

Predicting Stock Returns

Contents

Introduction	1
Exercises	2
1. Exploring the Dataset	2
Problem 1.1	2
Problem 1.2	3
Problem 1.3	3
Problem 1.4	4
2. Initial Logistic Regression Model	4
Problem 2.1	4
Problem 2.2	5
Problem 2.3	5
3. Clustering Stocks	5
Problem 3.2	6
Problem 3.3	6
Problem 3.4	6
Problem 3.5	7
4. Cluster-Specific Predictions	7
Problem 4.1	7
Problem 4.2	8
Problem 4.3	10
Problem 4.4	11

Introduction

In the second lecture sequence this week, we heard about cluster-then-predict, a methodology in which you first cluster observations and then build cluster-specific prediction models. In the lecture sequence, we saw how this methodology helped improve the prediction of heart attack risk. In this assignment, we'll use cluster-then-predict to predict future stock prices using historical stock data.

When selecting which stocks to invest in, investors seek to obtain good future returns. In this problem, we will first use clustering to identify clusters of stocks that have similar returns over time. Then, we'll use logistic regression to predict whether or not the stocks will have positive future returns.

For this problem, we'll use `StocksCluster.csv`, which contains monthly stock returns from the NASDAQ stock exchange. The NASDAQ is the second-largest stock exchange in the world, and it lists many technology companies. The stock price data used in this problem was obtained from `infochimps`, a website providing access to many datasets.

Each observation in the dataset is the monthly returns of a particular company in a particular year. The years included are 2000-2009. The companies are limited to tickers that were listed on the exchange for the entire period 2000-2009, and whose stock price never fell below \$1. So, for example, one observation is for Yahoo in 2000, and another observation is for Yahoo in 2001. Our goal will be to predict whether or not the stock return in December will be positive, using the stock returns for the first 11 months of the year.

This dataset contains the following variables:

- **ReturnJan** : the return for the company's stock during January (in the year of the observation).
- **ReturnFeb** : the return for the company's stock during February (in the year of the observation).
- **ReturnMar** : the return for the company's stock during March (in the year of the observation).
- **ReturnApr** : the return for the company's stock during April (in the year of the observation).
- **ReturnMay** : the return for the company's stock during May (in the year of the observation).
- **ReturnJune** : the return for the company's stock during June (in the year of the observation).
- **ReturnJuly** : the return for the company's stock during July (in the year of the observation).
- **ReturnAug** : the return for the company's stock during August (in the year of the observation).
- **ReturnSep** : the return for the company's stock during September (in the year of the observation).
- **ReturnOct** : the return for the company's stock during October (in the year of the observation).
- **ReturnNov** : the return for the company's stock during November (in the year of the observation).
- **PositiveDec** : whether or not the company's stock had a positive return in December (in the year of the observation). This variable takes value 1 if the return was positive, and value 0 if the return was not positive.

For the first 11 variables, the value stored is a proportional change in stock value during that month. For instance, a value of 0.05 means the stock increased in value 5% during the month, while a value of -0.02 means the stock decreased in value 2% during the month.

Exercices

1. Exploring the Dataset

Problem 1.1 Load `StocksCluster.csv` into a data frame called "stocks".

```
## 'data.frame': 11580 obs. of 12 variables:
## $ ReturnJan : num 0.0807 -0.0107 0.0477 -0.074 -0.031 ...
## $ ReturnFeb : num 0.0663 0.1021 0.036 -0.0482 -0.2127 ...
## $ ReturnMar : num 0.0329 0.1455 0.0397 0.0182 0.0915 ...
## $ ReturnApr : num 0.1831 -0.0844 -0.1624 -0.0247 0.1893 ...
## $ ReturnMay : num 0.13033 -0.3273 -0.14743 -0.00604 -0.15385 ...
## $ ReturnJune : num -0.0176 -0.3593 0.0486 -0.0253 -0.1061 ...
## $ ReturnJuly : num -0.0205 -0.0253 -0.1354 -0.094 0.3553 ...
## $ ReturnAug : num 0.0247 0.2113 0.0334 0.0953 0.0568 ...
## $ ReturnSep : num -0.0204 -0.58 0 0.0567 0.0336 ...
## $ ReturnOct : num -0.1733 -0.2671 0.0917 -0.0963 0.0363 ...
## $ ReturnNov : num -0.0254 -0.1512 -0.0596 -0.0405 -0.0853 ...
## $ PositiveDec: int 0 0 0 1 1 1 1 0 0 0 ...
```

How many observations are in the dataset?

Answer : 11580

Explanation :

You can load the dataset with the read.csv function:

And see how many observations are included with either the str or nrow function:

Both tell us that there are 11580 observations in this dataset.

Problem 1.2 What proportion of the observations have positive returns in December?

```
## [1] 0.546114
```

Answer : 0.546114

Explanation :

You can compute the proportion of observations with positive returns by using the table function:

It tells us that 6324 observations have PositiveDec = 1, so $6324/11580 = 0.546$ of the observations have positive returns in December.

Alternatively, you could use the mean function to compute the proportion:

Problem 1.3 What is the maximum correlation between any two return variables in the dataset?

You should look at the pairwise correlations between ReturnJan, ReturnFeb, ReturnMar, ReturnApr, ReturnMay, ReturnJune, ReturnJuly, ReturnAug, ReturnSep, ReturnOct, and ReturnNov.

```
##          ReturnJan ReturnFeb ReturnMar ReturnApr ReturnMay
## ReturnJan  1.000000000  0.06677458 -0.090496798 -0.037678006 -0.044411417
## ReturnFeb  0.066774583  1.000000000 -0.155983263 -0.191351924 -0.095520920
## ReturnMar -0.090496798 -0.15598326  1.000000000  0.009726288 -0.003892789
## ReturnApr -0.037678006 -0.19135192  0.009726288  1.000000000  0.063822504
## ReturnMay -0.044411417 -0.09552092 -0.003892789  0.063822504  1.000000000
## ReturnJune  0.092238307  0.16999448 -0.085905486 -0.011027752 -0.021074539
## ReturnJuly -0.081429765 -0.06177851  0.003374160  0.080631932  0.090850264
## ReturnAug -0.022792019  0.13155979 -0.022005400 -0.051756051 -0.033125658
## ReturnSep -0.026437153  0.04350177  0.076518327 -0.028920972  0.021962862
## ReturnOct  0.142977229 -0.08732427 -0.011923758  0.048540025  0.017166728
## ReturnNov  0.067632333 -0.15465828  0.037323535  0.031761837  0.048046590
## PositiveDec 0.004728518 -0.03817318  0.022408661  0.094353528  0.058201934
##          ReturnJune ReturnJuly ReturnAug ReturnSep ReturnOct
## ReturnJan  0.09223831 -0.081429765 -0.0227920187 -0.0264371526  0.14297723
## ReturnFeb  0.16999448 -0.0617785094  0.1315597863  0.0435017706 -0.08732427
## ReturnMar -0.08590549  0.0033741597 -0.0220053995  0.0765183267 -0.01192376
## ReturnApr -0.01102775  0.0806319317 -0.0517560510 -0.0289209718  0.04854003
## ReturnMay -0.02107454  0.0908502642 -0.0331256580  0.0219628623  0.01716673
## ReturnJune  1.000000000 -0.0291525996  0.0107105260  0.0447472692 -0.02263599
## ReturnJuly -0.02915260  1.00000000000  0.0007137558  0.0689478037 -0.05470891
## ReturnAug  0.01071053  0.0007137558  1.00000000000  0.0007407139 -0.07559456
## ReturnSep  0.04474727  0.0689478037  0.0007407139  1.00000000000 -0.05807924
## ReturnOct -0.02263599 -0.0547089088 -0.0755945614 -0.0580792362  1.000000000
## ReturnNov -0.06527054 -0.0483738369 -0.1164890345 -0.0197197998  0.19167279
## PositiveDec 0.02340975  0.0743642097  0.0041669657  0.0416302863 -0.05257496
```

```
##          ReturnNov PositiveDec
## ReturnJan  0.06763233  0.004728518
## ReturnFeb -0.15465828 -0.038173184
## ReturnMar  0.03732353  0.022408661
## ReturnApr  0.03176184  0.094353528
## ReturnMay  0.04804659  0.058201934
## ReturnJune -0.06527054  0.023409745
## ReturnJuly -0.04837384  0.074364210
## ReturnAug -0.11648903  0.004166966
## ReturnSep -0.01971980  0.041630286
## ReturnOct  0.19167279 -0.052574956
## ReturnNov  1.00000000 -0.062346556
## PositiveDec -0.06234656  1.000000000
```

Answer : 0.19167279 (between ReturnOct and ReturnNov)

Explanation :

From :

We see the largest correlation coefficient is 0.19167279, between ReturnOct and ReturnNov.

Problem 1.4 Which month (from January through November) has the largest mean return across all observations in the dataset?

```
##      ReturnJan  ReturnFeb  ReturnMar  ReturnApr  ReturnMay  ReturnJune
## 0.012631602 -0.007604784  0.019402336  0.026308147  0.024736591  0.005937902
##      ReturnJuly  ReturnAug  ReturnSep  ReturnOct  ReturnNov  PositiveDec
## 0.003050863  0.016198265 -0.014720768  0.005650844  0.011387440  0.546113990
```

Answer : April

Which month (from January through November) has the smallest mean return across all observations in the dataset?

Answer : September

Explanation :

These can be determined using the summary function :

If you look at the mean value for each variable, you can see that April has the largest mean value (0.026308), and September has the smallest mean value (-0.014721).

2. Initial Logistic Regression Model

Problem 2.1 Run the following commands to split the data into a training set and testing set, putting 70% of the data in the training set and 30% of the data in the testing set:

Then, use the stocksTrain data frame to train a logistic regression model (name it StocksModel) to predict PositiveDec using all the other variables as independent variables. Don't forget to add the argument family=binomial to your glm command.

What is the overall accuracy on the training set, using a threshold of 0.5?

```
## [1] 0.5711818
```

Answer : 0.5711818

Explanation :

We can train the model with:

Then, we can compute our predictions on the training set with:

And construct a classification matrix with the table function:

The overall accuracy of the model is $(990 + 3640)/(990 + 2689 + 787 + 3640) = 0.571$.

Problem 2.2 Now obtain test set predictions from StocksModel.

What is the overall accuracy of the model on the test, again using a threshold of 0.5?

```
## [1] 0.5670697
```

Answer : 0.5670697

Explanation :

You can compute predictions on the test set using the predict function:

Then, you can compute the classification matrix on the test set with the table function:

The overall accuracy of the model on the test set is $(417 + 1553)/(417 + 1160 + 344 + 1553) = 0.567$

Problem 2.3 What is the accuracy on the test set of a baseline model that always predicts the most common outcome (PositiveDec = 1)?

```
##  
##      FALSE TRUE  
##    0    417 1160  
##    1    344 1553
```

```
## [1] 0.5460564
```

Answer : 0.5460564

Explanation :

This can be computed by making a table of the outcome variable in the test set:

The baseline model would get all of the PositiveDec = 1 cases correct, and all of the PositiveDec = 0 cases wrong, for an accuracy of $1897/(1577 + 1897) = 0.5460564$.

3. Clustering Stocks

Now, let's cluster the stocks. The first step in this process is to remove the dependent variable using the following commands:

Why do we need to remove the dependent variable in the clustering phase of the cluster-then-predict methodology?

1. Leaving in the dependent variable might lead to unbalanced clusters
2. Removing the dependent variable decreases the computational effort needed to cluster
3. **Needing to know the dependent variable value to assign an observation to a cluster defeats the purpose of the methodology**

Explanation :

In cluster-then-predict, our final goal is to predict the dependent variable, which is unknown to us at the time of prediction. Therefore, if we need to know the outcome value to perform the clustering, the methodology is no longer useful for prediction of an unknown outcome value.

This is an important point that is sometimes mistakenly overlooked. If you use the outcome value to cluster, you might conclude your method strongly outperforms a non-clustering alternative. However, this is because it is using the outcome to determine the clusters, which is not valid.

Problem 3.2 In the market segmentation assignment in this week's homework, you were introduced to the `preProcess` command from the `caret` package, which normalizes variables by subtracting by the mean and dividing by the standard deviation.

In cases where we have a training and testing set, we'll want to normalize by the mean and standard deviation of the variables in the training set. We can do this by passing just the training set to the `preProcess` function:

What is the mean of the `ReturnJan` variable in `normTrain`?

```
## [1] 2.100586e-17
```

Answer : 2.100586e-17

What is the mean of the `ReturnJan` variable in `normTest`?

```
## [1] -0.0004185886
```

Answer : -0.0004185886

Explanation :

After running the provided normalization commands, we can read the means with

Problem 3.3 Why is the mean `ReturnJan` variable much closer to 0 in `normTrain` than in `normTest`?

1. Small rounding errors exist in the normalization procedure
2. **The distribution of the `ReturnJan` variable is different in the training and testing set**
3. The distribution of the dependent variable is different in the training and testing set

Explanation :

From :

We see that the average return in January is slightly higher in the training set than in the testing set. Since `normTest` was constructed by subtracting by the mean `ReturnJan` value from the training set, this explains why the mean value of `ReturnJan` is slightly negative in `normTest`.

Problem 3.4 Set the random seed to 144 (it is important to do this again, even though we did it earlier). Run k-means clustering with 3 clusters on `normTrain`, storing the result in an object called `km`.

Which cluster has the largest number of observations?

```
##  
##      1      2      3  
## 2479 4731  896
```

1. Cluster 1
2. **Cluster 2**
3. Cluster 3

Explanation :

We can set the seed and run the k-means algorithm with:

From

We can see that cluster 2 has the largest number of observations. Alternatively, you can see the number of observations in each cluster by typing `km$size` in your console.

Problem 3.5 Recall from the recitation that we can use the `flexclust` package to obtain training set and testing set cluster assignments for our observations (note that the call to `as.kcca` may take a while to complete):

How many test-set observations were assigned to Cluster 2?

```
## [1] 2029
```

Answer : 2080

Explanation :

After running the provided commands, we can obtain the breakdown of the testing set clusters with

4. Cluster-Specific Predictions

Problem 4.1 Using the `subset` function, build data frames `stocksTrain1`, `stocksTrain2`, and `stocksTrain3`, containing the elements in the `stocksTrain` data frame assigned to clusters 1, 2, and 3, respectively (be careful to take subsets of `stocksTrain`, not of `normTrain`). Similarly build `stocksTest1`, `stocksTest2`, and `stocksTest3` from the `stocksTest` data frame.

Which training set data frame has the highest average value of the dependent variable?

```
## [1] 0.6103267
```

```
## [1] 0.5250476
```

```
## [1] 0.4799107
```

```
## [[1]]
## [1] 0.6103267
##
## [[2]]
## [1] 0.5250476
##
## [[3]]
## [1] 0.4799107
```

1. **stocksTrain1**
2. `stocksTrain2`
3. `stocksTrain3`

Explanation :

We can obtain the necessary subsets with:

From:

We see that `stocksTrain1` has the observations with the highest average value of the dependent variable.

Problem 4.2 Build logistic regression models `StocksModel1`, `StocksModel2`, and `StocksModel3`, which predict `PositiveDec` using all the other variables as independent variables. `StocksModel1` should be trained on `stocksTrain1`, `StocksModel2` should be trained on `stocksTrain2`, and `StocksModel3` should be trained on `stocksTrain3`.

```
## [[1]]
##
## Call:
## glm(formula = PositiveDec ~ ., family = "binomial", data = stocksTrainKM[[x]])
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.7220  -1.2879   0.8679   1.0096   1.7170
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.20739    0.08020   2.586  0.00971 **
## ReturnJan    0.12448    0.31474   0.396  0.69246
## ReturnFeb   -0.46307    0.32713  -1.416  0.15692
## ReturnMar    0.55465    0.24804   2.236  0.02534 *
## ReturnApr    1.08354    0.25005   4.333 1.47e-05 ***
## ReturnMay    0.30487    0.24993   1.220  0.22253
## ReturnJune   0.00172    0.33525   0.005  0.99591
## ReturnJuly  -0.02763    0.30216  -0.091  0.92714
## ReturnAug    0.40299    0.34570   1.166  0.24373
## ReturnSep    0.70779    0.32611   2.170  0.02998 *
## ReturnOct   -1.33254    0.29055  -4.586 4.51e-06 ***
## ReturnNov   -0.78944    0.30583  -2.581  0.00984 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 3314.9  on 2478  degrees of freedom
## Residual deviance: 3244.0  on 2467  degrees of freedom
## AIC: 3268
##
## Number of Fisher Scoring iterations: 4
##
##
## [[2]]
##
## Call:
## glm(formula = PositiveDec ~ ., family = "binomial", data = stocksTrainKM[[x]])
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```



```

## -2.2268 -1.2086 0.9698 1.1294 1.6769
##
## Coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.1165255  0.0335061   3.478 0.000506 ***
## ReturnJan    0.4489316  0.2414090   1.860 0.062938 .
## ReturnFeb   -0.1385458  0.1514734  -0.915 0.360373
## ReturnMar    0.4657754  0.2442129   1.907 0.056488 .
## ReturnApr    0.7839726  0.2558668   3.064 0.002184 **
## ReturnMay    0.7709831  0.2625217   2.937 0.003316 **
## ReturnJune   0.5764584  0.2191809   2.630 0.008537 **
## ReturnJuly   0.8654737  0.2869778   3.016 0.002563 **
## ReturnAug    0.0177313  0.2317020   0.077 0.939001
## ReturnSep    1.0464947  0.2684133   3.899 9.67e-05 ***
## ReturnOct   -0.0001062  0.2426989   0.000 0.999651
## ReturnNov   -0.3716212  0.2603210  -1.428 0.153421
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 6546.7  on 4730  degrees of freedom
## Residual deviance: 6481.4  on 4719  degrees of freedom
## AIC: 6505.4
##
## Number of Fisher Scoring iterations: 4
##
##
## [[3]]
##
## Call:
## glm(formula = PositiveDec ~ ., family = "binomial", data = stocksTrainKM[[x]])
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8938 -1.0863 -0.5241  1.0874  2.1892
##
## Coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.4504     0.1700   2.649 0.008071 **
## ReturnJan    0.3816     0.2742   1.392 0.163927
## ReturnFeb    0.3943     0.4612   0.855 0.392490
## ReturnMar   -1.7220     0.4372  -3.938 8.20e-05 ***
## ReturnApr    0.5214     0.3275   1.592 0.111316
## ReturnMay    1.1301     0.4274   2.644 0.008194 **
## ReturnJune   1.5889     0.4423   3.592 0.000328 ***
## ReturnJuly   1.3800     0.4602   2.999 0.002709 **
## ReturnAug    0.4146     0.4824   0.860 0.390054
## ReturnSep   -0.0148     0.4754  -0.031 0.975167
## ReturnOct   -0.6820     0.3044  -2.241 0.025039 *
## ReturnNov   -1.7175     0.4133  -4.156 3.24e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##

```

```
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1240.7  on 895  degrees of freedom
## Residual deviance: 1131.5  on 884  degrees of freedom
## AIC: 1155.5
##
## Number of Fisher Scoring iterations: 4
```

Which variables have a positive sign for the coefficient in at least one of `StocksModel1`, `StocksModel2`, and `StocksModel3` and a negative sign for the coefficient in at least one of `StocksModel1`, `StocksModel2`, and `StocksModel3`?

Select all that apply.

1. **ReturnJan**
2. **ReturnFeb**
3. **ReturnMar**
4. ReturnApr
5. ReturnMay
6. **ReturnJune**
7. ReturnJuly
8. **ReturnAug**
9. ReturnSep
10. **ReturnOct**
11. ReturnNov

Explanation :

We can build the models with:

From:

ReturnJan, ReturnFeb, ReturnMar, ReturnJune, ReturnAug, and ReturnOct differ in sign between the models.

Problem 4.3 Using `StocksModel1`, make test-set predictions called `PredictTest1` on the data frame `stocksTest1`. Using `StocksModel2`, make test-set predictions called `PredictTest2` on the data frame `stocksTest2`. Using `StocksModel3`, make test-set predictions called `PredictTest3` on the data frame `stocksTest3`.

```
## [[1]]
## [1] 0.6446125
##
## [[2]]
## [1] 0.5367176
##
## [[3]]
## [1] 0.625323
```

What is the overall accuracy of `StocksModel1` on the test set `stocksTest1`, using a threshold of 0.5?

Answer : 0.6194145

What is the overall accuracy of `StocksModel2` on the test set `stocksTest2`, using a threshold of 0.5?

Answer : 0.5504808

What is the overall accuracy of StocksModel3 on the test set stocksTest3, using a threshold of 0.5?

Answer : 0.6458333

Explanation :

The predictions can be obtained with:

And the classification matrices can be computed with:

The overall accuracy of StocksModel1 is $(30 + 774)/(30 + 471 + 23 + 774) = 0.6194145$, the overall accuracy of StocksModel2 is $(388 + 757)/(388 + 626 + 309 + 757) = 0.5504808$, and the overall accuracy of StocksModel3 is $(49 + 13)/(49 + 13 + 21 + 13) = 0.6458333$.

Problem 4.4 To compute the overall test-set accuracy of the cluster-then-predict approach, we can combine all the test-set predictions into a single vector and all the true outcomes into a single vector:

What is the overall test-set accuracy of the cluster-then-predict approach, again using a threshold of 0.5?

```
## [1] 0.5794473
```

Answer : 0.5788716

Explanation :

After combining the predictions and outcomes with the provided code, we can compute the overall test-set accuracy by creating a classification matrix:

Which tells us that the overall accuracy is $(467 + 1544)/(467 + 1110 + 353 + 1544) = 0.5788716$.

We see a modest improvement over the original logistic regression model. Since predicting stock returns is a notoriously hard problem, this is a good increase in accuracy. By investing in stocks for which we are more confident that they will have positive returns (by selecting the ones with higher predicted probabilities), this cluster-then-predict model can give us an edge over the original logistic regression model.