean    :0   Mean    :0   Mean    :0   Mean    : 4975   Mean    : 4279
3rd Qu.:0   3rd Qu.:0   3rd Qu.:0   3rd Qu.: 2600   3rd Qu.: 2000
Max.    :0   Max.    :0   Max.    :0   Max.    :402000   Max.    :400000

      class
Min.    :0.00000
1st Qu.:0.00000
Median :0.00000
Mean    :0.06579
```

3rd Qu.:0.00000
Max. :1.00000

Dimensions of Data Matrix

[1] 2584 19

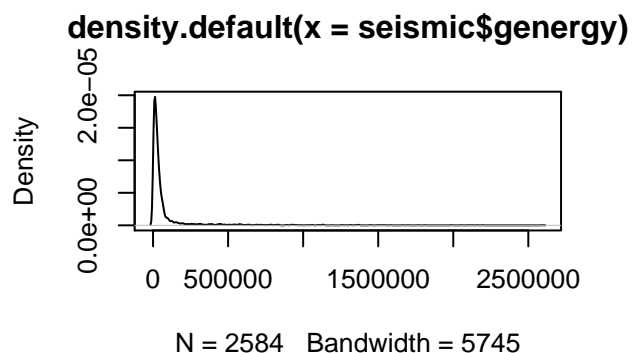
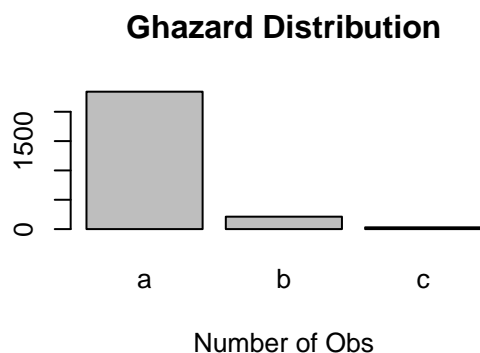
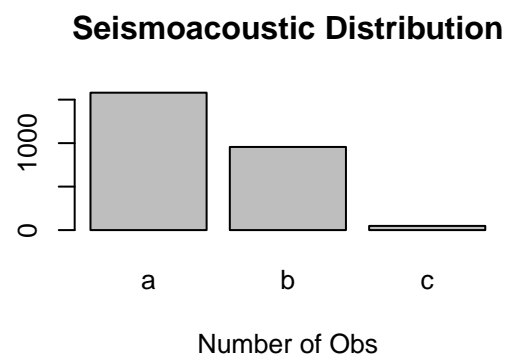
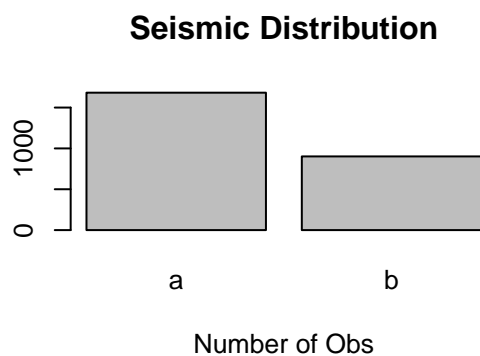
Check for Normality of Data

```
par(mfrow=c(2,2))

# Barplots of Factor Variables
counts <- table(seismic$seismic)
barplot(counts, main="Seismic Distribution",
        xlab="Number of Obs")
counts <- table(seismic$seismoacoustic)
barplot(counts, main="Seismoacoustic Distribution",
        xlab="Number of Obs")
counts <- table(seismic$ghazard)
barplot(counts, main="Ghazard Distribution",
        xlab="Number of Obs")

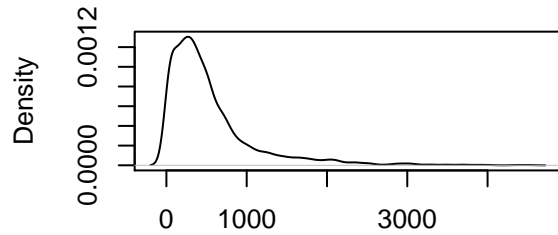
## Have a look at the densities

plot(density(seismic$genergy));plot(density(seismic$gpuls))
```



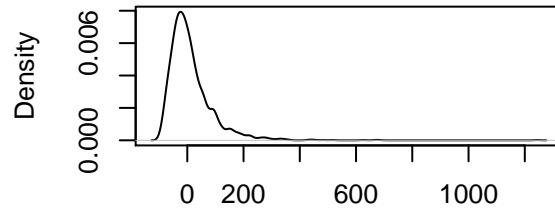
```
plot(density(seismic$gdenergy));plot(density(seismic$gdpuls))
plot(density(seismic$maxenergy));plot(density(seismic$nbumps))
```

density.default(x = seismic\$gdpuls)



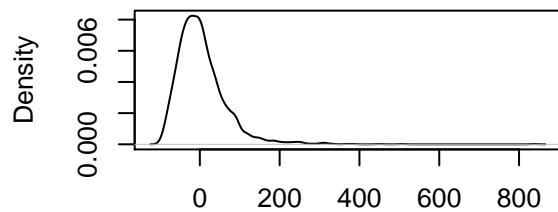
N = 2584 Bandwidth = 66.84

density.default(x = seismic\$gdenergy)



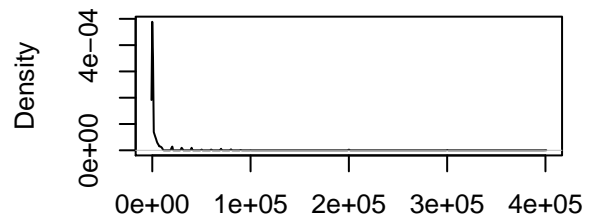
N = 2584 Bandwidth = 10.47

density.default(x = seismic\$gdpuls)



N = 2584 Bandwidth = 9.244

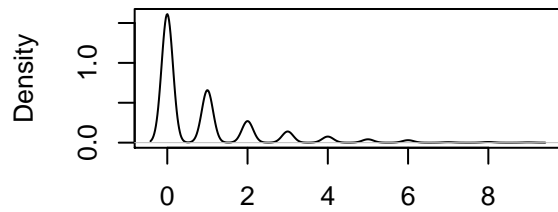
density.default(x = seismic\$maxenergy)



N = 2584 Bandwidth = 279.1

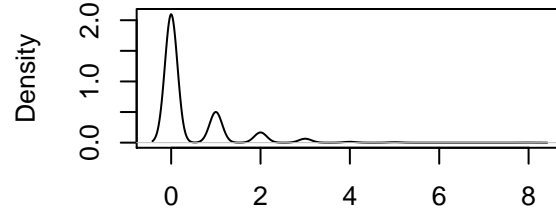
```
plot(density(seismic$nbumps2));plot(density(seismic$nbumps3))
plot(density(seismic$nbumps4));plot(density(seismic$nbumps5))
```

density.default(x = seismic\$nbumps)



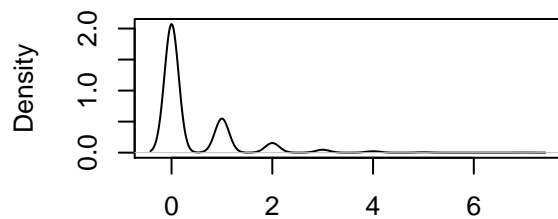
N = 2584 Bandwidth = 0.1395

density.default(x = seismic\$nbumps2)



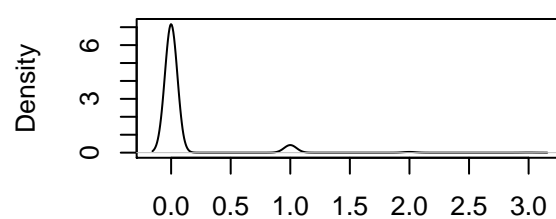
N = 2584 Bandwidth = 0.1395

density.default(x = seismic\$nbumps3)



N = 2584 Bandwidth = 0.1395

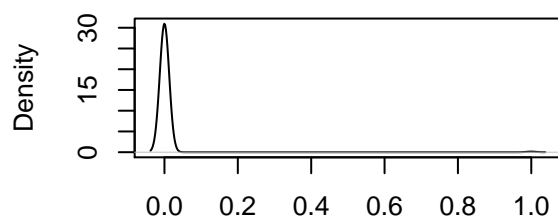
density.default(x = seismic\$nbumps4)



N = 2584 Bandwidth = 0.05218

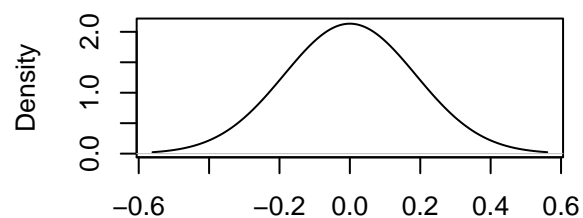
```
plot(density(seismic$nbumps6));plot(density(seismic$nbumps7))
plot(density(seismic$nbumps89));plot(density(seismic$energy))
```

density.default(x = seismic\$nbumps5)



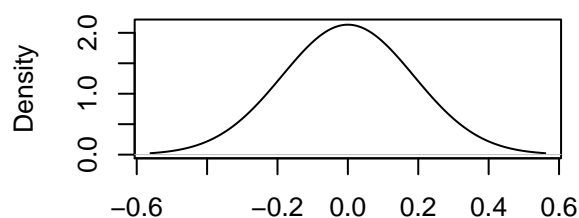
N = 2584 Bandwidth = 0.01271

density.default(x = seismic\$nbumps6)



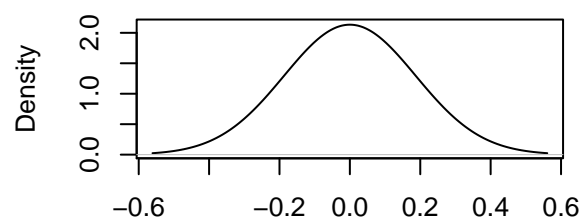
N = 2584 Bandwidth = 0.187

density.default(x = seismic\$nbumps7)



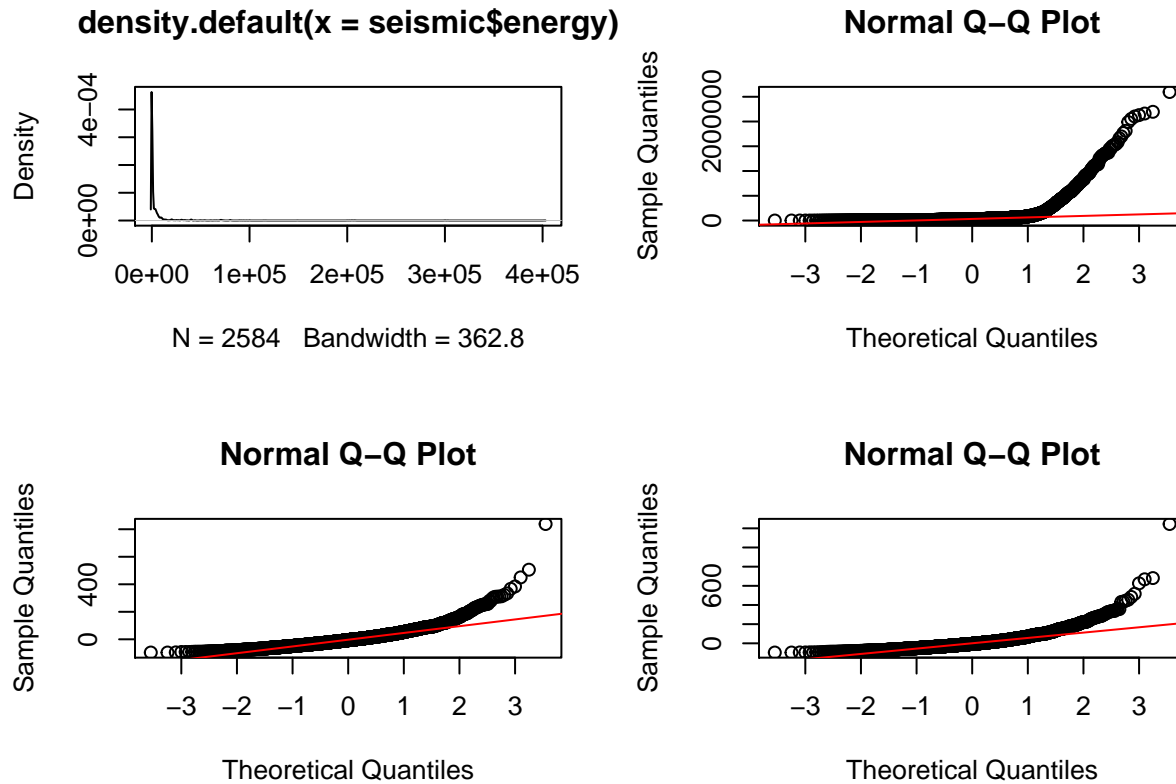
N = 2584 Bandwidth = 0.187

density.default(x = seismic\$nbumps89)

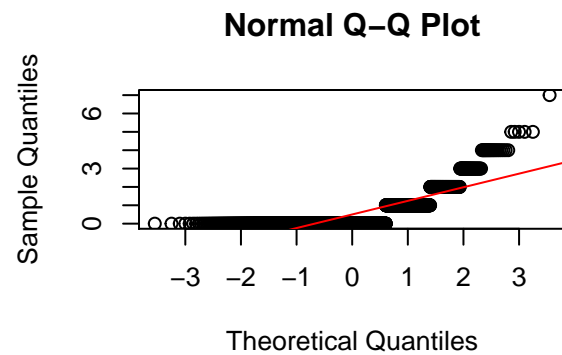
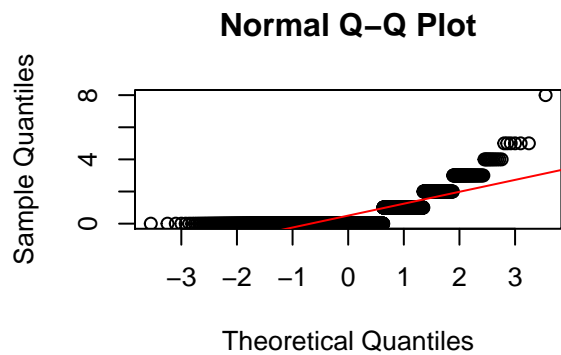
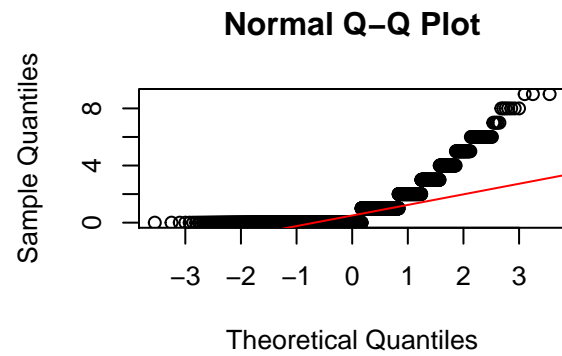
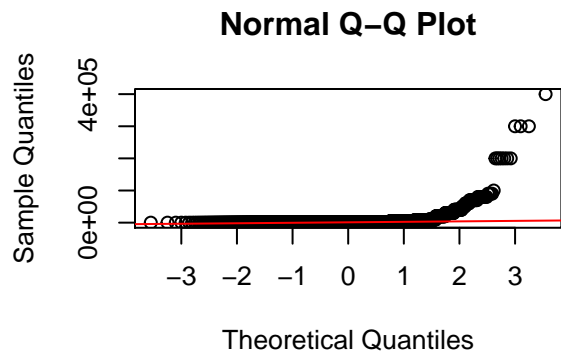


N = 2584 Bandwidth = 0.187

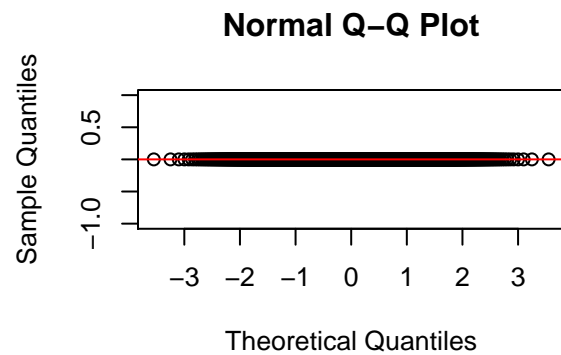
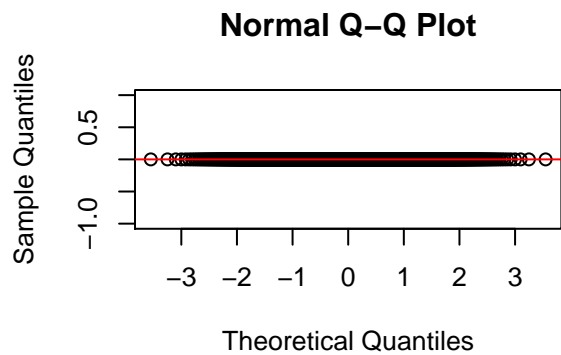
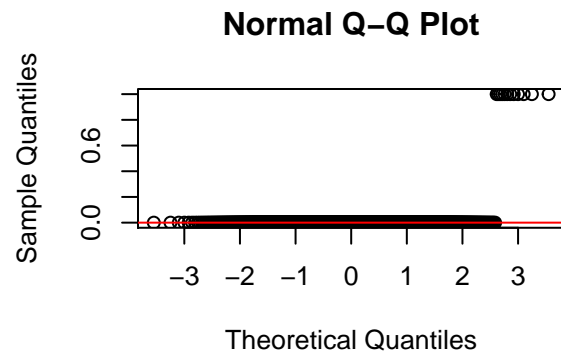
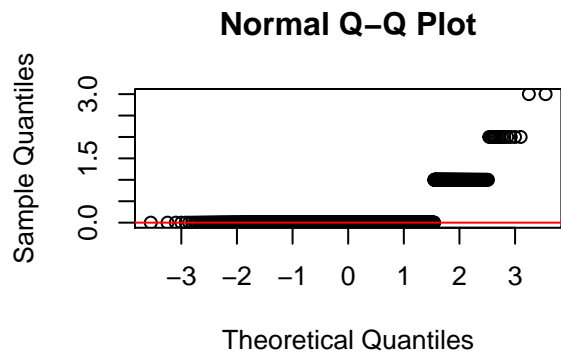
```
## Plot using a qqplot
qqnorm(seismic$energy);qqline(seismic$energy, col = 2)
qqnorm(seismic$gdpuls);qqline(seismic$gdpuls, col = 2)
qqnorm(seismic$gdenergy);qqline(seismic$gdenergy, col = 2)
```



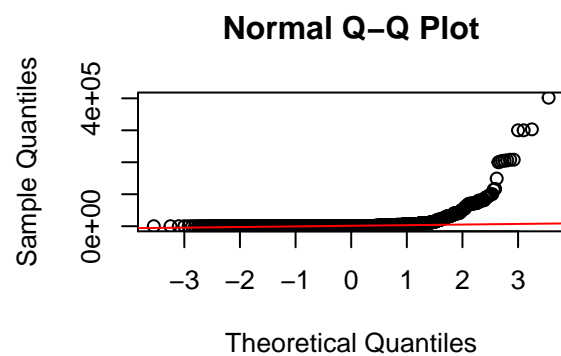
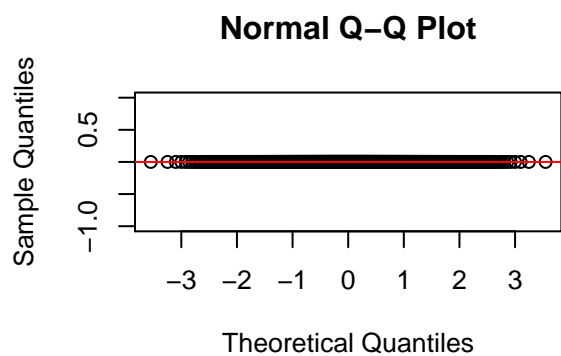
```
qqnorm(seismic$maxenergy);qqline(seismic$maxenergy, col = 2)
qqnorm(seismic$nbumps);qqline(seismic$nbumps, col = 2)
qqnorm(seismic$nbumps2);qqline(seismic$nbumps2, col = 2)
qqnorm(seismic$nbumps3);qqline(seismic$nbumps3, col = 2)
```



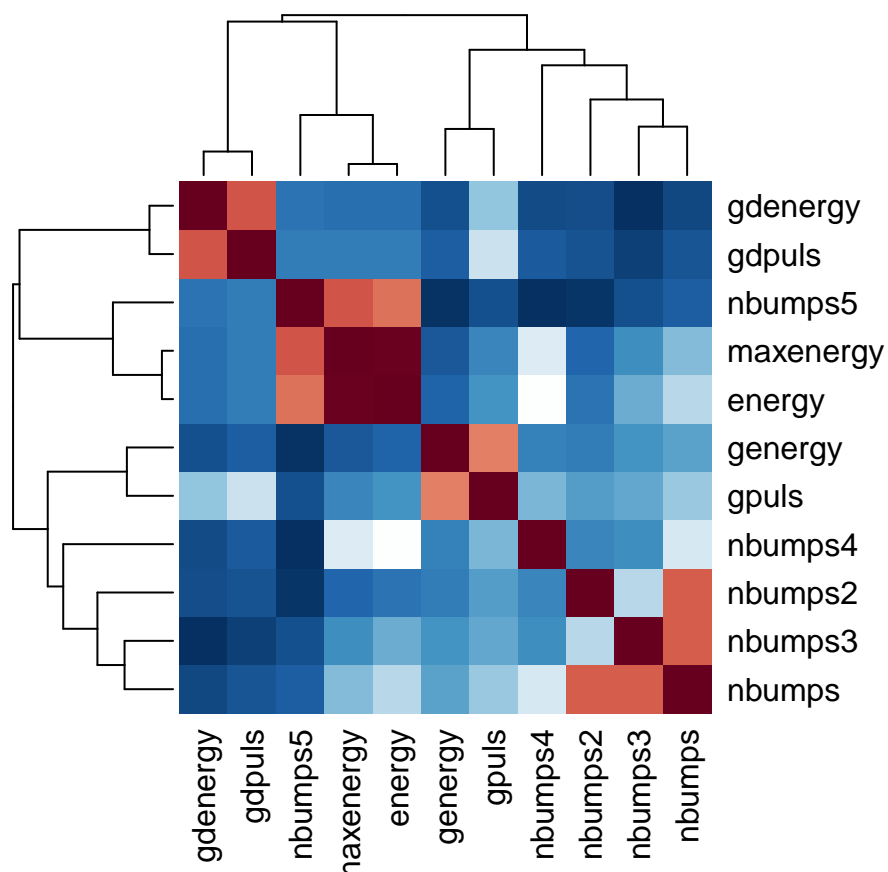
```
qqnorm(seismic$nbumps4);qqline(seismic$nbumps4, col = 2)
qqnorm(seismic$nbumps5);qqline(seismic$nbumps5, col = 2)
qqnorm(seismic$nbumps6);qqline(seismic$nbumps6, col = 2)
qqnorm(seismic$nbumps7);qqline(seismic$nbumps7, col = 2)
```



```
qqnorm(seismic$nbumps89);qqline(seismic$nbumps89, col = 2)
qqnorm(seismic$energy);qqline(seismic$energy, col = 2)
```



Correlation of the Variables

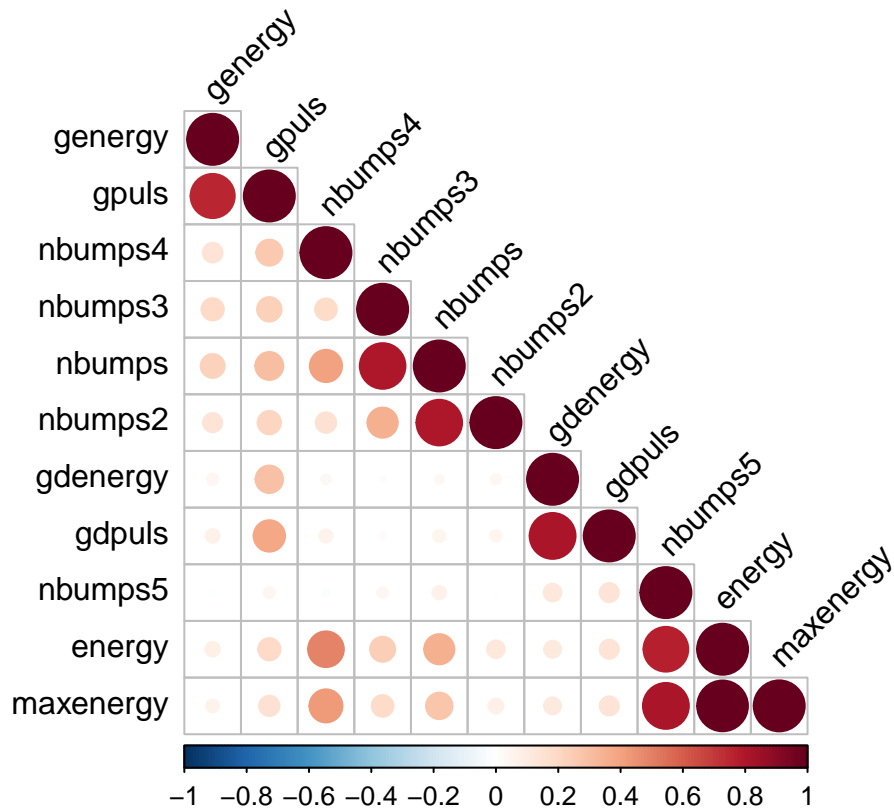


\$r	row	column	cor	p
1	genergy	gpuls	0.7500	0.0e+00
2	genergy	nbumps4	0.1500	1.4e-14
3	gpuls	nbumps4	0.2600	0.0e+00
4	genergy	nbumps3	0.1900	0.0e+00
5	gpuls	nbumps3	0.2300	0.0e+00
6	nbumps4	nbumps3	0.1800	0.0e+00
7	genergy	nbumps	0.2200	0.0e+00
8	gpuls	nbumps	0.3000	0.0e+00
9	nbumps4	nbumps	0.4000	0.0e+00
10	nbumps3	nbumps	0.8000	0.0e+00
11	genergy	nbumps2	0.1400	2.2e-13
12	gpuls	nbumps2	0.2100	0.0e+00
13	nbumps4	nbumps2	0.1600	0.0e+00
14	nbumps3	nbumps2	0.3500	0.0e+00
15	nbumps	nbumps2	0.8000	0.0e+00
16	genergy	gdenergy	0.0490	1.4e-02
17	gpuls	gdenergy	0.2900	0.0e+00
18	nbumps4	gdenergy	0.0370	6.1e-02
19	nbumps3	gdenergy	-0.0120	5.4e-01
20	nbumps	gdenergy	0.0300	1.3e-01
21	nbumps2	gdenergy	0.0410	3.6e-02
22	genergy	gdpuls	0.0720	2.7e-04

23	gpuls	gdpuls	0.3800	0.0e+00
24	nbumps4	gdpuls	0.0660	7.6e-04
25	nbumps3	gdpuls	0.0150	4.5e-01
26	nbumps	gdpuls	0.0580	3.2e-03
27	nbumps2	gdpuls	0.0510	9.4e-03
28	gdenergy	gdpuls	0.8100	0.0e+00
29	genergy	nbumps5	-0.0099	6.2e-01
30	gpuls	nbumps5	0.0490	1.2e-02
31	nbumps4	nbumps5	-0.0170	4.0e-01
32	nbumps3	nbumps5	0.0460	1.8e-02
33	nbumps	nbumps5	0.0700	4.0e-04
34	nbumps2	nbumps5	-0.0053	7.9e-01
35	gdenergy	nbumps5	0.1200	3.3e-10
36	gdpuls	nbumps5	0.1400	5.9e-13
37	genergy	energy	0.0810	3.9e-05
38	gpuls	energy	0.1900	0.0e+00
39	nbumps4	energy	0.4900	0.0e+00
40	nbumps3	energy	0.2400	0.0e+00
41	nbumps	energy	0.3500	0.0e+00
42	nbumps2	energy	0.1200	2.0e-10
43	gdenergy	energy	0.1100	6.7e-08
44	gdpuls	energy	0.1400	2.5e-13
45	nbumps5	energy	0.7700	0.0e+00
46	genergy	maxenergy	0.0640	1.1e-03
47	gpuls	maxenergy	0.1600	0.0e+00
48	nbumps4	maxenergy	0.4200	0.0e+00
49	nbumps3	maxenergy	0.1800	0.0e+00
50	nbumps	maxenergy	0.2700	0.0e+00
51	nbumps2	maxenergy	0.0850	1.5e-05
52	gdenergy	maxenergy	0.1100	3.2e-08
53	gdpuls	maxenergy	0.1400	2.2e-13
54	nbumps5	maxenergy	0.8100	0.0e+00
55	energy	maxenergy	0.9900	0.0e+00

\$p
NULL

\$sym
NULL



\$r	row	column	cor	p
1	genergy	gpuls	0.7500	0.0e+00
2	genergy	nbumps4	0.1500	1.4e-14
3	gpuls	nbumps4	0.2600	0.0e+00
4	genergy	nbumps3	0.1900	0.0e+00
5	gpuls	nbumps3	0.2300	0.0e+00
6	nbumps4	nbumps3	0.1800	0.0e+00
7	genergy	nbumps	0.2200	0.0e+00
8	gpuls	nbumps	0.3000	0.0e+00
9	nbumps4	nbumps	0.4000	0.0e+00
10	nbumps3	nbumps	0.8000	0.0e+00
11	genergy	nbumps2	0.1400	2.2e-13
12	gpuls	nbumps2	0.2100	0.0e+00
13	nbumps4	nbumps2	0.1600	0.0e+00
14	nbumps3	nbumps2	0.3500	0.0e+00
15	nbumps	nbumps2	0.8000	0.0e+00
16	genergy	gdenergy	0.0490	1.4e-02
17	gpuls	gdenergy	0.2900	0.0e+00
18	nbumps4	gdenergy	0.0370	6.1e-02
19	nbumps3	gdenergy	-0.0120	5.4e-01
20	nbumps	gdenergy	0.0300	1.3e-01
21	nbumps2	gdenergy	0.0410	3.6e-02
22	genergy	gdpuls	0.0720	2.7e-04
23	gpuls	gdpuls	0.3800	0.0e+00
24	nbumps4	gdpuls	0.0660	7.6e-04
25	nbumps3	gdpuls	0.0150	4.5e-01
26	nbumps	gdpuls	0.0580	3.2e-03

27	nbumps2	gdpuls	0.0510	9.4e-03
28	gdenergy	gdpuls	0.8100	0.0e+00
29	genergy	nbumps5	-0.0099	6.2e-01
30	gpuls	nbumps5	0.0490	1.2e-02
31	nbumps4	nbumps5	-0.0170	4.0e-01
32	nbumps3	nbumps5	0.0460	1.8e-02
33	nbumps	nbumps5	0.0700	4.0e-04
34	nbumps2	nbumps5	-0.0053	7.9e-01
35	gdenergy	nbumps5	0.1200	3.3e-10
36	gdpuls	nbumps5	0.1400	5.9e-13
37	genergy	energy	0.0810	3.9e-05
38	gpuls	energy	0.1900	0.0e+00
39	nbumps4	energy	0.4900	0.0e+00
40	nbumps3	energy	0.2400	0.0e+00
41	nbumps	energy	0.3500	0.0e+00
42	nbumps2	energy	0.1200	2.0e-10
43	gdenergy	energy	0.1100	6.7e-08
44	gdpuls	energy	0.1400	2.5e-13
45	nbumps5	energy	0.7700	0.0e+00
46	genergy	maxenergy	0.0640	1.1e-03
47	gpuls	maxenergy	0.1600	0.0e+00
48	nbumps4	maxenergy	0.4200	0.0e+00
49	nbumps3	maxenergy	0.1800	0.0e+00
50	nbumps	maxenergy	0.2700	0.0e+00
51	nbumps2	maxenergy	0.0850	1.5e-05
52	gdenergy	maxenergy	0.1100	3.2e-08
53	gdpuls	maxenergy	0.1400	2.2e-13
54	nbumps5	maxenergy	0.8100	0.0e+00
55	energy	maxenergy	0.9900	0.0e+00

\$p
NULL

\$sym
NULL

Logistic Regression on the Training and Test Sets

Call:

```
glm(formula = y.train ~ ., family = binomial, data = seismic.train)
```

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-3.6758	-0.8347	-0.5605	0.9798	3.0316

Coefficients: (3 not defined because of singularities)

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.982e+00	1.147e-01	-17.278	< 2e-16 ***
shiftW	7.341e-01	1.369e-01	5.363	8.19e-08 ***
genergy	-4.212e-06	5.599e-07	-7.522	5.39e-14 ***
gpuls	1.864e-03	2.329e-04	8.002	1.22e-15 ***
gdenergy	2.632e-03	1.141e-03	2.307	0.02106 *

```

gdpuls      -4.068e-03  1.578e-03  -2.578  0.00995 **
ghazardb     7.853e-01  1.918e-01   4.093  4.25e-05 ***
ghazardc    -1.245e+00  6.014e-01  -2.070  0.03844 *
nbumps      -1.119e+01  3.247e+02  -0.034  0.97250
nbumps2     1.105e+01  3.247e+02   0.034  0.97287
nbumps3     1.112e+01  3.247e+02   0.034  0.97269
nbumps4     1.279e+01  3.247e+02   0.039  0.96857
nbumps5     8.235e+00  3.248e+02   0.025  0.97977
nbumps6              NA          NA      NA      NA
nbumps7              NA          NA      NA      NA
nbumps89             NA          NA      NA      NA
energy       1.396e-05  5.053e-05   0.276  0.78233
maxenergy    -9.175e-07  5.006e-05  -0.018  0.98538
class        4.398e-01  2.133e-01   2.062  0.03918 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

(Dispersion parameter for binomial family taken to be 1)

```

Null deviance: 2480.8  on 1937  degrees of freedom
Residual deviance: 2080.5  on 1922  degrees of freedom
AIC: 2112.5

```

Number of Fisher Scoring iterations: 11

The predictors that are significant in our logistic model are genergy, gpuls and ghazardb and a couple more. The predictors nbumps6, nbumps7 and nbumps89 are not defined due to singularities, which may indicated collinearity.

```

      y.train
glm.pred  a    b
a 1161  416
b  121  240

```

```
[1] 0.7229102
```

The diagonal elements of the confusion matrix indicate correct predictions, while the off-diagonals represent incorrect predictions. Hence our model on the training data set correctly predicted that the seismic activity would be of no harzard on 1176 observations and that it would be a low hazard on 230 observations, for a total of $1176 + 230 = 1406$ correct predictions. The `mean()` function can be used to compute the fraction of seismic activity for which the prediction was correct. In this case, logistic regression correctly predicted the movement of the market 73 percent of the time.

```

##      y.test
## glm.pred  a    b
##      a 345 130
##      b  55 116

```

```
## [1] 0.7136223
```

The diagonal elements of the confusion matrix indicate correct predictions, while the off-diagonals represent incorrect predictions. Hence our model on the testing data set correctly predicted that the seismic activity

would be of no hazard on 352 observations and that it would be a low hazard on 110 days, for a total of $352 + 110 = 462$ correct predictions. The `mean()` function can be used to compute the fraction of seismic activity for which the prediction was correct. In this case, logistic regression correctly predicted the movement of the market 71.5 percent of the time.

Recall that the logistic regression model had only 7ish predictors that were significant from an available 17. Perhaps by removing the variables that appear not to be helpful in predicting seismic hazard, we can obtain a more effective model. After all, using predictors that have no relationship with the response tends to cause a deterioration in the test error rate (since such predictors cause an increase in variance without a corresponding decrease in bias), and so removing such predictors may in turn yield an improvement [straight from the book]