

Overview

Finance and algorithmic trading heavily use linear algebra. In this project, we will investigate how you can directly create a portfolio (set of allocations for companies in the market) by using the singular value decomposition (SVD) that you learned in class. In essence, we will use the past daily prices for over 1000 companies over a "training" period to generate right singular vectors and normalize these singular vectors to get allocations for each of the 1000 companies over the "test" period.

In particular, we outline below how to create these "eigen"-portfolios:

1. Import necessary packages and data.
1. Pre-process, clean up, and plot the data using pandas and matplotlib.
1. Transform the data to use company returns and normalized returns rather than close prices.
1. Split the data into a training set and a test set.
1. Compute the SVD of the training data (you will complete this portion).
1. Plot the SVD to get a sense of the market.
1. Compile the SVD into the eigen-portfolios.
1. Compute returns and cumulative returns.
1. Compute performances of portfolios with the Sharpe ratio.
1. Plot performances of best eigen-portfolios.

The only aspects of these steps that you need to actually do here will be to compute the SVD and compile the SVD into the eigen-portfolios.

```
In [1]: # import packages
import pandas as pd
# clean data
import numpy as np
import os
import datetime
import matplotlib.pyplot as plt
%matplotlib inline

df_path = 'close_prices.csv'

c:\Users\math\anaconda3\lib\site-packages\pandas\core\computation\expressions.py:20: UserWarning: Pa
nadas requires version '2.7.3' or newer of 'numexpr' (version '2.7.1' currently installed).
  from pandas.core.computation.check import NUMEXPR_INSTALLED
```

Uploading data into Google Colaboratory

To actually run this jupyter notebook, you can run jupyter notebook on your own computer if you have it set up; however, you may also use Google Colaboratory to create and run this notebook. I've implemented the first method of how to import the data into google colab. The details are outlined [here](#). If you are using Google colab, uncomment the lines of code after the first in the next cell.

```
In [2]: # only uncomment the next lines if using google colab (takes some time)
# import io
# from google.colab import files
# df_path = files.upload()
# df_path = io.StringIO(uploaded['close_prices.csv']).decode('utf-8')
```

Data Wrangling

We clean the data here to make sure that any of the dates we use have at least 500 data points. In particular, we get rid of dates that have too many NaN values. We also transform the dates into pandas datetime indices to be able to easily manage data. We finally take a look at some portion of the data frame.

```
In [3]: # import data
close_prices = pd.read_csv(df_path)
# clean data
close_prices['date'] = close_prices['date'].apply(lambda x: x.split()[0])
close_prices = close_prices.set_index(['date'])
close_prices = close_prices[close_prices.index.duplicated(keep='first')]
close_prices = close_prices[close_prices.isnull().sum(axis=1) < 500]
dts = pd.to_datetime(close_prices.index)
close_prices.index = dts
close_prices.name = 'prices'
close_prices.head(40)
```

```
Out[3]:
```

	AAC	AAL	AAON	AAP	AAPL	AB	ABB	ABV	ABC	ABCB	...	Y	YELL	YNDX	
date															
1999-11-01	17.215847	NaN	0.963988	0.595872	5.832890	NaN	NaN	2.646925	5.628370	...	150.909670	108830.902945	NaN	11.01	
1999-11-02	18.253337	NaN	0.965558	NaN	6.16062	6.378645	NaN	NaN	2.621790	5.602393	...	151.689773	108438.483824	NaN	11.01
1999-11-03	18.026386	NaN	1.009507	NaN	6.625658	6.289377	NaN	NaN	2.536717	5.76416	...	151.897299	110400.579430	NaN	11.01
1999-11-04	18.447866	NaN	0.995558	NaN	6.419343	6.555154	NaN	NaN	2.535151	5.76416	...	150.909670	110008.160309	NaN	10.91
1999-11-05	19.452934	NaN	1.005102	NaN	6.779377	6.620076	NaN	NaN	2.440043	5.76416	...	149.724516	109223.322067	NaN	10.81
1999-11-08	18.836925	NaN	0.991153	NaN	7.398612	6.441539	NaN	NaN	2.453578	5.736608	...	149.724516	107195.823274	NaN	10.81
1999-11-09	18.447866	NaN	0.963988	NaN	6.887993	6.429366	NaN	NaN	2.513915	5.628370	...	149.329464	107588.242395	NaN	10.91
1999-11-10	17.831857	NaN	0.945633	NaN	7.01965	6.327925	NaN	NaN	2.610189	5.76416	...	147.354207	110400.579430	NaN	11.11
1999-11-11	17.831857	NaN	0.940494	NaN	7.08183	6.175762	NaN	NaN	2.501915	5.602393	...	146.666817	111185.417673	NaN	11.01
1999-11-12	20.684954	NaN	0.936089	NaN	6.959670	6.530808	NaN	NaN	2.536717	5.602393	...	145.378949	112362.675036	NaN	10.91
1999-11-15	24.316168	NaN	0.954443	NaN	6.866612	6.417193	NaN	NaN	2.561852	5.602393	...	146.189052	111970.255915	NaN	10.91
1999-11-16	23.505629	NaN	0.963988	NaN	7.00046	6.682970	NaN	NaN	2.501915	5.628370	...	147.749258	113212.916465	NaN	10.81
1999-11-17	23.311100	NaN	0.945633	NaN	6.892380	6.569356	NaN	NaN	2.476779	5.628370	...	148.491955	113409.126026	NaN	10.81
1999-11-18	22.095090	NaN	0.954443	NaN	6.887993	6.542961	NaN	NaN	2.670127	5.844846	...	152.395064	113212.916465	NaN	11.01
1999-11-19	24.737648	NaN	0.954443	NaN	6.079642	6.518635	NaN	NaN	2.846925	5.901130	...	150.514619	111577.836794	NaN	11.21
1999-11-22	23.116570	NaN	0.954443	NaN	6.959670	6.352271	NaN	NaN	2.583610	5.684654	...	149.717934	113605.335587	NaN	11.31
1999-11-23	21.495493	NaN	0.936089	NaN	7.12482	6.301550	NaN	NaN	2.380106	5.844846	...	148.179345	113605.335587	NaN	11.31
1999-11-24	20.879483	NaN	0.945633	NaN	7.260915	6.250829	NaN	NaN	2.380106	5.736608	...	148.586768	112362.675036	NaN	11.41
1999-11-26	20.879483	NaN	0.917734	NaN	7.29755	6.277204	NaN	NaN	2.380106	5.736608	...	147.796665	109615.741188	NaN	11.31
1999-11-29	19.258405	NaN	0.936089	NaN	7.25917	6.086494	NaN	NaN	2.368505	5.684654	...	148.539361	109223.322067	NaN	11.31
1999-11-30	19.258405	NaN	0.927278	NaN	7.51327	6.112668	NaN	NaN	2.393640	5.684654	...	151.689773	107195.823274	NaN	11.31
1999-12-01	17.831857	NaN	0.936089	NaN	7.91169	6.149387	NaN	NaN	2.405241	5.710631	...	152.884928	107195.823274	NaN	11.61
1999-12-02	18.026386	NaN	0.936089	NaN	8.45905	6.125041	NaN	NaN	2.345033	5.736608	...	152.884928	107588.242395	NaN	12.31
1999-12-03	18.253337	NaN	0.931683	NaN	8.828380	6.265031	NaN	NaN	2.271831	5.684654	...	157.279979	107588.242395	NaN	12.31
1999-12-06	17.442798	NaN	0.917734	NaN	8.890507	6.137214	NaN	NaN	2.308567	5.76416	...	161.378535	106803.404153	NaN	12.11
1999-12-07	17.637327	NaN	0.908924	NaN	9.04402	6.289377	NaN	NaN	2.283432	5.684654	...	158.959555	100066.875905	NaN	12.21
1999-12-08	19.647464	NaN	0.931683	NaN	8.847967	6.417193	NaN	NaN	2.283432	5.680324	...	158.858106	92766.957991	NaN	11.71
1999-12-09	18.642396	NaN	0.936089	NaN	8.00982	6.352271	NaN	NaN	2.295033	5.628370	...	157.230494	96763.315790	NaN	11.91
1999-12-10	19.063876	NaN	0.954443	NaN	7.790709	6.289377	NaN	NaN	2.295033	5.602393	...	153.327385	93657.363592	NaN	11.71
1999-12-13	18.836925	NaN	0.945633	NaN	7.60002	6.289377	NaN	NaN	2.405241	6.104619	...	152.094825	94834.620956	NaN	11.51
1999-12-14	19.647464	NaN	0.940494	NaN	7.28296	6.214310	NaN	NaN	2.525116	5.912389	...	150.909670	94049.782714	NaN	11.31
1999-12-15	19.452934	NaN	0.945633	NaN	7.444648	5.997225	NaN	NaN	2.525116	6.302045	...	157.279979	96584.862385	NaN	11.41
1999-12-16	21.074012	NaN	0.950038	NaN	7.54705	6.086494	NaN	NaN	2.416842	6.104619	...	152.537282	94049.782714	NaN	11.51
1999-12-17	21.289542	NaN	0.954443	NaN	7.67678	6.125041	NaN	NaN	2.476779	5.912389	...	148.642075	96862.119749	NaN	11.51
1999-12-20	20.263473	NaN	1.009507	NaN	7.52325	6.125041	NaN	NaN	2.453578	5.844846	...	148.144310	96077.281507	NaN	11.31
1999-12-21	20.263473	NaN	1.027862	NaN	7.768870	6.149387	NaN	NaN	2.525116	5.782503	...	146.189052	102028.971511	NaN	11.21
1999-12-22	18.836925	NaN	1.027862	NaN	7.67218	6.098667	NaN	NaN	2.490314	5.325306	...	147.006562	103206.228875	NaN	11.31
1999-12-23	19.874415	NaN	1.009507	NaN	7.794547	6.112668	NaN	NaN	2.585054	5.325306	...	150.024755	105233.272668	NaN	11.21
1999-12-27	19.874415	NaN	1.027862	NaN	7.62381	6.137214	NaN	NaN	2.683661	5.553905	...	147.899378	110008.160309	NaN	11.11
1999-12-28	19.452934	NaN	1.046217	NaN	7.537783	6.098667	NaN	NaN	2.925345	5.664941	...	146.564104	113997.754708	NaN	11.11

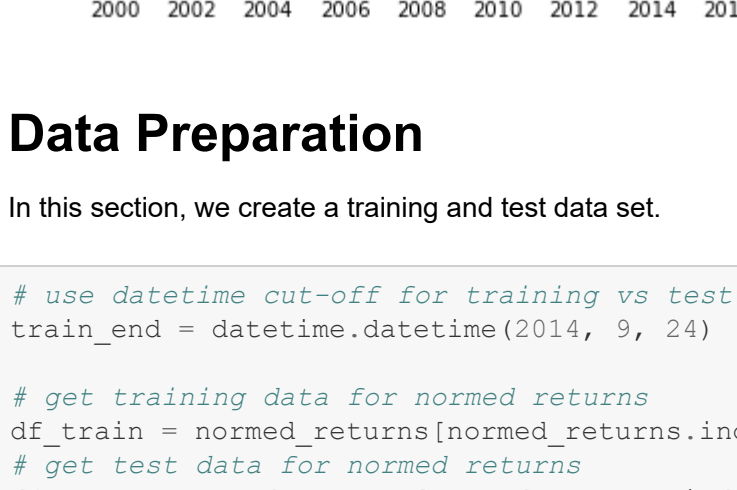
40 rows × 1256 columns

Plotting functionality

Just to get a visual idea of the data, we plot the prices for a particular symbol. We will finally plot this for the asset returns later on as well. Take a look at the plot for AAPL.

```
In [4]: # plotting function
def plot_symbol(symbol, df, csum=False):
    # csum denotes cumulative summation (useful for returns)
    yvals = df[symbol]
    if csum: yvals = np.cumsum(yvals)
    plt.plot(df.index, yvals)
    title = symbol + ' ' + df.name
    plt.title(title)

# check price chart for
plot_symbol('AAPL', close_prices)
```

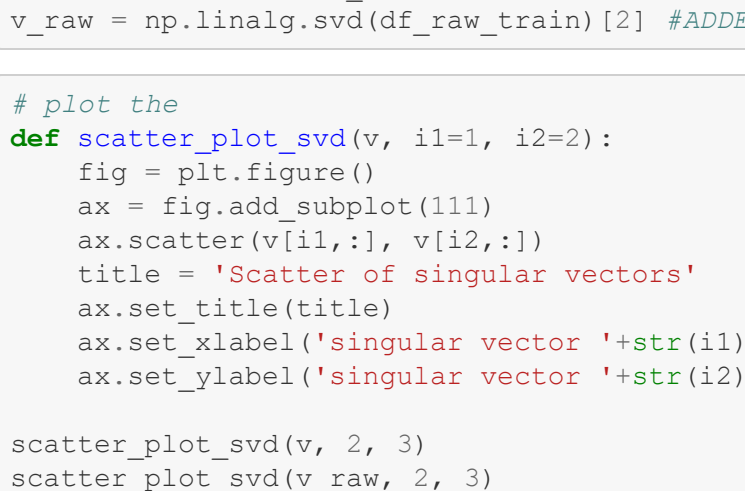


Asset Returns Transform

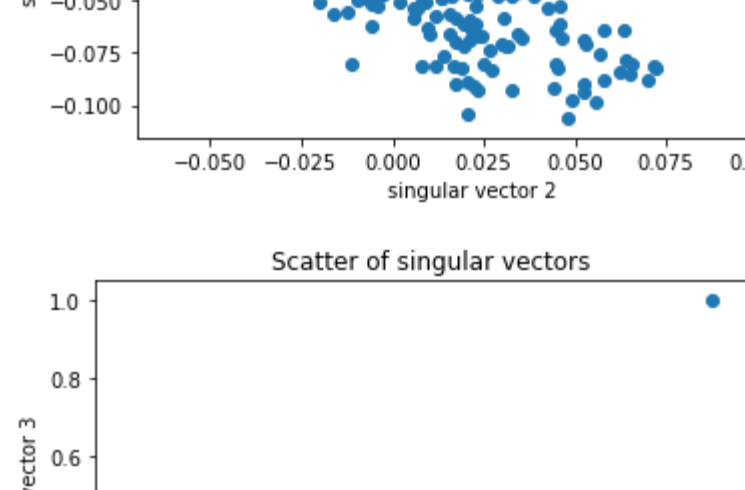
We calculate the returns by the day-to-day percent change of an asset. Once we have the returns, we normalize them by normalizing each individual asset's mean and standard deviation. Why would we want to normalize the returns in this manner?

```
In [5]: # calculate the percent change of each asset (pandas as an easy way to do this....)
returns = close_prices.pct_change().dropna(axis=0, how='all')
normed_returns = (returns - returns.mean()) / returns.std()
normed_returns = normed_returns.dropna(axis=0, how='all')
returns.name = 'returns'
normed_returns.name = 'normalized returns'
```

```
In [6]: # plot returns
plot_symbol('AMZN', returns, csum=True)
```



```
In [7]: # plot normalized returns
plot_symbol('AMZN', normed_returns, csum=True)
```



Data Preparation

In this section, we create a training and test data set.

```
In [8]: # use datetime cut-off for training vs test data set
train_end = datetime.datetime(2014, 9, 24)

# get training data for normed returns
df_train = normed_returns[normed_returns.index <= train_end].copy().dropna(axis=1, how='any')
# get test data for normed returns
df_test = normed_returns[normed_returns.index > train_end].copy().dropna(axis=1, how='any')
df_test = df_test[df_train.columns] # retain same tickers in test data as in training

# get training data for regular returns
df_raw_train = returns[returns.index <= train_end].copy().dropna(axis=1, how='any')
# get test data for regular returns
df_raw_test = returns[returns.index > train_end].copy().dropna(axis=1, how='any')
df_raw_test = df_raw_test[df_train.columns] # retain same tickers in test data as in training

print('Train dataset:', df_train.shape)
print('Test dataset:', df_test.shape)

Train dataset: (257, 1073)
Test dataset: (3742, 1073)
```

Computing SVD of training data

Consider our training data matrix as a $T \times N$ matrix X with N samples (our tickers) and T variables (our dates). If we assume that X is normalized as we have done above, we can calculate the empirical correlation matrix of our data by

$$C = \frac{1}{T-1} X^T X \in \mathbb{R}^{N \times N}$$

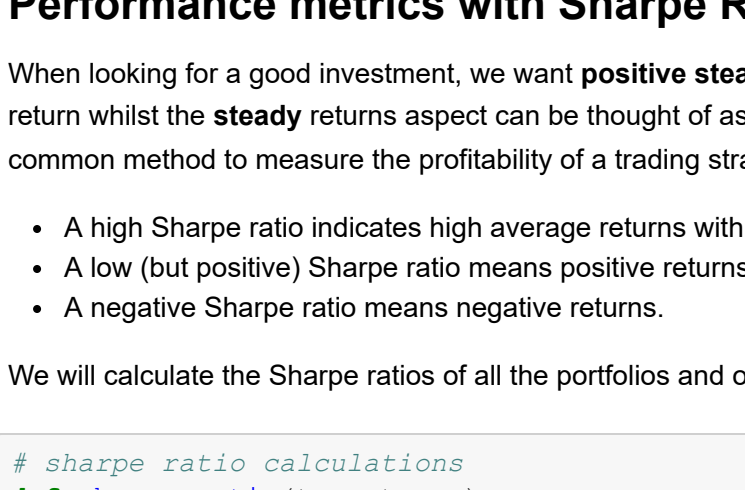
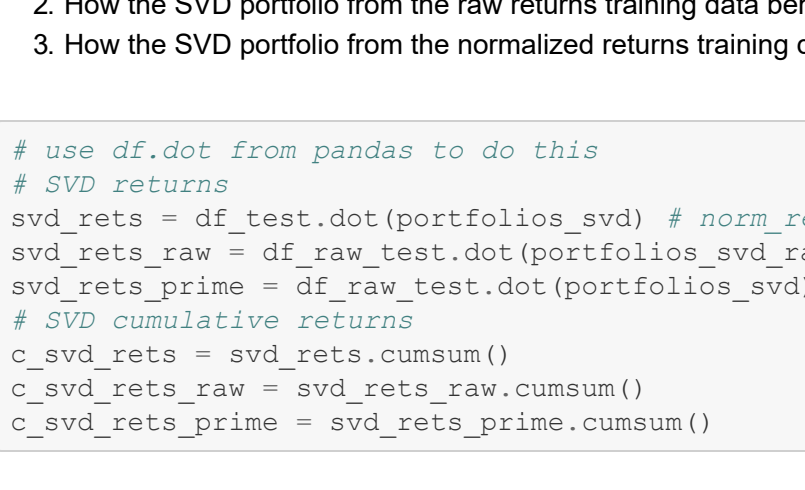
with eigendecomposition $C = V \Lambda V^T$. The eigenvectors of C should tell us how the stocks correlate to each other. Notice, however, that the eigenvectors V of C are just the right singular vectors of $X = U \Sigma V^T$. This means that we only need to compute the SVD of our data to be able to investigate the correlation of the stocks. Let's calculate the SVD using `numpy`. Call the right singular vectors v for

after gathering the singular vectors, we proceed to create a scatter plot of the singular vectors. Try to plot different singular vectors and comment on the behavior of the singular vectors for the normalized vs raw return SVDs.

```
In [9]: # calculate SVD here
v = np.linalg.svd(df_train)[2] # ADDED THIS
v_raw = np.linalg.svd(df_raw_train)[2] # ADDED THIS
```

```
In [13]: # plot the
def scatter_plot_svd(v, i1=1, i2=2):
    fig = plt.figure()
    ax = fig.add_subplot(111)
    ax.scatter(v[i1,:], v[i2,:])
    title = 'Scatter of singular vectors'
    ax.set_title(title)
    ax.set_xlabel('singular vector ' + str(i1))
    ax.set_ylabel('singular vector ' + str(i2))

scatter_plot_svd(v, 2, 3)
scatter_plot_svd(v_raw, 2, 3)
```



Generate the eigenportfolios by normalizing

We will define the j th eigenportfolio $Q^{(j)}$ by simply the normalized j th right singular vector $v^{(j)}$ so that the resultant $Q^{(j)}$ sums to 1. In particular, compute

$$Q^{(j)} = \frac{1}{\sum_{k=1}^N v_k^{(j)}} v^{(j)}.$$

We do this in the function below and compile the portfolios into a single pandas dataframe.

```
In [15]: # function to get all eigenportfolios up to jmax
def j_eigPortfolio(a_vecs, tickers=df_train.columns):
    # make empty portfolio list
    portfolios = []
    for j in range(len(tickers)):
        # normalize singular vector to sum to 1
        # calculate the jth eigenportfolio and call it j_port
        j_port = 1/np.sum(a_vecs[j])*a_vecs[j] # ADDED THIS
        j_port = pd.DataFrame(j_port,index=tickers, columns=['Q_'+str(j+1)])
        portfolios.append(j_port)
    portfolios = pd.concat(portfolios, axis=1)
    return portfolios

portfolios_svd = j_eigPortfolio(v) # SVD portfolio computed from normalized returns
portfolios_svd_raw = j_eigPortfolio(v_raw) # SVD portfolio computed from raw returns
```

```
In [16]: # view data
display(portfolios_svd_raw.head(10))
display(portfolios_svd.head(10))
```

	mean returns	cumulative returns	vol	sharpe
Q_1020	0.049062	12.363619	0.214973	0.228224
Q_644	0.189340	47.713793	0.898191	0.210802
Q_300	0.019358	4.878180	0.096674	0.200238
Q_872	0.022231	5.602198	0.115524	0.199337
Q_909	0.042788	10.782517	0.231836	0.184560
SVD prime sharpe				
	mean returns	cumulative returns	vol	sharpe