

Part I Signal Generation

In this part we look at how a trading signal might be developed for a given stock when there is only one factor available. Specifically, we will look at Microsoft stock (ticker MSFT) and we will use the XLK ETF (the Technology Select Sector SPDR ETF) as our factor. The data set covers the years 2018 and 2019. We want to start trading on March 31st, 2018 (three months after the start of the data set) and stop trading on Dec 31st, 2019

Factor Models and Pairs-Trading

One common technique for predicting the behavior of a stock is to use factor models. In this setting, historical data of the stock price or returns are regressed on certain factors; when the factors are observed in the future they can be used to make predictions about the stock's behavior. Mathematically we write this as

$$r_{stock} = \beta^T F$$

where F is a vector of factors and β is a vector of *loadings*; the β s can be determined through regression.

Alternatively, this model can be arrived at by studying *pairs-trading*. Suppose S and F are two stocks (or a stock and a ETF) with similar characteristics, then one expects the returns of the two to be highly correlated and that once this has been corrected for, the discrepancy between the returns can be suitably modeled. In a continuous form this can be stated as

$$\frac{dS_t}{S_t} = \alpha dt + \beta \frac{dF(t)}{F(t)} + dX_t$$

where $d\tilde{X}_t = \alpha dt + dX_t$ is the residual process, and X_t is assumed to be a mean-reverting process. A large value of X_t suggests that the return of S is uncharacteristically large and likely to drop (relative to F) so

one should go short in S and long in F .

For this assignment we have chosen only one factor, the returns of the ETF XLK.

What does the sign of β indicate?

Beta is indicative of the correlation between a stock with a comparative factor (in this case, XLK ETF). Simply put, it's like a bridge that measures the linkage of two stocks. A positive sign of beta indicates that both are moving in same direction. A negative sign indicates movement are opposite. Basically, we can say that the sign of the correlation is determinant of sign of beta.

What do you think β will be (positive or negative) for the Microsoft stock and the XLK factor? Why?

We know that XLK ETF represents the Technology Select Sector SPDR ETF. Because Microsoft is a tech firm and XLK ETF is fund consisting of technology sectors, we could generally say that the beta would be positive, as both would ideally move in the same direction—positively correlated.

What is a market-neutral trading strategy? Describe dollar-neutral trading strategies. What is β -neutral trading?

Market neutral trading strategy aims at minimizing market risks by creating zero correlation for the portfolio to the market. According to Avellaneda and Lee, a market neutral portfolio has a property of having the sum of product of [dollar amounts in stocks] and [beta between factors] be zero. Regarding dollar-neutral trading strategy, it aims at managing risks while generating reasonable returns by means of keeping the correlation low compared to market averages. Regarding beta-neutral trading, it also aims at creating a portfolio with weighted average beta of zero such that market neutrality is accomplished to its best (meaning that market risks are avoided as much as possible without too much exposure).

In terms of market-neutral trading, what do you think is better, a β with large absolute value, or small absolute value? Why?

If we want to achieve market-neutral, we'd better have very small beta in absolute value (hopefully very close to zero). This is because stocks that perfectly correspond the market has beta of exactly 1; stocks that are perfectly uncorrelated has -1 beta. Now we don't want exposure to market in order to minimize market risk. That's why beta with small absolute value is ideal.

Using the Given Financial Data

We have provided you with the daily historical prices for both MSFT (Microsoft stock) and the XLK ETF (Technology Select Sector SPDR ETF) in CSV files, downloaded from Yahoo Finance.

Suppose you want to use three months¹ of historical data in the regression to determine the β (via regression) for each day in the trading window, March 31st, 2018 to Dec 31st, 2019. What date range should you select to generate β 's for the entire time horizon?

2018 Trading Days Calendar																															Monthly		Quarterly					
Month	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	Trading days	Trading hours	Trading days	Trading hours			
Jan	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	21	136.5					
Feb	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W			19	123.5						
Mar	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	21	136.5	61	396.5			
Apr	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M		21	136.5					
May	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	22	143.0					
Jun	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S		21	136.5	64	416.0			
Jul	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	21	133.5					
Aug	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	23	149.5					
Sep	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S		19	123.5	63	406.5			
Oct	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	23	149.5					
Nov	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F		21	133.5					
Dec	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	20	127.0	64	410.0			
Jan 1	New Year's Day										Jul 4	Independence Day										total	252	1629	252	1629.0												
Jan 15	Martin Luther King, Jr. Day										Sep 3	Labor Day										avg	21	135.8	63	407.3												
Feb 19	Washington's Birthday										Nov 22	Thanksgiving Day																										
Mar 30	Good Friday										Nov 23	markets close early at 1:00																										
May 28	Memorial Day										Dec 24	markets close early at 1:00																										
Jul 3	markets close early at 1:00										Dec 25	Christmas																										
note: Stock markets were closed Wednesday, December 5, 2018 in observance of a national day of mourning for President George H. W. Bush.																																						

From the trading days calendar above, to determine beta for each day beginning March 31st, we need to get approximately 60 trading days prior to March 31st, 2018. This could be seen from the first three rows of calendar above.

What is a better fundamental piece of data: the daily return or the daily stock price? What events could trigger a large, yet algorithmically insignificant, discrepancy between the price on two contiguous days? How would you ensure that all of your stock prices are in the same units (reference the same price at a given time)?

I believe the daily return, instead of daily stock price, is a better fundamental piece of data. Sometimes, events like stock market crash would generates tremendous discrepancy between stock price from day to day. If I use daily return, because return is dependent only on percentage rather than price, the data wouldn't be so dramatically significant since price changes are algorithmically insignificant under market neutrality. The way to ensure all prices are in the same units is to divide all stock prices with a constant. This is why we would prefer daily return to daily stock price.

Write a R, Python or MATLAB script to read in the data you downloaded and regress the (daily log-) returns of MSFT on the (daily log-) returns of the XLK ETF to estimate the β for each day from March 31st, 2018 to Dec 31st, 2019, using the opening prices. The regression should use 60 days of historical data, so as not to depend on future information.

```
In [1]: # This function is for usage in the next block
# At each iteration of every 60 trading days, input two stock vectors
with
# corresponding beta and intercept, output X_t for this iteration
get_Xk <- function(vec_msft,vec_XLK,beta_val,intercept) {
  for(i in 1:length(vec_msft)){
    epsilon = vec_msft - beta_val*vec_XLK
    X_k = rep(0,60)
    X_k[1] = epsilon[1]
    for(i in 2:60){
      X_k[i] = X_k[i-1] + epsilon[i] - intercept
    }
  }
  return(X_k)
}
```

A Model For the Residual Process

Consider a portfolio that holds \$1 of the stock and β of the ETF at time t (where β has been determined by regression based historical data). The return of such a portfolio over a small time window is approximately

$$\Delta r_t dt = (r_{stock} - \beta r_{ETF}) dt = d\tilde{X}_t.$$

where $d\tilde{X}_t = \alpha dt + dX_t$ is the residual process. For simplicity, we will discuss the case where $\alpha = 0$, as α is often small in practice. Note that the process X_t corresponds to the cointegrated process computed in the previous section, which we hypothesized as being a mean reverting process. Intuitively, we expect the

behavior of the stock to fluctuate about the factor.

Suppose we model the residual process with the most basic mean-reverting stochastic process, the Ornstein Uhlenbeck process, with dynamics

$$dX_t = \kappa(m - X_t)dt + \sigma dW_t$$

where dW_t is a standard Brownian motion. Under this model, it says that

$$\Delta r_t dt = dX_t = \kappa(m - X_t)dt + \sigma dW_t$$

which has an expected return of $\kappa(m - X_t)dt$. The implications of this for a trading strategy will be discussed below; however, it should be clear that estimation of the parameters of mean reversion is important.

Questions:

1. Why does it make economic sense to assume that X_t is mean reverting? What does that mean about the underlying prices/returns?

The above function gives us the X_t (X_k in my function) required by the question. It makes sense to assume X_t to be mean reverting because sometimes, the price of microsoft stock price might be overpriced or underpriced due to company-related factors. Under EMH, in a highly efficient market, other stock prices of tech sector would also be affected. Hence, with mean-reversion time of the residual process, potential difference due to the incidents would be offset.

2. Write a R, Python or MATLAB function that estimates the parameters of an OU process. Hint: The discrete OU process can be written as an AR-1 model where:

$$X_{n+1} = a + bX_n + \zeta_{n+1}$$

where

$$b = e^{-\kappa\Delta t}, \quad a = m \cdot (1 - e^{-\kappa\Delta t}), \quad \text{var}(\zeta) = \sigma^2 \frac{1 - e^{-2\kappa\Delta t}}{2\kappa}$$

and Δt is the time between observations in years. Here you regress for a , b and ζ and back out the parameters κ , m and σ . In theory the ζ_n are independent Gaussian random variables with mean zero and variance as given above. See also the appendix of Avellaneda and Lee for further details. Fit the OU parameters to each set of 60 days used in the regressions.

```
In [2]: library(forecast)
```

```
Registered S3 method overwritten by 'xts':
  method      from
as.zoo.xts zoo
```

```
Registered S3 method overwritten by 'quantmod':
  method      from
as.zoo.data.frame zoo
```

```
Registered S3 methods overwritten by 'forecast':
  method      from
fitted.fracdiff fracdiff
residuals.fracdiff fracdiff
```

```

In [3]: data_msft <- read.csv(file = 'MSFT.csv')
data_XLK <- read.csv(file = 'XLK.csv')

Beta=c() #capture series of beta for part 2.1.q2 plotting

## Variables in this section is for usage in Part 3 and Part 4
b_vec=c() # b in AR-1 Model
s_val=c() # s scales the displacement from mean by average size of deviation
s_list = c() #capture series of s_val
mean_rev=c() # half-life of process

#iterate through every 60 trading days
for(i in 1:(length(data_msft$Open)-60)){

  #calculate daily return using open price for two stocks
  vec_price_msft = data_msft[c(i:(i+60)), "Open"]
  vec_logreturn_msft <- log(vec_price_msft[-1]/vec_price_msft[-length(vec_price_msft)])
  vec_price_XLK <- data_XLK[c(i:(i+60)), "Open"]
  vec_logreturn_XLK <- log(vec_price_XLK[-1]/vec_price_XLK[-length(vec_price_XLK)])

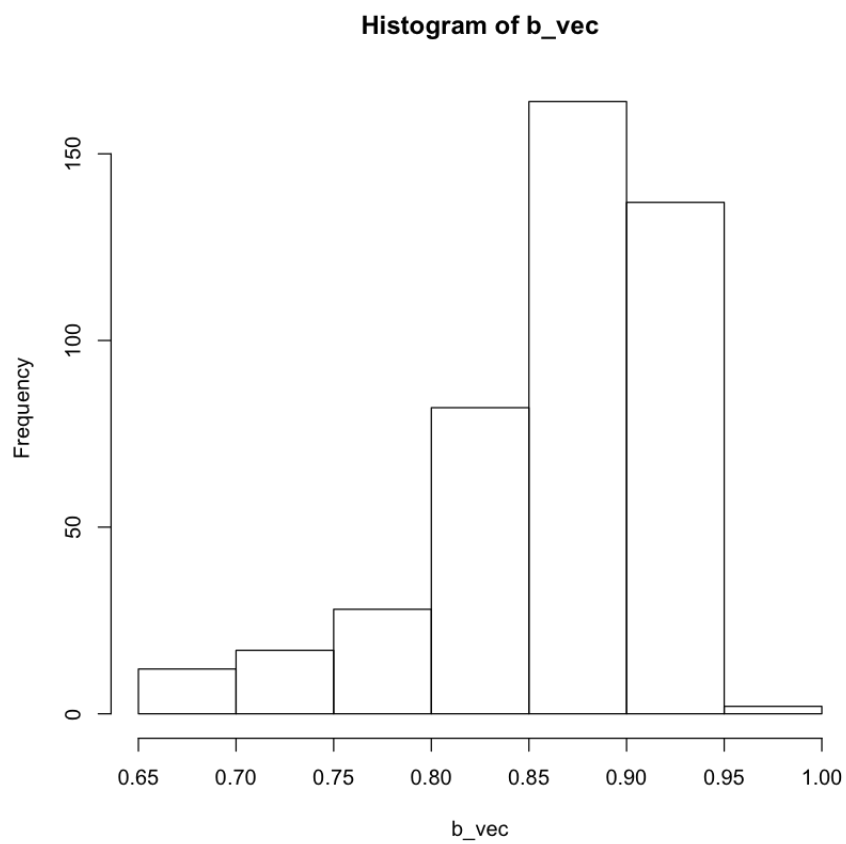
  #regression
  model <- lm(vec_logreturn_msft ~ vec_logreturn_XLK)
  beta = model$coefficients[2]
  intercept = model$coefficients[1]

  #append information
  X_k=c()
  X_k <- c(X_k, get_Xk(vec_logreturn_msft, vec_logreturn_XLK, beta, intercept))
  Beta <- c(Beta, beta)

  #perform AR(1)
  model = arima(X_k, order=c(1,0,0))
  a = model$coef[2]
  b = model$coef[1]
  b_vec <- c(b_vec, b)
  k = -log(b)*252
  mean_rev <- c(mean_rev, 1/k)
  m = a/(1-b)
  sigma_2 = model$sigma2[1]
  sigma = sqrt(sigma_2*2*k/(1-b^2))
  sigma_eq = sigma/sqrt(2*k)
  s_val= -m/sigma_eq
  s_list <- c(s_list, s_val)
}

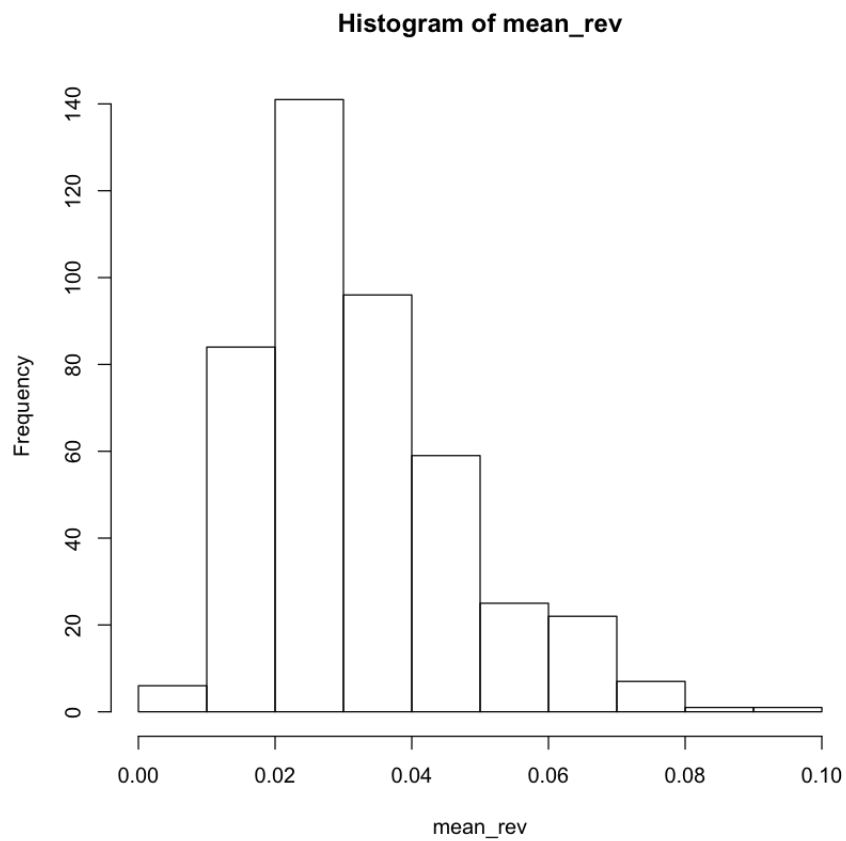
hist(b_vec)

```



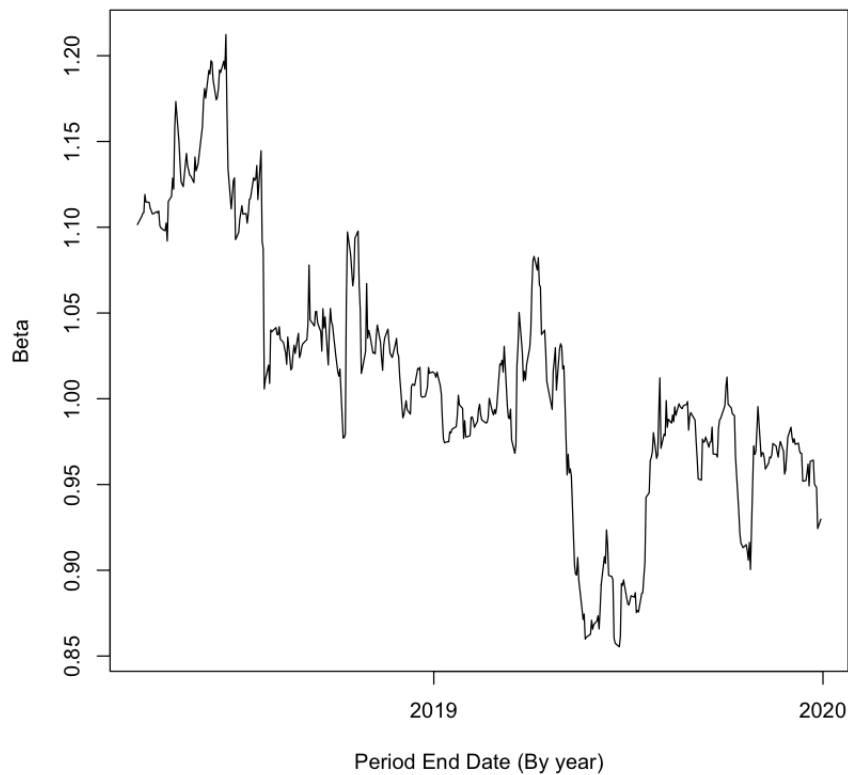
3. Produce a histogram of the “mean-reversion times” $1/\kappa$. Do you think a 60-day data window was reasonable?

```
In [4]: hist(mean_rev)
```



4. Produce a histogram of b and compare it with the histogram of $1/\kappa$. Which is less skewed? Why?


```
In [5]: D = data_msft[c(61:(length(data_msft$Open))), "Date"]  
plot(Beta~as.Date(D, "%Y-%m-%d"), type="l", xlab="Period End Date (By year)", ylab="Beta")
```



From above, the mean_rev one is more skewed than beta one. This implies that the hypothesis in our mean reversion fails more than the hypothesis in beta (not close to 1).

Generating a Trading Signal

Suppose we invest (long) \$1 in the stock and short β of the factor at time zero, then the profit of this investment is given by

$$e^{\int_0^t r_s(t)dt} - \beta e^{\int_0^t r_f(t)dt} - (1 - \beta)$$

which for small times t is approximately

$$(1 + r_s(0)t) - \beta(1 + r_f(0)dt) - (1 - \beta) = r_s(0)dt - \beta r_f(0)dt = d\tilde{X}_t(0)$$

which has expectation (assuming $\alpha = 0$)

$$\kappa(m - X_t)dt.$$

What this says is that if you were to go long a dollar in the stock and short β dollars of the factor, your (expected) instantaneous rate of return would be $\kappa(m - X_t)dt$. So, clearly you want to do this when $m - X_t$ is positive. Furthermore, as the accumulated rate of return from the start of entering the position is the integral

$$\int_0^t \kappa(m - X_t)dt + \int_0^t \sigma dW_t$$

you want to remain in this position as long as the expected additional rate of return is positive, that is, while $m - X_t > 0$. The closer to zero $m - X_t$ becomes the more dominant the stochastic fluctuations are, and the more likely you are to start losing money (by switching to a region where the difference in the two positions becomes negative), therefore you want to enter this position at a time when $m - X_t$ is sufficiently positive that you will remain in the position long enough to account for transaction costs, but exit as some point before $m - X_t$ becomes negative.

The parameter $1/\kappa = \tau$ is the half-life of the process, the time it would take the deterministic system to move closer to m by a factor of e . In this sense, κ determines how quickly the process will return to its mean, and thus how quickly $m - X_t$ goes to zero. In the stochastic setting, it gives a measure of how frequent the excursions from the mean are, and how long they will last; this relationship between κ and the average fluctuation size is evident by the role it plays in the equilibrium standard deviation of the process:

$$\sigma_{eq} = \frac{\sigma}{\sqrt{2\kappa}}.$$

σ_{eq} can be used to non-dimensionalize the process, by measuring the size of $m - X_t$, in terms of standard deviations. Taking for example:

$$s = \frac{X_t - m}{\sigma_{eq}}$$

scales the displacement from the mean by the average size of such deviations, allowing one to detect large excursions easily across different processes. Notice that the regression requires $X_{60} = 0$, so we get that

$$s = \frac{-m}{\sigma_{eq}}$$

A large and positive value of s suggests that the return of the stock will fall relative to the return of the factor so shorting the stock and buying the factor is a strong position; when s becomes smaller (but still positive), it is a good time to exit such a position. Conversely, when s is large and negative, taking a long position in the stock and a short position in the factor is correct. See Figure 1. When computing the s value, updated estimates of β , κ , σ and m are used, based on the last 60 days.

A trading strategy can thus be defined by four numbers: $s_{long,enter}$, $s_{long,exit}$, $s_{short,enter}$ and $s_{short,exit}$ where $s_{long,enter} < s_{long,exit} < 0$ and $s_{short,enter} > s_{short,exit} > 0$.

1. In your own words, explain the trading strategy described above and why it works. Give explicit rules for when to enter/exit long/short positions in terms of the numbers $s_{long,enter}$, $s_{long,exit}$, $s_{short,enter}$, $s_{short,exit}$ and the trading signal $s(t)$.

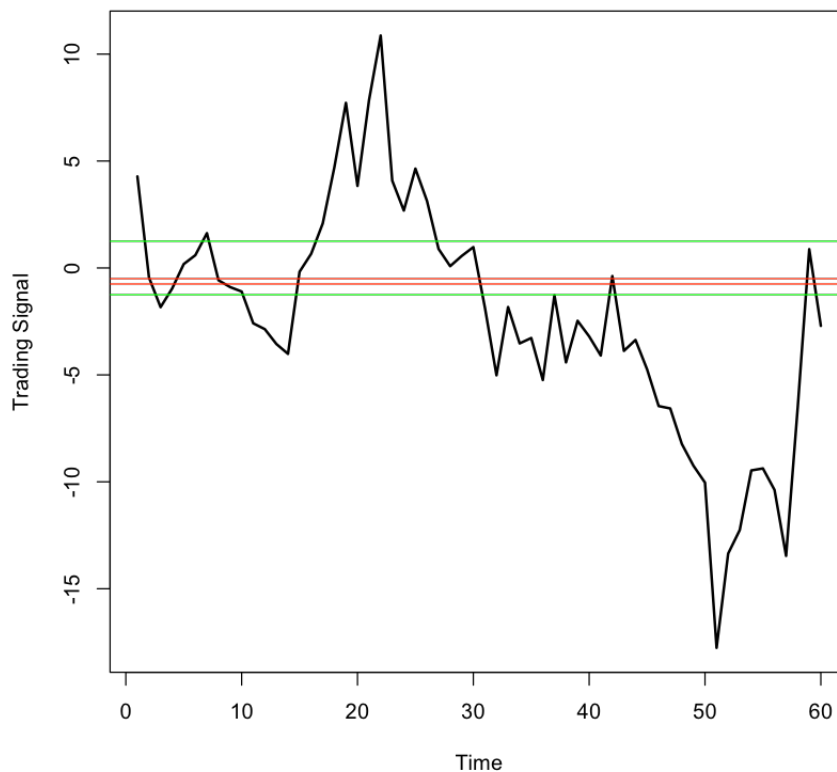
Suppose we want to buy a dollar in microsoft and sell beta amount of ETF, the instantaneous rate is supposedly $k(m - X_t)d_t$. As a rule of thumb, we would want to remain in a position so long as the $m - X_t$ is positive, i.e. expected additional rate of return is positive. Specifically speaking, when $m - X_t$ is positive enough, we would maintain in the position to account for transaction costs; upon becoming negative, we would exit as some point. A positive signal indicates that shorting the MSFT and buying the ETF is a strong decision; a negative one indicates the other way around. For s , we're supposed to have $[short\ enter > short\ exit > 0]$ and $[long\ enter < long\ exit < 0]$.

2. Using the data, generate the s value for each 60-day window. For the values:

$$s_{long,enter} = -1.25 \quad s_{long,exit} = -0.50 \quad s_{short,enter} = 1.25 \quad \text{and} \quad s_{short,exit} = 0.75$$

plot them with first 60 days of s values and indicate where you enter and exit positions; the plot should, stylistically, look like Figure 1. How many times do you enter a long position? a short position?

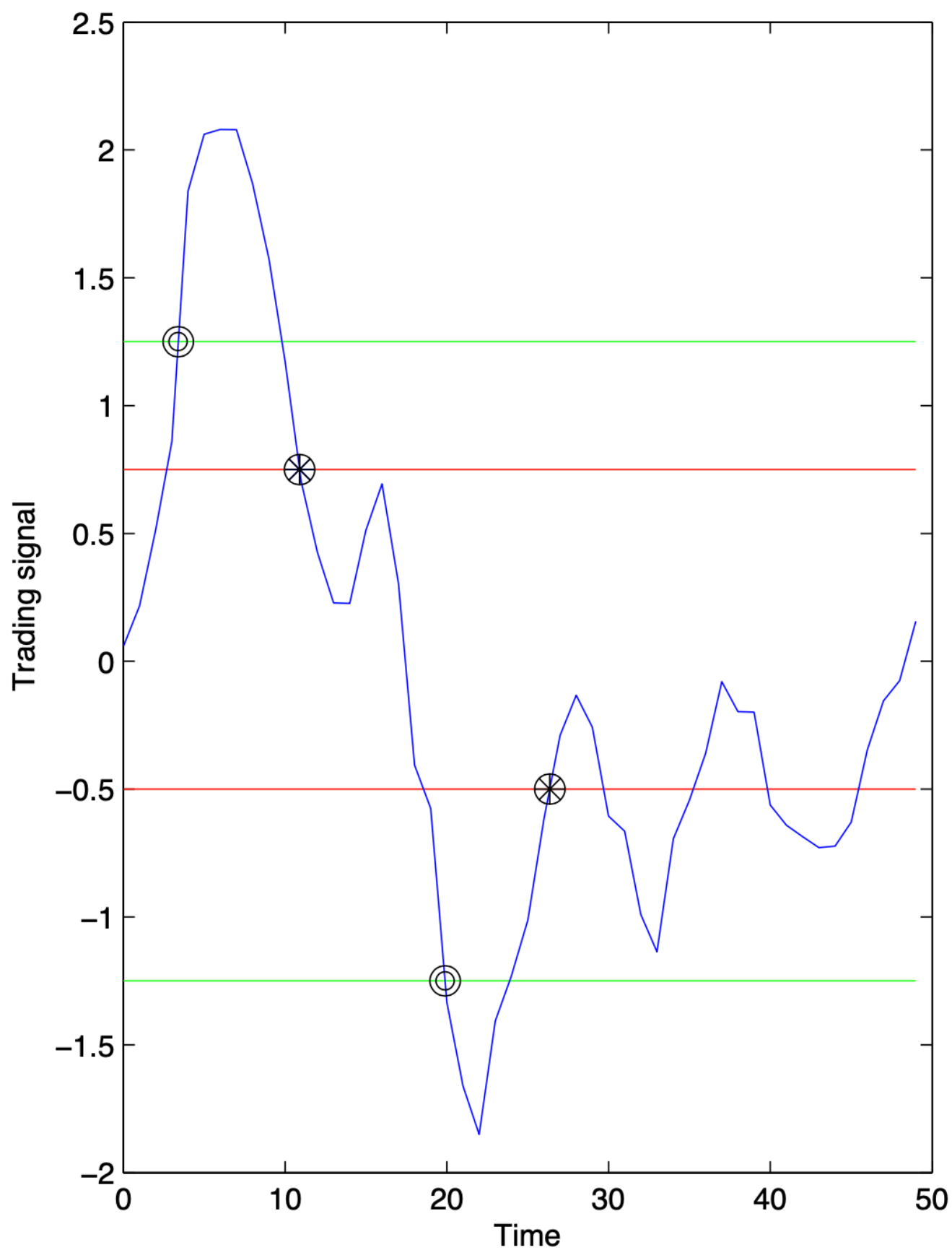
```
In [6]: plot(s_list[1:60],type='l', col="black", lwd = 2, xlab="Time",ylab="Trading Signal") +  
abline(h=1.25, col="green") + abline(h=-1.25, col="green") +  
abline(h=-0.5, col="red") + abline(h=-0.75, col="red")
```



I entered long position for 5 times and entered short position for 2 times.

3. Execute the trading strategy (you will now need price data). How does it perform? What is the Sharpe ratio for the last 5 months (using a 1 month averaging period) computed against 2% constant growth? You may assume that any money left over or lost at the end of the day is put into / taken from a bank account with an interest rate of 2%. Qualitatively, how does this strategy perform?

Evaluating Sharpe ratio for the profit by the strategy, assuming a 2% constant growth, we'd have at least an at least 3-times-better ratio compared to trading without pair-trading-signal.



Part II Empirical Factors

Empirical Correlation Matrices

Suppose we want to generate empirical factors on day J ; we will use a $n_{d,ep} = 252$ day (one year) window to compute our statistics. We start by standardizing the the daily returns: if r_{ij} is the return of the i th stock on the j th day, then the standardized return is

$$Y_{i,J-j} = \frac{r_{i,J-j} - \bar{r}_i}{s_i}, j = 0, \dots, n_{d,ep} - 1$$

where

$$\bar{r}_i = \frac{1}{n_{d,ep}} \sum_{j=0}^{n_{d,ep}-1} r_{i,J-j} \quad \text{and} \quad s_i^2 = \frac{1}{n_{d,ep} - 1} \sum_{j=0}^{n_{d,ep}-1} (r_{i,J-j} - \bar{r}_i)^2.$$

The rows of the matrix $Y = \{Y_{i,J-j}\}_{i=1, \dots, n_{stocks}, j=0, \dots, n_{d,ep}-1}$ are now mean zero and variance one. The correlation matrix can then be computed as

$$\rho_J = \frac{1}{n_{d,ep} - 1} Y Y^T.$$

It is the principal components of this matrix that will form the basis of our empirical factors. If $\rho_J = V \Lambda V^T$ is the eigen-decomposition of the correlation matrix,

$$V = \begin{bmatrix} | & | & & | \\ v^{\{1\}} & v^{\{2\}} & \dots & v^{\{n_{stocks}\}} \\ | & | & & | \end{bmatrix}$$

where $v^{\{k\}}$ is the k th eigenvector, corresponding to the k th eigenvalue λ_k and

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{n_{stocks}} \geq 0.$$

Observe that

$$\sum_{i=1}^{n_{stocks}} \lambda_i = \text{trace}(\rho_J) = n_{stocks}.$$

The *eigen-portfolios*, or empirical risk factors, are vectors $p^{\{k\}} = \{p_i^{\{k\}}\}_{i=1, \dots, n_{stocks}}$, where

$$p_i^{\{k\}} = \frac{v_i^{\{k\}}}{s_i}$$

where $v_i^{\{k\}}$ is the i th weighting of stock i in the k th principal component and s_i is its empirical standard deviation. The interpretation of $p_i^{\{k\}}$ is that in the k th empirical risk factor, we invest $p_i^{\{k\}}$ dollars in stock i (if it is negative, we short the stock). The return of the k th eigen-portfolio on day j is

$$F_{kj} = \sum_{i=1}^{n_{stocks}} \hat{p}_i^{\{k\}} r_{ij}$$

where $\hat{p}_i^{\{k\}} = p_i^{\{k\}} / e^T p^{\{k\}}$ (in other words, the weights have been normalized to sum to one). One crucial fact about the returns of the eigen-portfolio returns is that they are empirically orthogonal. Finally, the “percentage of explained variance” of a given set \mathcal{K} of eigen-portfolios is defined by

$$\frac{\sum_{k \in \mathcal{K}} \lambda_k}{n_{stocks}} \times 100\%.$$

Note that this term might be a bit misleading as it does not correspond to variance explained (the PCA factors of the covariance matrix would do this) but rather correlation explained.

Questions:

Write a R, Python or MATLAB script to compute the eigen-portfolios on a given day.

```
In [7]: library(R.matlab)
```

```
R.matlab v3.6.2 (2018-09-26) successfully loaded. See ?R.matlab for help.
```

```
Attaching package: 'R.matlab'
```

```
The following objects are masked from 'package:base':
```

```
getOption, isOpen
```



```
In [ ]: final_proj_data <- readMat("finalproject_data.mat")

#process data
num_stocks <- final_proj_data$nStocks[1,1]

num_days <- final_proj_data$nDays[1,1]

return_dat <- final_proj_data$returns
dat <- as.data.frame(t(as.matrix(return_dat)))

stock_name_list <- c()
for(i in 1:422){
  each_stock_name= toString(final_proj_data$names[[i]])
  #print(each_stock_name)
  stock_name_list <- c(stock_name_list,each_stock_name)
}

stock_ticker_list <- c()
for(i in 1:422){
  each_ticker_name= toString(final_proj_data$tickers[[i]])
  names(dat)[i] <- each_ticker_name
  stock_ticker_list <- c(stock_ticker_list,each_ticker_name)
}
```

```

In [ ]: Fld_vec <- c()

# Beginning from day 252, up to day 502
for(d in (252:num_days)){
  interim = return_dat[, (d-251):d]

  rbar_vec = rowSums(interim)/252
  s_square = rowSums((interim - rbar_vec)^2)/251

  Y = interim / sqrt(s_square)

  #get matrix rho
  rho = Y%*%t(Y)/251
  #retrieve eigen values and vectors
  V = eigen(rho)$vectors
  lambda_val = eigen(rho)$values
  lambda_val_precentage = lambda_val/num_stocks
  lambda = diag(lambda_val)

  P = V
  P_hat = P

  PR = 0

  for(k in (1:num_stocks)){
    P[,k] = V[,k]/sqrt(s_square)
    P_hat[,k] = P[,k]/sum(P[,k])

    # Calculating cumulative return
    PR = PR + (P_hat[k,1]*return_dat[k,d])

  }
  Fld_vec <- c(Fld_vec,PR)
}

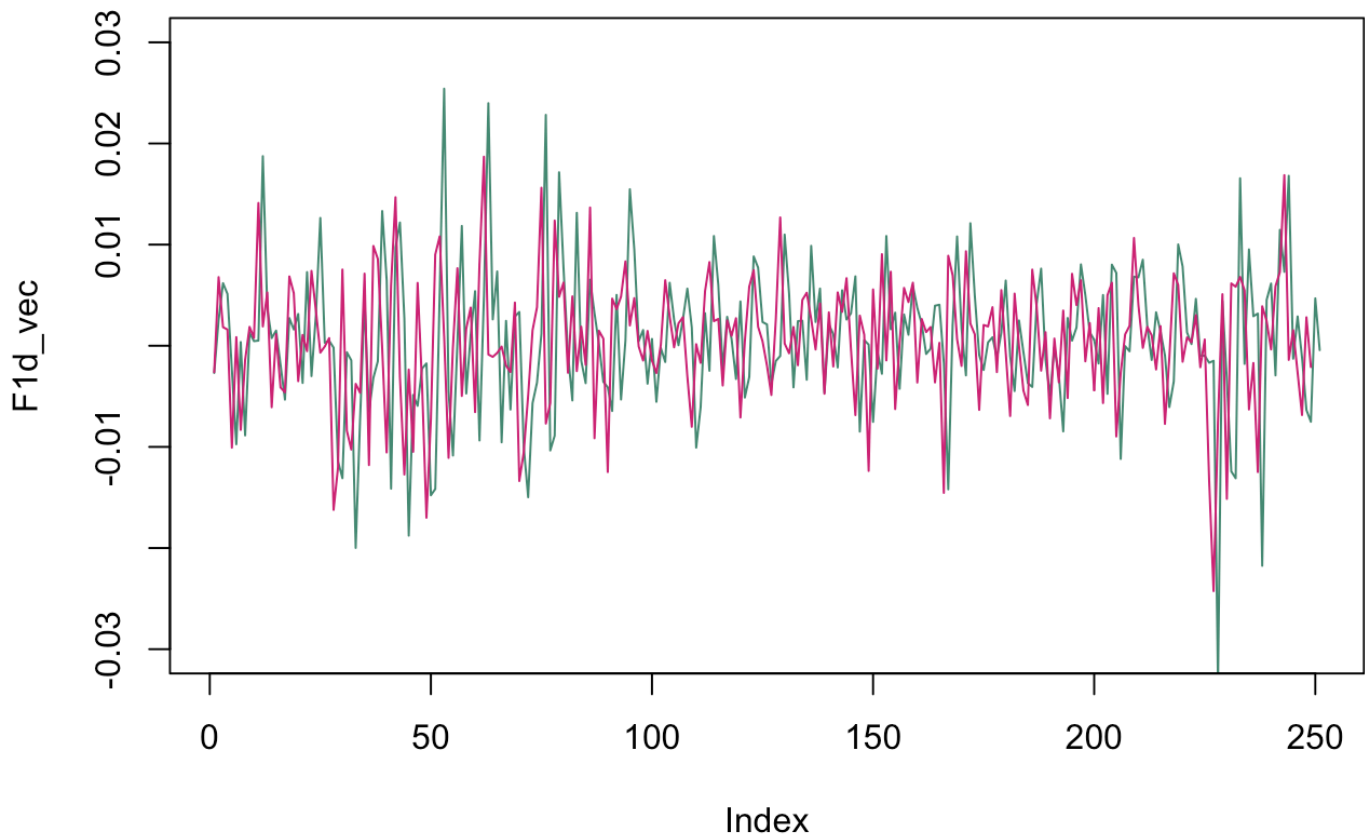
```

```

In [ ]: data_SPY_yahoo <- read.csv(file = 'SPY_from_YahooFinance.csv')
SPY_price_yahoo <- data_SPY_yahoo$Open
SPY_return_yahoo <- diff(SPY_price_yahoo)/SPY_price_yahoo[-length(SPY_
price_yahoo)]

plot(Fld_vec,type='l',ylim=c(-0.03,0.03),col='aquamarine4')
lines(SPY_return_yahoo,type='l',col='deeppink3')

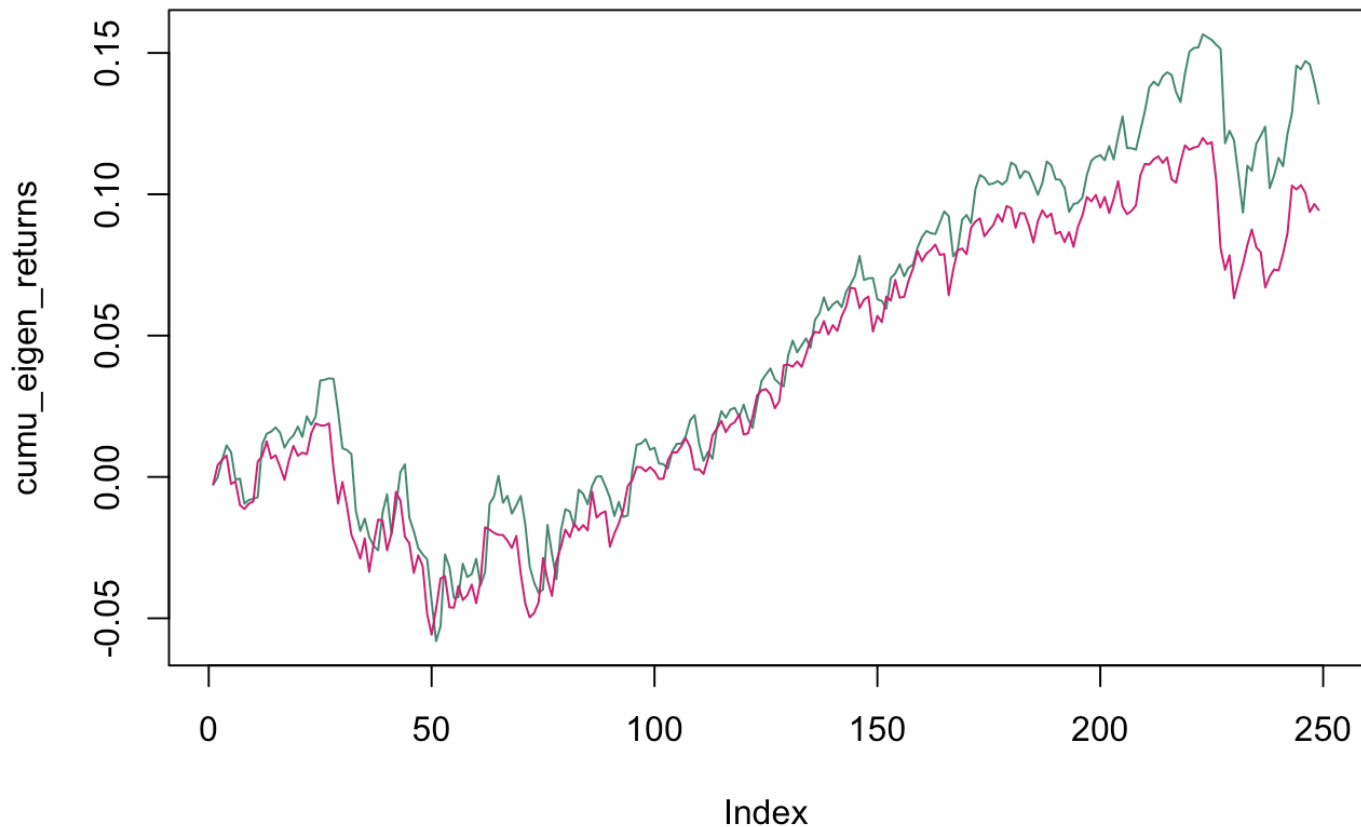
```



Plot the cumulative return of the principal eigen-portfolio over the course of time. How does it compare to the market? (You can use the SPY (an index based on the S&P 500) as a representation of the market; you need only look at a plot of it on, say, Google's finance page.)

```
In [ ]: cumu_eigen_returns = c()
        cumu_SPY_returns = c()
        for(i in (1:length(SPY_return_yahoo))){
            cumu_eigen_returns = c(cumu_eigen_returns,sum(F1d_vec[1:i]))
            cumu_SPY_returns = c(cumu_SPY_returns,sum(SPY_return_yahoo[1:i]))
        }

        plot(cumu_eigen_returns, type = 'l', col='aquamarine4')
        lines(cumu_SPY_returns, type='l', col='deeppink3')
```



From the plot above, the cumulative return aligns with the market well.

Carefully read Avellaneda and Lee's discussion about eigen-portfolios. Explain, in your own words, the financial implications of the weighting given to each stock in different eigen-portfolios. For example, if stock i has a large positive value in the second eigen-portfolio but a large negative value in the third, and stock j has a large positive value in both eigen-portfolios, what can be said about the following two positions: going long in both stock i and stock j and going long in stock j and shorting stock i ?

According to Avellaneda and Lee, the weights in the dominant eigenportfolio are inversely correlated to the stock volatility. In this paper, we also know that if we relabel $evec$ in descending order, then the nearest stocks tend to be in the same industry sector. As number increases, this becomes not as true. Consider stock i has large positive 2nd and large negative 3rd eval; j has large positive in both 2nd and 3rd, we would have the 2nd eigenportfolio associating to longing both stocks yet 3rd eignportfolio for shorting i and longing j .

The Marcenko-Pastur Distribution for the Distribution of Eigen- values

Suppose that X is a matrix with m rows and n columns and that each entry X_{ij} is an iid mean-zero-variance-one Gaussian random variables. The Marcenko-Pastur distribution of the eigenvalues of the covariance matrix

$$\frac{1}{n-1}XX^T = V\Lambda V^T$$

have a limiting density

$$\rho(\lambda) = \frac{Q}{2\pi} \frac{\sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}}{\lambda}, \lambda \in [\lambda_-, \lambda_+]; \quad \lambda_{\pm} = \left(1 \pm \sqrt{Q^{-1}}\right)^2.$$

where $Q = n/m$ and the limit is as $m, n \rightarrow \infty$ so that Q is a constant. As $n \rightarrow \infty$ one can note that the covariance matrix in this vanilla case also converges to the correlation matrix, and thus the same result hold for the correlation matrix.

In other words

$$\Pr(\{\lambda_k\} \cap (\lambda + d\lambda) \neq \emptyset) d\lambda = \rho(\lambda) d\lambda \approx \frac{\# \text{ of eigenvalues in } (\lambda, \lambda + d\lambda)}{m}.$$

Let $n = 5000$ and $m = 1000$ and suppose the X_{ij} are mean zero, variance one Gaussian random variables. Using R, Python or MATLAB, simulate a data matrix X , compute the eigen-values of its covariance matrix, and plot the distribution of the eigenvalues. On top of this, plot the limiting Marcenko-Pastur density. Comment on the results.

```
In [9]: #Number of rows
m = 1000
#Number of columns
n = 5000

Q = n/m

#Generate matrix
X <- matrix(rnorm(m*n,mean=0,sd=1), m, n)
A <- 1/(n-1)*X%*%t(X)

#Perform Eigenvalue decomposition
E <- eigen(A)$vectors
d <- eigen(A)$values
D <- diag(d)

interval_6a <- c(seq(0,2.5,by=0.05))
histo_6a <- hist(d, plot=FALSE, breaks=interval_6a)
density_6a <- c(histo_6a$counts)
experiment_6a = density/length(d)
max(d)
```

2.09849881358075

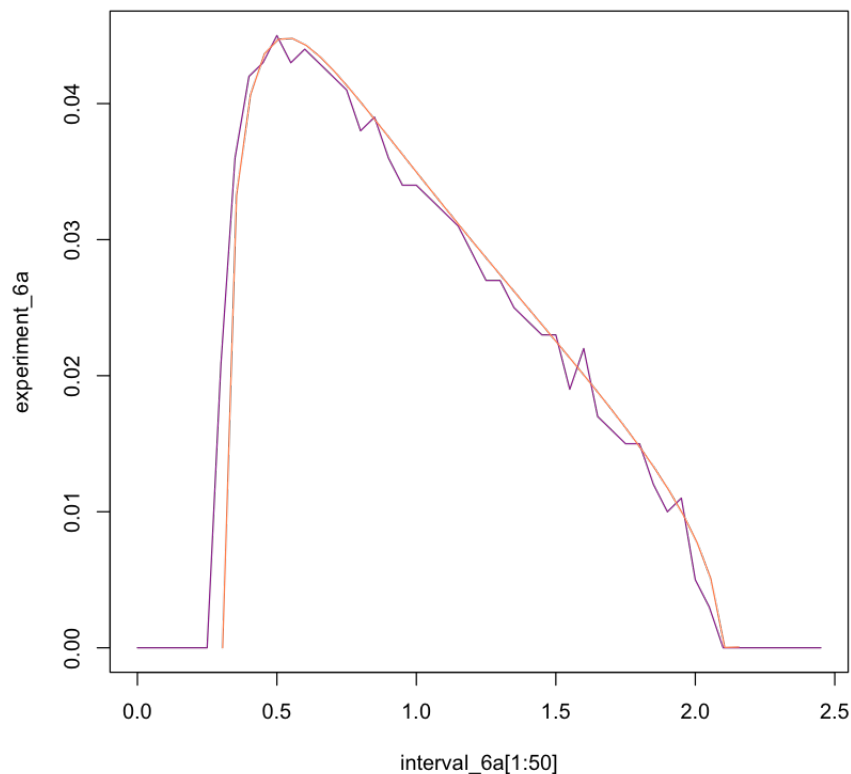
```
In [10]: min(d)
```

0.305813459481512

```
In [11]: lambda_plus = (1+sqrt(1/Q))^2
lambda_minus = (1-sqrt(1/Q))^2

lambda <- c(seq(lambda_minus,2.2,by=0.05))
scalar = Q/2/pi
i = lambda_plus - lambda
j = lambda - lambda_minus
k = sqrt(as.complex(i*j))
rho = scalar * k / lambda
rho = rho/sum(rho)
#experimental
plot(interval_6a[1:50], experiment_6a, type='l', col='darkmagenta')
#Marcenko-Pastur
lines(lambda, rho, type='l', col='coral')
```

Warning message in xy.coords(x, y):
"imaginary parts discarded in coercion"



Both lines align with each other well with the defined bin size.

Often times stock returns are modeled as iid Gaussian random variables. If this were the case, the distribution of the eigenvalues of the correlation matrix of the stock returns should look like the Marcenko-Pastur Distribution. Standardize the entire history of stock returns so that the mean for each stock is zero and the variance one (over the entire 502 day period). Compute the eigenvalues of the data's correlation matrix. Plot the density of the eigenvalues as well as the theoretical distribution given by the Marcenko-Pastur law. How do they compare?

```
In [12]: library(R.matlab)
#final_proj_data <- readMat("finalproject_data.mat")
return_dat <- final_proj_data$returns
data_mat <- data.matrix(return_dat)

#Perform standardization
for(row in 1:422){
  data_mat[row,] = data_mat[row,] - mean(data_mat[row,])
  data_mat[row,] = data_mat[row,]/sd(data_mat[row,])
}

#Correlation matrix
rho_6b = (1/501)*data_mat%*%t(data_mat)

#Eigenvalues
lambda_6b <- eigen(rho_6b)$values
D <- diag(lambda_6b)

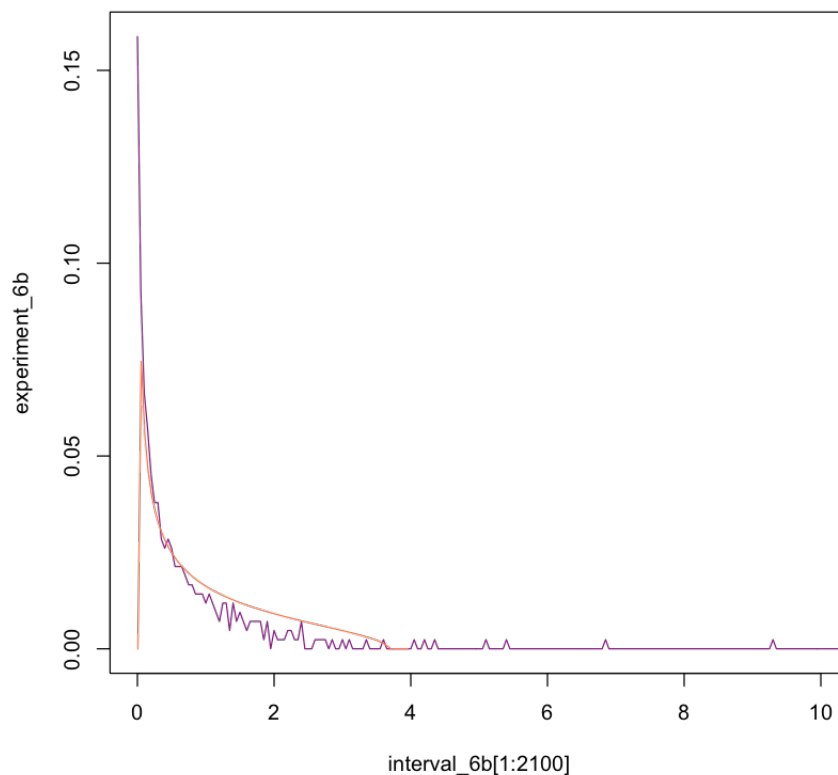
#Create density
interval_6b <- c(seq(0,105,by=0.05))
histo_6b <- hist(lambda_6b, plot=FALSE, breaks=interval_6b)
density_6b <- c(histo_6b$counts)
experiment_6b = density_6b/length(lambda_6b)
```

```
In [13]: #Create Marcenko-Pastur density
m=422
n=502
Q = n/m;
lambda_plus = (1+sqrt(1/Q))^2
lambda_minus = (1-sqrt(1/Q))^2
lambda_MP = c(seq(lambda_minus,4,by=0.05))

scalar = Q/2/pi
a = lambda_plus - lambda_MP
b = lambda_MP - lambda_minus
c = sqrt(as.complex(a*b))

densMP = scalar * c / lambda_MP
densMP = densMP/sum(densMP)
plot(interval_6b[1:2100],experiment_6b, type='l', col='darkmagenta', x
lim=c(0,10))
lines(lambda_MP, densMP, type='l', col='coral')
```

Warning message in xy.coords(x, y):
 “imaginary parts discarded in coercion”



Although shapes are generally the same, for the values especially in the first half of interval, the density doesn't agree too much as it flattens throughout the interval.

Now take the matrix Y of standardized returns used in the previous section and apply a (different) random permutation to each time series of returns so as to remove any time correlation in the returns. Compute the correlation matrix again, and plot the distribution of its eigenvalues. Compare this with the Marcenko-Pastur distribution; how has it changed? How might you use this to determine if your returns are iid?

```
In [14]: # Create random permutations
for(row in 1:422){
  data_mat[row,] <- sample(data_mat[row,])
}

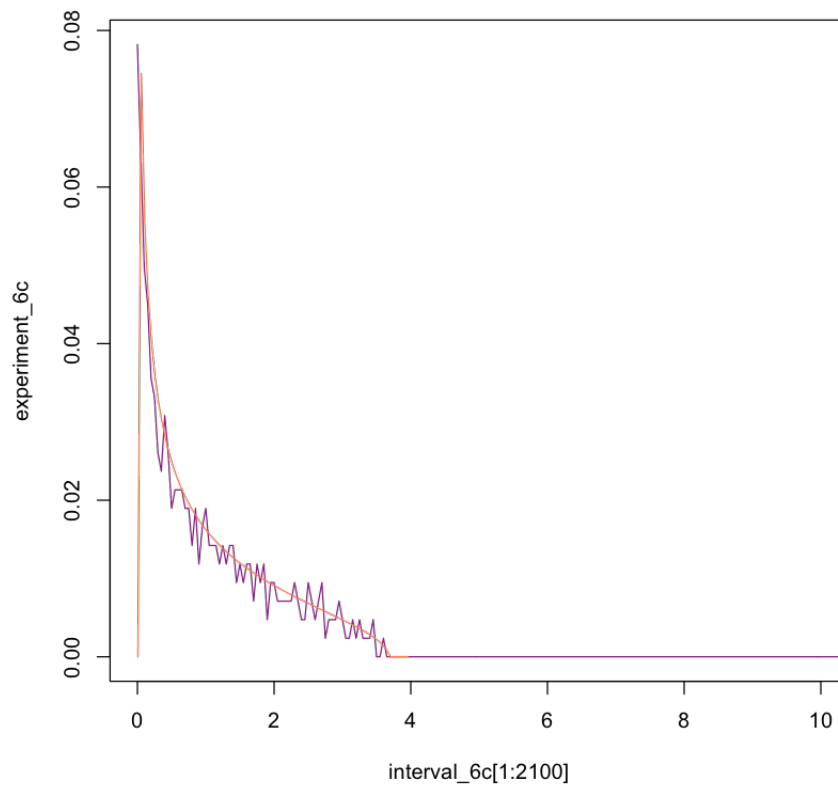
#Correlation matrix
rho_6c = (1/501)*data_mat%*%t(data_mat)

#Compute eigenvalues of the matrix
lambda <- eigen(rho_6c)$values

#Create density
interval_6c <- c(seq(0,105,by=0.05))
histo_6c <- hist(lambda, plot=FALSE, breaks=interval_6c)
density_6c <- c(histo_6c$counts)
experiment_6c = density_6c/length(lambda)

plot(interval_6c[1:2100],experiment_6c, type='l', col='darkmagenta', x
lim=c(0,10))
lines(lambda_MP, densMP, type='l', col='coral')
```

```
Warning message in xy.coords(x, y):  
"imaginary parts discarded in coercion"
```



Compared to 6.2, this seems to perform better by removing time correlation in the returns. From this, we see that by randomizing each series in the data independently (using permutation), the observed result aligns with theoretical Marcenko-Pastur one. This indicates that observed returns are actually not IID – correlations exist in the data and play an important role.