

Lab Course Machine Learning

Exercise Sheet 2

November 19th, 2025

Syed Wasif Murtaza-Jafri-311226

Part A: (interesting stats)

Maximum sale recorded

```
In [2]: import pandas as pd

#store with max sales
maxStoreId = td['Sales'].idxmax()
print(td.iloc[maxStoreId],[0,1,3,4])

Store DayOfWeek Sales Customers
44393 909 1 4151 1721
```

least possible and maximum possible competition distance

```
In [4]: #store with max competitionDistance
maxcompId = td['CompetitionDistance'].idxmax()
print('maximum possible competition distance')
print(td.iloc[maxcompId],[0,1,3])

maximum possible competition distance
Store StoreType CompetitionDistance
452 4933 a 75860.0

In [8]: #store with min competitionDistance
print('least possible competition distance')
minId = td['CompetitionDistance'].idxmin()
print('maximum possible competition distance')
print(td.iloc[min],[0,1,3])

Least possible competition distance
maximum possible competition distance
Store StoreType CompetitionDistance
515 516 a 20.0
```

Difference in the mean of sales (across all stores) when offering a Promo and not?

```
In [10]: promoData = td[td["Promo"]==1]
promaData.groupby(["Store"]).agg(["mean", 'count'])
#mean sales when promotion is running
promaData = td[td["Promo"]==1]
meanPromot = promoData['Sales'].mean()

noPromoData = td[td["Promo"]==0]
meanNoPromot = noPromoData['Sales'].mean()

print('Difference in the mean of sales:', (meanPromot-meanNoPromot))

Difference in the mean of sales:: 3585.1012408091174
```

Are there any anomalies in the data as in where the store was "Open" but had no sales recorded?

```
In [13]: nd = td.query('Open == 1 & Sales == 0')
nd = nd.reindex(sorted(nd.columns), axis=1)
nd.sort_values(by="Date").head()

Out [13]: Store DayOfWeek Date Sales Customers Open Promo StateHoliday SchoolHoliday
995016 762 4 2013-01-17 0 0 1 0 0 0 0
990681 232 4 2013-01-24 0 0 1 1 0 0 0
984098 339 3 2013-01-30 0 0 1 0 0 0 0
982983 339 4 2013-01-31 0 0 1 0 0 0 0
975098 259 4 2013-02-07 0 0 1 1 0 0 0
```

```
In [14]: nd = td.query('Open == 0 & Sales != 0')
nd = nd.reindex(sorted(nd.columns), axis=1)
nd.sort_values(by="Date").head()

Out [14]: Store DayOfWeek Date Sales Customers Open Promo StateHoliday SchoolHoliday
995016 762 4 2013-01-17 0 0 1 0 0 0 0
990681 232 4 2013-01-24 0 0 1 1 0 0 0
984098 339 3 2013-01-30 0 0 1 0 0 0 0
982983 339 4 2013-01-31 0 0 1 0 0 0 0
975098 259 4 2013-02-07 0 0 1 1 0 0 0
```

Yes, there exist anomaly, when the store is open and there were no sales recored, but anomaly for vice versa not exist)

Which store type ('a','b' etc.) has had the most sales?

```
In [17]: td = pd.read_csv('train.csv', low_memory=False)
sd = pd.read_csv('store.csv', low_memory=False)
df = sd.append(td, sort=False)
df = pd.merge(td, sd, how='left', on='Store')
print(df.groupby(['StoreType'])['Sales'].sum())

StoreType
a 3165394859
b 159231395
c 76321426
d 1765592943
Name: Sales, dtype: int64
```

Part B: (Plotting)

Monthly mean of sales

```
In [96]: td["Month"] = pd.DatetimeIndex(td["Date"]).month

In [18]: td

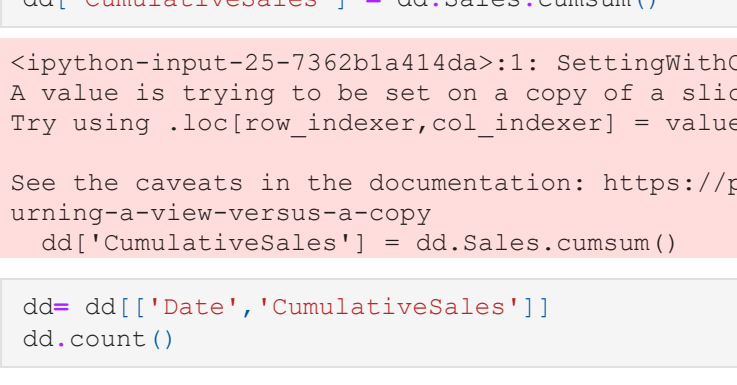
Out [18]: Store DayOfWeek Date Sales Customers Open Promo StateHoliday SchoolHoliday
0 1 5 2015-07-31 5263 555 1 1 0 1
1 2 5 2015-07-31 6064 625 1 1 0 1
2 3 5 2015-07-31 8314 821 1 1 0 1
3 4 5 2015-07-31 13995 1498 1 1 0 1
4 5 5 2015-07-31 4822 559 1 1 0 1
... ..
1017204 1111 2 2013-01-01 0 0 0 0 a 1
1017205 1112 2 2013-01-01 0 0 0 0 a 1
1017206 1113 2 2013-01-01 0 0 0 0 a 1
1017207 1114 2 2013-01-01 0 0 0 0 a 1
1017208 1115 2 2013-01-01 0 0 0 0 a 1
1017209 rows x 9 columns
```

```
In [98]: import matplotlib as plt
import matplotlib.pyplot as plt

monthData = pd.DataFrame(td.groupby(td["Month"])["Sales"].agg(["mean"]))
monthData.plot(kind='line')
monthData.head()

Out [98]: mean

monthYear
2013-01 5211.555578
2013-02 5494.37397
2013-03 5480.349168
2013-04 5823.749836
2013-05 5364.127383
```

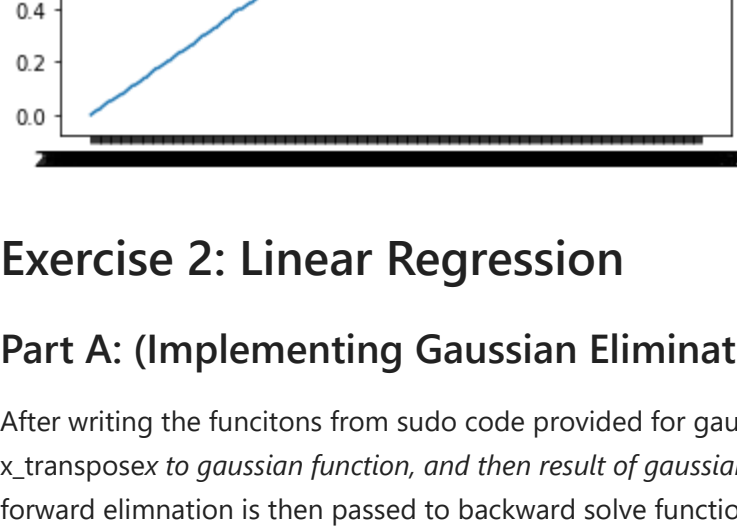


Daily mean of sales

```
In [99]: monthData = pd.DataFrame(td.groupby(td["Date"])["Sales"].agg(["mean"]))
monthData.plot(kind='line')
monthData.head()

Out [99]: mean

Date
2013-01-01 87.284560
2013-01-02 6233.030493
2013-01-03 5693.112108
2013-01-04 59542.18834
2013-01-05 5337.751570
```



Cumulative sales for the first year.

```
In [24]: td["Year"] = td["Date"].astype("datetime64[ns]").dt.strftime('%Y')
td["Day"] = td["Date"].astype("datetime64[ns]").dt.strftime('%d')

ss = td.query("Store==1 & Year == '"+td["Year"].min()+"")
pd = ss.sort_values(by="Date")
dd = pd["Date", "Sales"]

#groupby(['Month', 'Day']).sum().groupby(level=0).cumsum()
#data['date'] = pd.to_datetime(data[['Month', 'Day']]).assign(YEAR=td["Year"].min())
#data.query("Month == 1 | Month == 2")

In [25]: dd["CumulativeSales"] = dd.Sales.cumsum()

<ipython-input-25-7362b1a414da>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using loc[row_index,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ret
urning-a-view-versus-a-copy
dd["CumulativeSales"] = dd.Sales.cumsum()

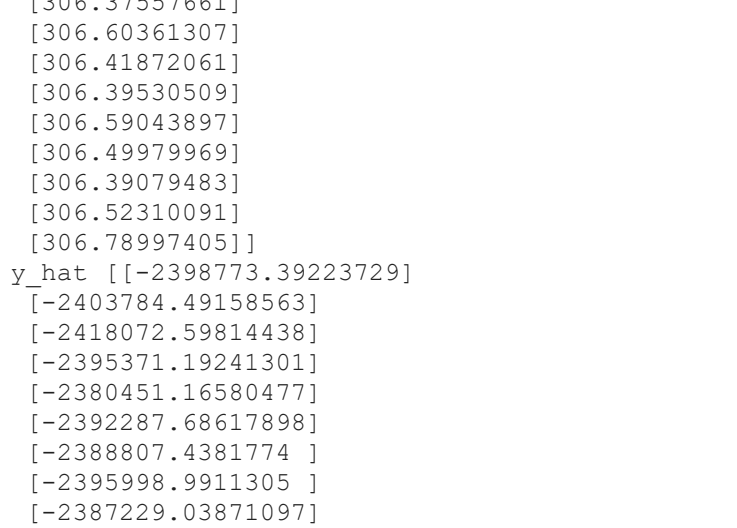
In [26]: dd = dd["Date", "CumulativeSales"]
dd.count()
```

```
Out [26]: Date 365
CumulativeSales 365
dtypes: object 1
```

```
In [19]: import matplotlib.pyplot as plt
type(ddaxes.style.hide_index())
fig = plt.figure()
ax = fig.add_subplot(1,1,1)
#to get a tick every 15 minutes
ax.xaxis.set_major_formatter(mdates.DateFormatter('%H:%M'))

In [27]: import matplotlib.pyplot as plt
plt.plot(dd.Date, dd.CumulativeSales)

Out [27]: [matplotlib.lines.Line2D at 0x2e4802a1910]
```



Exercise 2: Linear Regression

Part A: (Implementing Gaussian Elimination)

After writing the functions from sudo code provided for gaussian elimination, forward elimination and backward solve, i passed x transpose to gaussian function, and then result of gaussian elimination and x transpose was passed to forward elimination, and result from forward elimination is then passed to backward solve function, which gives the β parameters matrix.

```
In [32]: import numpy as np
rows = 100
cols = 10
mean, sd = 2, 0.01
def transpose(matrix):
    return [[row[i] for row in matrix] for i in range(len(matrix[0]))]

def multiplyMatrix(X, Y):
    return [[sum(a*b for a,b in zip(X_row,Y_col)) for Y_col in zip(*Y)] for X_row in X]

x = np.random.normal(mean, sd, (rows,cols))
y = np.random.uniform(1,2, (rows,1))
bias_column = np.ones(shape=(rows,1))
A = np.append(bias_column,x,axis=1)

x_tran = np.transpose(A)
x_tran_x = np.dot(x_tran,A)

x_tran_y = np.dot(x_tran,y)
print(x_tran_y)

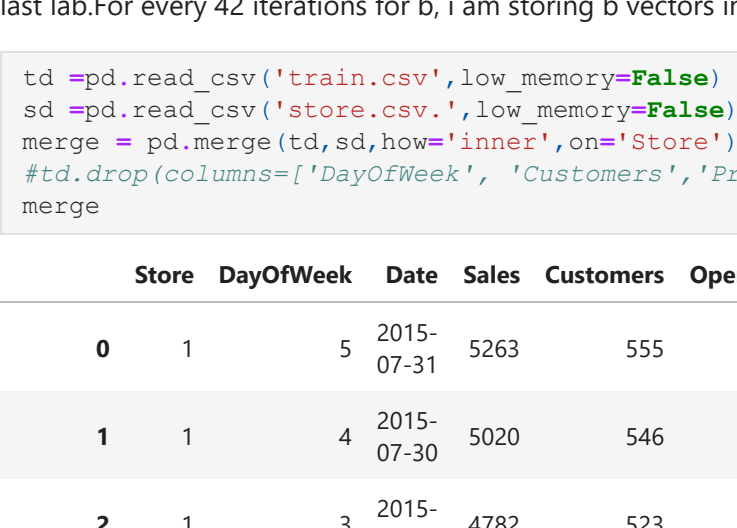
def Gauss(mat):
    r,c = mat.shape
    if(r==c):
        for k in range(0,r-1):
            for i in range(k+1,r):
                mat[i][k] = mat[i][k]/mat[k][k]
            for j in range(k+1,r):
                mat[i][j] = mat[i][j]-mat[i][k]*mat[k][j]
    else:
        print('Dimensions of matrices are not correct for Gaussian Elimination')
    return mat
def forward(mat,b):
    r,c = mat.shape
    for k in range(0,r-1):
        for i in range(k+1,r):
            b[i] = b[i]-mat[i][k]*b[k]
    return b
def backward(mat,b):
    r,c = mat.shape
    x = np.zeros(shape=(r,1))
    for i in range(r-1,-1,-1):
        s = b[i]
        for j in range(i+1,r):
            s = s-mat[i][j]*x[j]
        x[i] = s/mat[i][i]
    return b
B=backward(Gaussian(x_tran_x),forward(Gaussian(x_tran_x),x_tran_y))

y_hat = np.dot(A,B)
print('y_hat',y_hat)

[[115.2729445]
 [306.1780161]
 [306.3755766]
 [306.6036130]
 [306.4187061]
 [306.3953050]
 [306.5904389]
 [306.4597969]
 [306.3907948]
 [306.5231091]
 [306.7899705]
 y_hat [[-2398773.39223729]
 [-2403784.49158563]
 [-2418072.59814438]
 [-2395371.19241301]
 [-2380451.16980477]
 [-232287.68617398]
 [-238807.4381774]
 [-2395998.9911305]
 [-2387229.03871097]
 [-2401302.3057517]
 [-2422441.867165]
 [-2407159.45915848]
 [-2384461.83013668]
 [-2404607.76824108]
 [-2384511.59261865]
 [-2386348.63171126]
 [-2413701.83099992]
 [-2394220.66988416]
 [-2390835.4564557]
 [-2425580.58916502]
 [-2403179.00413475]
 [-242614.76968767]
 [-2392677.6845118]
 [-2404985.43616723]
 [-239569.6569907]
 [-2388076.2449239]
 [-2396098.75123243]
 [-2375818.7913636]
 [-2407797.6858879]
 [-2395032.78079963]
 [-2407903.53513572]
 [-2398052.43754703]
 [-2391324.98440729]
 [-2391231.41576494]
 [-2406592.56437425]
 [-2393798.26711161]
 [-2396643.12315594]
 [-2397382.17766941]
 [-2419269.58364105]
 [-2417307.59994118]
 [-2397680.64369447]
 [-2417899.24682773]
 [-2393141.24132567]
 [-2409383.66866789]
 [-2390608.4150301]
 [-2402092.12800221]
 [-2374266.58980215]
 [-24134.6620304]
 [-2399531.27005681]
 [-2387730.64251461]
 [-2385490.33016484]
 [-2403136.04836872]
 [-2371759.7022767]
 [-239322.7824001]
 [-2414006.6839411]
 [-238446.13034541]
 [-239891.53525747]
 [-2398104.42029671]
 [-2395247.41835434]
 [-238609.5576359]
 [-2401126.2957042]
 [-2415182.59316823]
 [-2393164.5605238]
 [-2409335.45551809]
 [-2392688.01692832]
 [-2395433.2055266]
 [-2375961.00367086]
 [-2415669.27989177]
 [-239569.6929945]
 [-2402014.91284942]
 [-2414939.73714255]
 [-2411916.2319089]
 [-240262.2133911]
 [-2398821.44482914]
 [-2378077.50323884]
 [-2400428.33545187]
 [-2421823.8081093]
 [-2380411.13406952]
 [-2387193.5110024]
 [-2394784.86629732]
 [-2379844.54765447]
 [-2389576.58900697]
 [-2416137.168585]
 [-2405247.89603803]
 [-2391015.55409316]
 [-2393182.86325146]
 [-2383477.12301829]
 [-2400502.08289589]
 [-2403170.34866473]
 [-2430444.79322164]
 [-2411023.4255056]
 [-240352.5458206]
 [-2403996.45349238]
 [-2390270.79340889]
 [-2405869.49791777]
 [-2387888.80213353]
 [-2407955.48708516]
 [-238145.37283155]
 [-2416177.94661667]
 [-2386160.19769473]]

In [33]: import matplotlib.pyplot as plt
axis = np.linspace(0, 1, 100)
plt.scatter(axis, y, c='b', marker='o', label='Y')
plt.scatter(axis, y_hat, c='g', marker='s', label='y_hat')
plt.show()

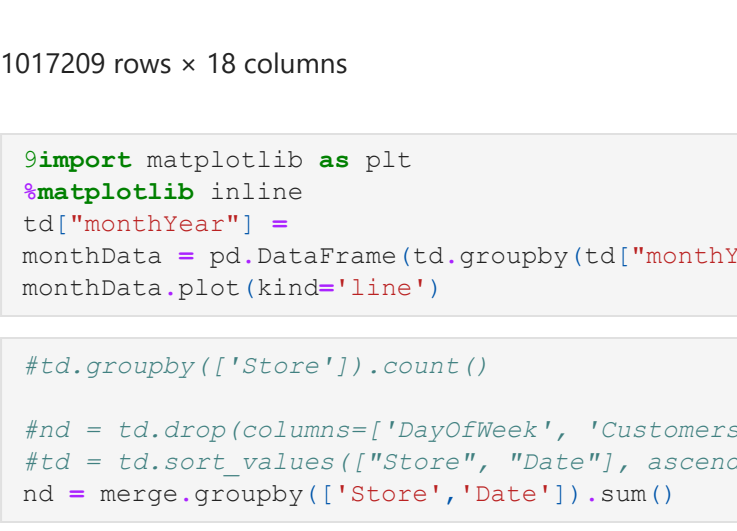
Out [33]:
```



```
In [34]: B = np.linalg.lstsq(A, y, rcond=None)[0]
y_hat = np.dot(A,B)

import matplotlib.pyplot as plt
axis = np.linspace(0, 1, 100)
plt.scatter(axis, y, c='b', marker='o', label='Y')
plt.scatter(axis, y_hat, c='g', marker='s', label='y_hat')
plt.show()

Out [34]:
```



Part B: Multiple Linear (Auto)Regression

First i merged the datasets of train and store, then i grouped the dataframe by store and date columns. I am then looping through for every store and converting dataframes of Sales columns for every store into first row of matrix which is already sorted by days. There were only 10 records which have in total 942 days sales. For store which have less than 942 records, I am appending zeros and in that way creating x matrix. After which creating x train, x test, y train, y test matrices from x and finding values of β matrix from solving for linear regression as last lab. For every 42 iterations for β , I am storing β vectors in an matrix. After that finding RMSE.

```
In [27]: td = pd.read_csv('train.csv', low_memory=False)
sd = pd.read_csv('store.csv', low_memory=False)
merge = pd.merge(td, sd, how='inner', on='Store')
#td.drop(columns=["DayOfWeek", 'Customers', 'Promo', 'SchoolHoliday', 'StateHoliday', 'Open'])
merge

Out [27]: Store DayOfWeek Date Sales Customers Open Promo StateHoliday SchoolHoliday StoreType CompetitionDistance
0 1 5 2015-07-31 5263 555 1 1 0 1 c a 1270
1 1 4 2015-07-30 5020 546 1 1 0 1 c a 1270
2 1 3 2015-07-29 4782 523 1 1 0 1 c a 1270
3 1 2 2015-07-28 5011 560 1 1 0 1 c a 1270
4 1 1 2015-07-27 6102 612 1 1 0 1 c a 1270
... ..
1017204 1115 6 2013-01-05 4771 339 1 0 0 1 d c 5350
1017205 1115 5 2013-01-04 4540 326 1 0 0 1 d c 5350
1017206 1115 4 2013-01-03 4297 300 1 0 0 1 d c 5350
1017207 1115 3 2013-01-02 3697 305 1 0 0 1 d c 5350
1017208 1115 2 2013-01-01 0 0 0 0 a 1 d c 5350
1017209 rows x 18 columns

In [ ]: import matplotlib.pyplot as plt
import matplotlib.pyplot as plt

monthData = pd.DataFrame(td.groupby(td["MonthYear"])["Sales"].agg(["mean"]))
monthData.plot(kind='line')

In [142]: #td.groupby(["Store","Date"])

#nd = td.drop(columns=["DayOfWeek", 'Customers', 'Promo', 'SchoolHoliday', 'Open'])
#nd = td.sort_values(["Store", "Date"], ascending = (True, True)).reset_index()
nd = merge.groupby(["Store", "Date"]).sum()

In [143]: nd

Out [143]: DayOfWeek Sales Customers Open Promo SchoolHoliday CompetitionDistance CompetitionOpenSinceMonth CompetitionOp
Store Date
1 2013-01-01 2 0 0 0 0 1 1270.0 9.0
2 2013-01-02 3 5530 668 1 0 1 1270.0 9.0
2013-01-03 4 4327 578 1 0 1 1270.0 9.0
2013-01-04 5 4486 619 1 0 1 1270.0 9.0
2013-01-05 6 4997 635 1 0 1 1270.0 9.0
... ..
1115 2015-07-27 1 10712 608 1 1 1 5350.0 0.0
2015-07-28 2 8093 500 1 1 1 5350.0 0.0
2015-07-29 3 7661 473 1 1 1 5350.0 0.0
2015-07-30 4 8405 502 1 1 1 5350.0 0.0
2015-07-31 5 8680 538 1 1 1 5350.0 0.0
1017209 rows x 12 columns

In [145]: rows = nd.index
x = np.zeros(shape=(rows,942))
print(rows)
index = 0
for store, new_df in nd.groupby(level=0):
    number_column = new_df.loc[:, 'Sales']
    B = np.zeros(shape=(1,942))
    numbers = number_column.values
    if (len(numbers)<942):
        a = np.zeros((1,942-len(numbers)))
        numbers = np.append(numbers,a)
    x[index,:] = numbers
    index+=1
print(x)

1017209
[[ 0. 5530. 4327. ... 4782. 5020. 5263.]
 [ 0. 4422. 4159. ... 6402. 5567. 6064.]
 [ 0. 6823. 5902. ... 7610. 8977. 8314.]
 ...
 [ 0. 0. 0. ... 0. 0. 0.]
 [ 0. 0. 0. ... 0. 0. 0.]
 [ 0. 0. 0. ... 0. 0. 0.]]

In [151]: x_train = x[:800,0:900]
y_train = y[:800,0:900]
x_test = x[800:1000,0:900]
y_test = y[800:1000,0:900]
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)

bias_column = np.ones(shape=(800,1))
A = np.append(bias_column,x_train,axis=1)
bias_column = np.ones(shape=(800,1))
B = np.zeros(shape=(42,901))

for i in range(0,42):
    aTrain = np.transpose(A)
    aTrain_a = np.matmul(aTrainpose,a)
    aTrain_inv = np.linalg.inv(aTrain_a)
    aTrain_y = np.matmul(aTrainpose,y_train[:,i])
    B[i,:] = np.matmul((aTrain_inv),(aTrain_y))

print(B.shape)

(800, 901)
(800, 42)
(200, 900)
(200, 42)
(42, 901)

In [179]: bias_column = np.ones(shape=(200,1))
A_train = np.append(bias_column,x_test,axis=1)
print(A_train.shape)
y_hat_test = np.zeros(shape=(200, 42))
for i in range(0,42):
    y_hat_test[:,i] = y

print(y_hat_test)

(200, 901)
[[ 68500.46969727 71262.44213867 43189.03271484 ... 59319.96826172
 83573.36230469 102244.1328125]
 [ 49287.37988281 9177.03100586 12270.60095215 ... 44640.62719727
 83267.0237952 112969.09106445]
 [-2809.5612793 40208.50805664 -30184.32879639 ... -34250.36120605
 12773.62561035 -2710.7222168]
 ...
 [ 56438.08959961 27814.15600586 41082.89086914 ... 68104.33866719
 97599.09545898 118237.27648926]
 [ 59640.3200781 12844.22563281 -35911.95501709 ... -24306.86535645
 31095.93347168 58813.05285645]
 [150206.48754883 57331.34912109 81992.12310791 ... 134882.63391113
 128521.21032715 163513.37634277]]

In [179]: y_minus_y_hat = y_test-y_hat_test
RMSE = np.ones(shape=(42,1))
print(y_minus_y_hat.shape)
for i in range(42):
    y = y_minus_y_hat[:,i]
    rmse = 0
    for j in range(200):
        rmse = rmse + (y[j]*y[j])
    rmse = (rmse ** 0.5)/2
    RMSE[i]=rmse

print(RMSE)

(200, 42)
[[1.0279270e+18]
 [4.7423602e+17]
 [1.10221415e+18]
 [1.03057074e+18]
 [1.05315770e+18]
 [1.11325267e+18]
 [1.05960075e+18]
 [5.01515025e+17]
 [1.49972188e+18]
 [1.43302062e+18]
 [1.38114417e+18]
 [1.31536365e+18]
 [1.06342052e+18]
 [5.4996458e+17]
 [1.20883618e+18]
 [1.07184859e+18]
 [1.04257319e+18]
 [1.09988145e+18]
 [1.09830368e+18]
 [9.0960075e+17]
 [5.26430864e+17]
 [1.65334672e+18]
 [1.44202199e+18]
 [1.57776634e+18]
 [1.40318692e+18]
 [1.62389346e+18]]
```