

```
In [4]: 1 import warnings
        2 warnings.filterwarnings('ignore')
        3 import pandas as pd, numpy as np
```

CQF Exam 3 - 2021/5/23

delegate: Lu, Zhenyang, [zzhenyanglu@gmail.com \(mailto:zzhenyanglu@gmail.com\)](mailto:zzhenyanglu@gmail.com), +1-312-804-2532

Section 0 - data preprocessing

Group A Data Source:

Hong Kong Heng Seng Index (HSI) daily price(open, high, low, close) from 1999/10/25 to 2021/05/21

Source: Proprietary database

for Group A we create the following features:

- 1. sign: sign of daily return
- 2. return_t: T-t return, where t = 1,2,3,4,5
- 3. momentum_t: price change period t days, where t= 1,2,3,4,5
- 4. ewma: from https://en.wikipedia.org/wiki/EWMA_chart (https://en.wikipedia.org/wiki/EWMA_chart), $ewma_t = P_t + \lambda * ewma_{t-1}$, where $\lambda = 0.8$
- 5. O-C: Open - Close
- 6. H-L: High - Low
- 7. ewma-price: sing of ewma-price

dependent variable (sign_1 in the actual data below): T+1 daily return's sign, 1 if postive return, 0 otherwise

below is quick peek at group A data

In [5]:

```
1 # group_a
2 group_a = pd.read_csv(r'./HSI.csv')
3 shift_days = 5
4 _lambda = 0.4 #EWMA model's argument
5
6 for i in range(1,1+shift_days):
7     group_a[f'return_{i}'] = np.log(group_a['close']/group_a['close'].shift(i))
8     group_a[f'momentum_{i}'] = group_a['close'] - group_a['close'].shift(i)
9     group_a[f'sign_{i}'] = group_a[f'return_{i}'] >= 0
10    #group_a[f'sign_{i}'] = group_a[f'sign_{i}'].astype(int)
11    group_a[f'sign_{i}'] = group_a[f'sign_{i}'].apply(lambda x:1 if x else 0)
12
13    group_a[f'return_{i}'] = group_a[f'return_{i}'].shift(1)
14    group_a[f'momentum_{i}'] = group_a[f'momentum_{i}'].shift(1)
15    group_a[f'sign_{i}'] = group_a[f'sign_{i}'].shift(1)
16
17 group_a['O-C'] = group_a['open'] - group_a['close']
18 group_a['H-L'] = group_a['high'] - group_a['low']
19
20 # EWMA
21 close_list = group_a['close'].tolist()
22 ewma_list = [close_list[0]]
23 for i in range(1, group_a['close'].shape[0]):
24     ewma_list.append(_lambda*group_a['close'].iloc[i] + (1-_lambda)*ewma_list[i-1])
25
26 group_a['ewma'] = ewma_list
27 group_a['ewma-price'] = (group_a['ewma'] - group_a['close']) >= 0
28 group_a['ewma-price'] = group_a['ewma-price'].astype(int)
29 group_a = group_a.iloc[shift_days+2:]
30 group_a.head(10)
```

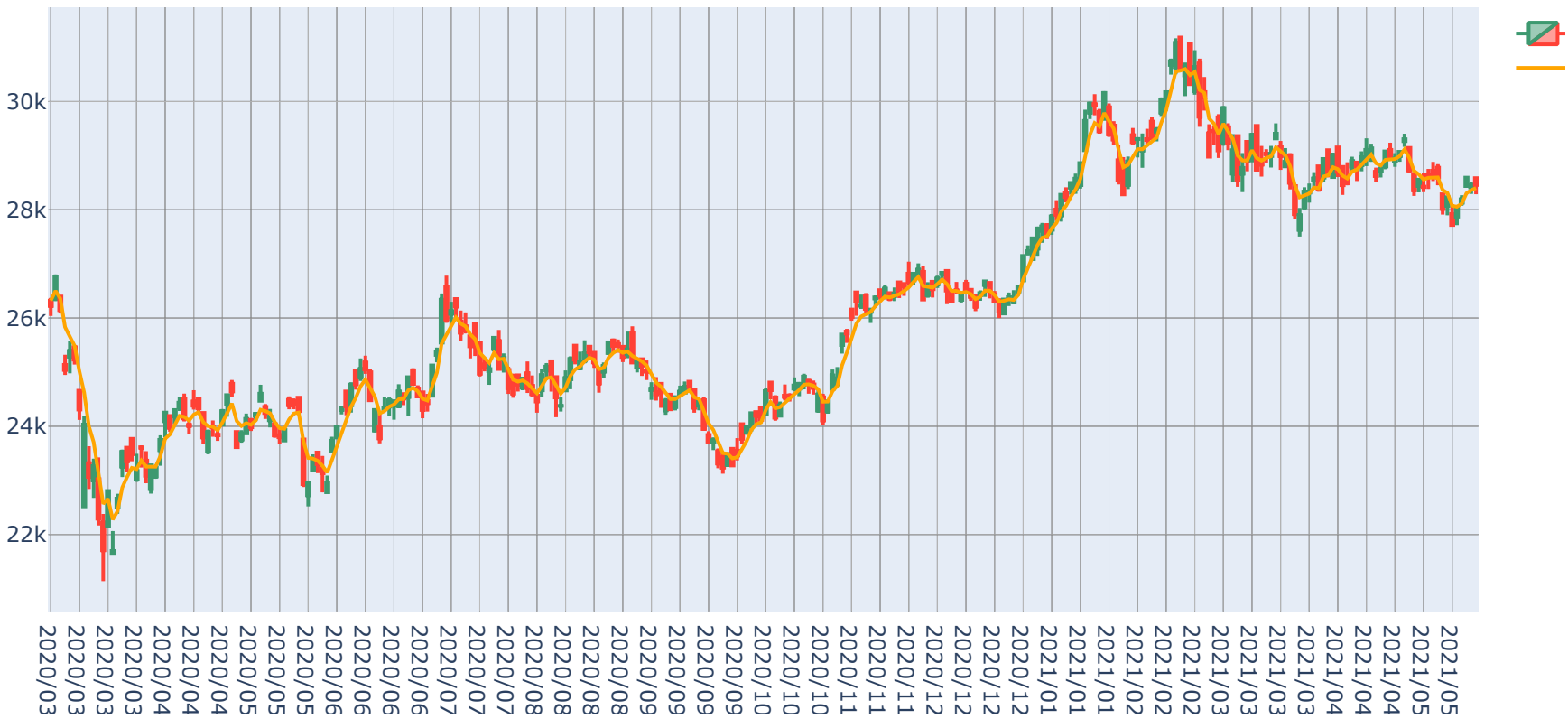
Out[5]:

	date	open	high	low	close	return_1	momentum_1	sign_1	return_2	momentum_2	...	return_4	momentum_4	sign
7	1999/11/03	13366.12	13497.36	13230.72	13257.33	0.001030	13.73	1.0	0.005933	78.89	...	0.048139	626.77	1
8	1999/11/04	13272.30	13693.39	13272.30	13651.51	-0.005905	-78.51	0.0	-0.004874	-64.78	...	0.038323	498.45	1
9	1999/11/05	13698.30	13730.89	13533.26	13610.27	0.029300	394.18	1.0	0.023395	315.67	...	0.029328	394.56	1
10	1999/11/08	13659.26	13756.84	13513.66	13521.11	-0.003025	-41.24	0.0	0.026274	352.94	...	0.021400	288.16	1
11	1999/11/09	13585.65	13735.07	13574.42	13669.70	-0.006572	-89.16	0.0	-0.009598	-130.40	...	0.013797	185.27	1
12	1999/11/10	13734.41	14005.26	13734.41	13975.54	0.010930	148.59	1.0	0.004357	59.43	...	0.030631	412.37	1
13	1999/11/11	14053.47	14217.36	13916.14	14105.71	0.022127	305.84	1.0	0.033056	454.43	...	0.023459	324.03	1
14	1999/11/12	14110.83	14268.39	13977.95	14189.67	0.009271	130.17	1.0	0.031398	436.01	...	0.035755	495.44	1
15	1999/11/15	14251.29	14634.99	14251.29	14562.22	0.005935	83.96	1.0	0.015206	214.13	...	0.048262	668.56	1
16	1999/11/16	14619.42	14726.65	14491.39	14689.46	0.025916	372.55	1.0	0.031851	456.51	...	0.063249	892.52	1

10 rows × 24 columns

GROUP A: below is a candlestick plot with EWMA average line of the most recent 500 trading days

```
In [6]: 1 import plotly.graph_objects as go
2 from plotly.subplots import make_subplots
3
4 fig = go.Figure(data=[go.Candlestick(x=group_a['date'][-300:],
5 open=group_a['open'][-300:],
6 high=group_a['high'][-300:],
7 low=group_a['low'][-300:],
8 close=group_a['close'][-300:], name = 'HSI'),
9 go.Scatter(x=group_a['date'][-300:], y=group_a['ewma'][-300:], line=dict(color='orange', width=2), r
10 )
11
12 fig.update_layout(xaxis_rangeslider_visible=False)
13 fig.show()
```



Group B:

St Louis Fed Relative Midwest Economy Index from 1990-01-01 to 2021-03-01

Source: Fred St Louis. <https://fred.stlouisfed.org/series/RMEIM683SFRBCHI> (<https://fred.stlouisfed.org/series/RMEIM683SFRBCHI>)

for Group B we create the following features:

- 1. sign: sign of monthly return
- 2. return_t: T-t return, where t=1,2,3,4,5,6
- 3. momentum_t: price change period t months, where t=1,2,3,4,5,6
- 4. ma_t: t month simple moving average, where t=1,2,3,4,5,6

dependent variable (return_1 in the actual data below): T+1 monthly return, with 10 bin bucketing (i.e. labels to be predicted is 0 thru 9) based on quantile

In [7]:

```
1 group_b = pd.read_csv(r'./RMEIM683SFRBCHI.csv')
2 shift_days = 6
3
4 for i in range(1,1+shift_days):
5     group_b[f'return_{i}'] = (group_b['value']/group_b['value'].shift(i)) -1
6     group_b[f'momentum_{i}'] = group_b['value'] - group_b['value'].shift(i)
7     group_b[f'sign_{i}'] = group_b[f'return_{i}'] >= 0
8     group_b[f'sign_{i}'] = group_b[f'sign_{i}'].apply(lambda x:1 if x else -1)
9     #group_b[f'sign_{i}'] = group_b[f'sign_{i}'].astype(int)
10    group_b[f'return_{i}']
11
12 # MA(5)
13 for i in range(1,1+shift_days):
14     group_b[f'ma_{i}'] = group_b['value'].rolling(i).mean()
15
16 group_b = group_b.iloc[shift_days:]
17 group_b.tail(10)
```

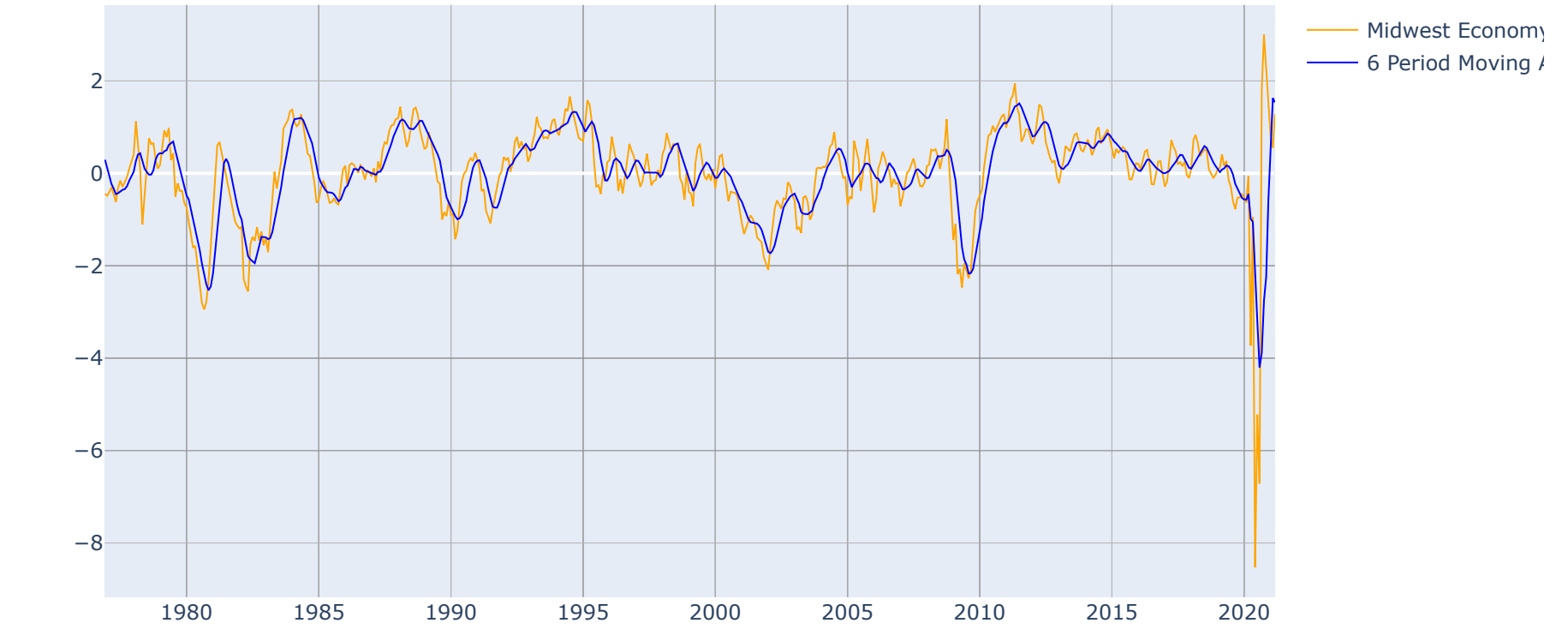
Out[7]:

	date	value	return_1	momentum_1	sign_1	return_2	momentum_2	sign_2	return_3	momentum_3	...	sign_5	return_6	momentu
528	2020-06-01	-8.52920	8.059161	-7.58770	1	1.287759	-4.80101	1	137.934680	-8.46781	...	1	17.967687	-8.0
529	2020-07-01	-5.22077	-0.387895	3.30843	-1	4.545162	-4.27927	1	0.400350	-1.49258	...	1	9.255102	-4.7
530	2020-08-01	-6.71652	0.286500	-1.49575	1	-0.212526	1.81268	-1	6.133850	-5.77502	...	1	9.647453	-6.0
531	2020-09-01	1.82300	-1.271420	8.53952	-1	-1.349182	7.04377	-1	-1.213736	10.35220	...	-1	-30.695390	1.8
532	2020-10-01	3.00375	0.647696	1.18075	1	-1.447218	9.72027	-1	-1.575346	8.22452	...	-1	-1.805686	6.7
533	2020-11-01	2.20800	-0.264919	-0.79575	-1	0.211190	0.38500	1	-1.328742	8.92452	...	-1	-3.345194	3.1
534	2020-12-01	1.38178	-0.374194	-0.82622	-1	-0.539982	-1.62197	-1	-0.242030	-0.44122	...	-1	-1.162006	9.9
535	2021-01-01	0.80070	-0.420530	-0.58108	-1	-0.637364	-1.40730	-1	-0.733433	-2.20305	...	-1	-1.153368	6.0
536	2021-02-01	0.55475	-0.307169	-0.24595	-1	-0.598525	-0.82703	-1	-0.748755	-1.65325	...	-1	-1.082595	7.2
537	2021-03-01	1.28082	1.308824	0.72607	1	0.599625	0.48012	1	-0.073065	-0.10096	...	-1	-0.297411	-0.5

10 rows × 26 columns

GROUP B - a quick plot of Midwest Economy Index and its 6 periods Moving Average line

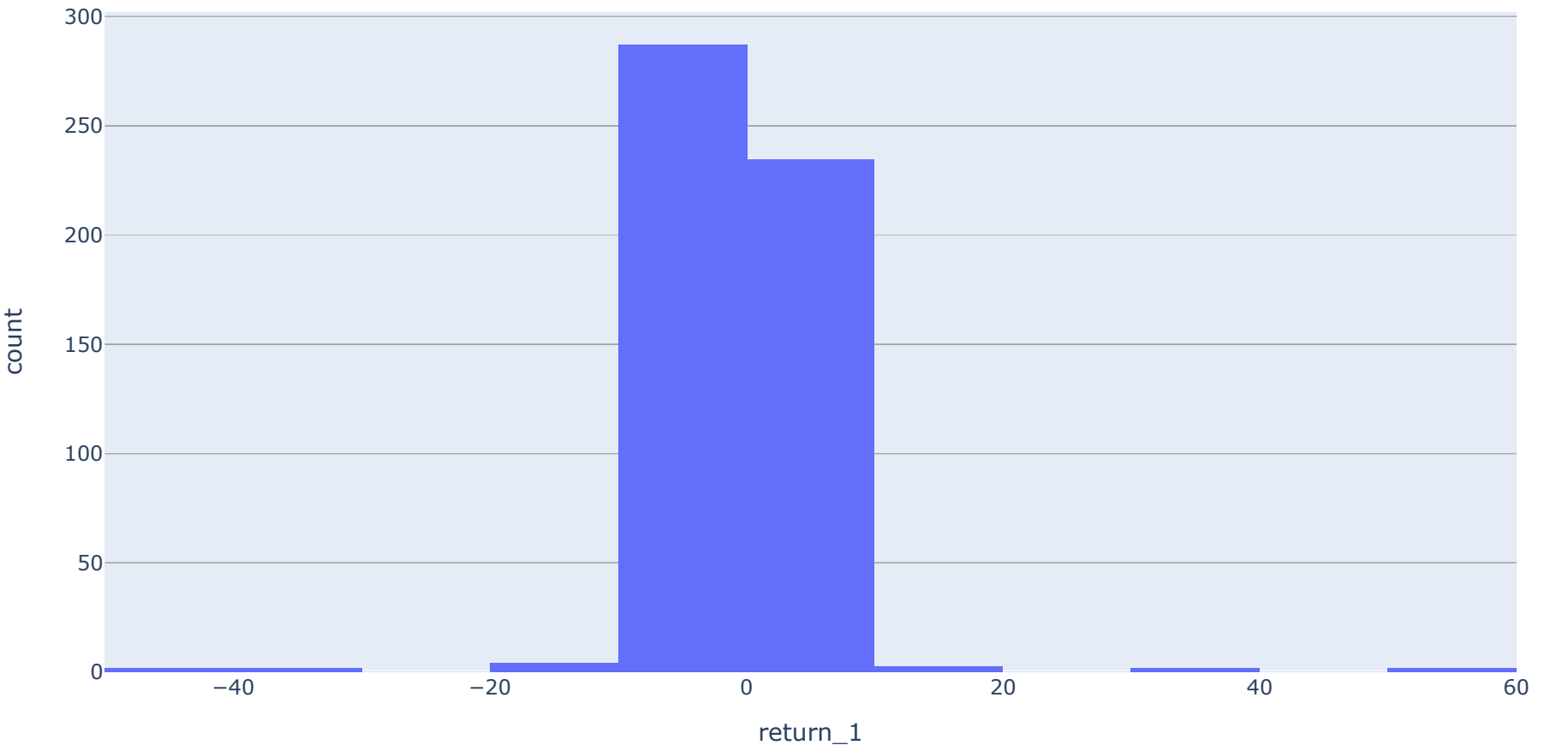
```
In [8]: 1 import plotly.graph_objects as go
2 from plotly.subplots import make_subplots
3
4 go.Figure(data=[
5     go.Scatter(x=group_b['date'], y=group_b['value'], line=dict(color='orange', width=1), name='Midwest Economy Inde
6     go.Scatter(x=group_b['date'], y=group_b['ma_6'], line=dict(color='blue', width=1), name='6 Period Moving Average
7 ]).show()
```



GROUP B - Histogram of 1 Period Return

reminder we are gonna have 10 bins based on quantile

```
In [9]: 1 import plotly.express as px
2 df = px.data.tips()
3 fig = px.histogram(group_b, x="return_1",nbins=20)
4 fig.show()
```



Section 1 - Classifier and Hyperparameters

1.1 Briefly present the maths of logistic classifier, use β for coefficients not θ , by writing down:

1.1.(a) the complete Loss Function (cost function) and relate it to MLE total log-likelihood function;

(Only briefly list down the math formula without a step by step induction)

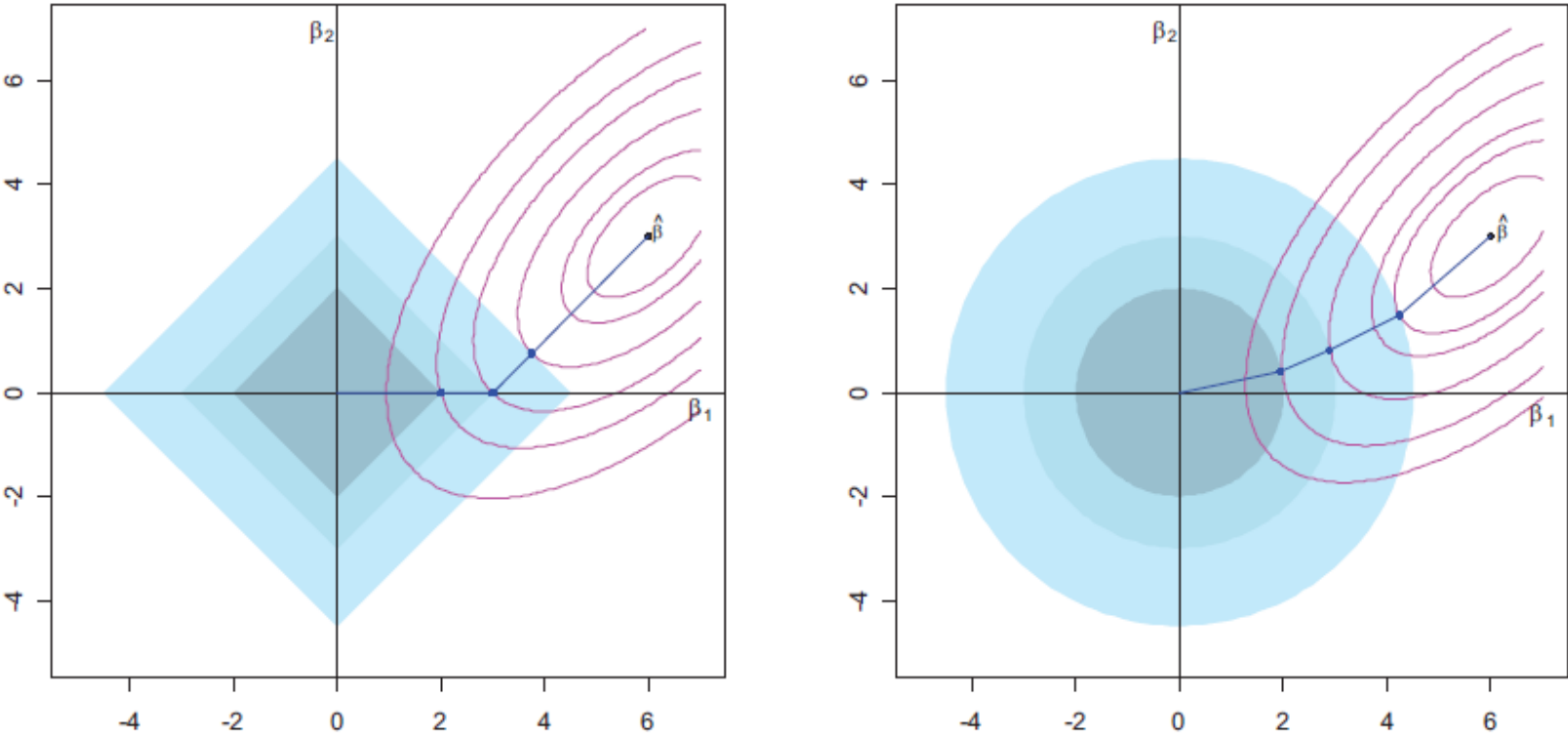
$$J(\beta) = \frac{1}{m} \sum_i^M y^i \log(h_{\beta}(x^i)) + (1 - y^i) \log(1 - h_{\beta}(x^i)), \text{ where } h_{\beta}(x) = \frac{1}{1+e^{-(\beta \cdot x)}} ,$$

so Log likeliness MLE $lik(\beta) = \prod f(x|\beta) = \sum_i^M y^i \log(h_{\beta}(x^i)) + (1 - y^i) \log(1 - h_{\beta}(x^i))$

1.1.(b) the penalty expressions for Lasso and Ridge regressions (mathematical and graphical);

$$J_{LASSO}(\beta) = J(\beta) + \lambda ||\beta|| = J(\beta) + \lambda \sum^k |\beta_k|$$

$$J_{RIDGE}(\beta) = J(\beta) + \lambda ||\beta||^2 = J(\beta) + \lambda \sum^k |\beta_k|^2$$



1.1.(c) RSS and MSE expressions, specific to logistic regression.

$$RSS = \sum_i^M (h_{\beta}(x^i) - y^i)^2$$

$$MSE = \frac{1}{n} \sum_i^N (h_{\beta}(x^i) - y^i)^2$$

In []:

1

1.2. Implement Ridge and Lasso logistic regressions using the suitable full set of features in each prediction task Group A, Group B

1.2.(a) produce a table comparing L1 and L2 type of penalisation: the impact made on regression coefficients;

Group A: HSI data

for group A data, setting daily return's sign(sign_1 column in group_a from the code below) (0 if return is negative and 1 otherwise) of daily return of dependent variable and the rest as independent variable, below is python code to do regression with L1 and L2 penlization and table comparing L1 and L2 penalization.

In conclusion, it looks like from the table and graph (each dot on the line represents a coefficient) below L2 boost L1's coefficients.

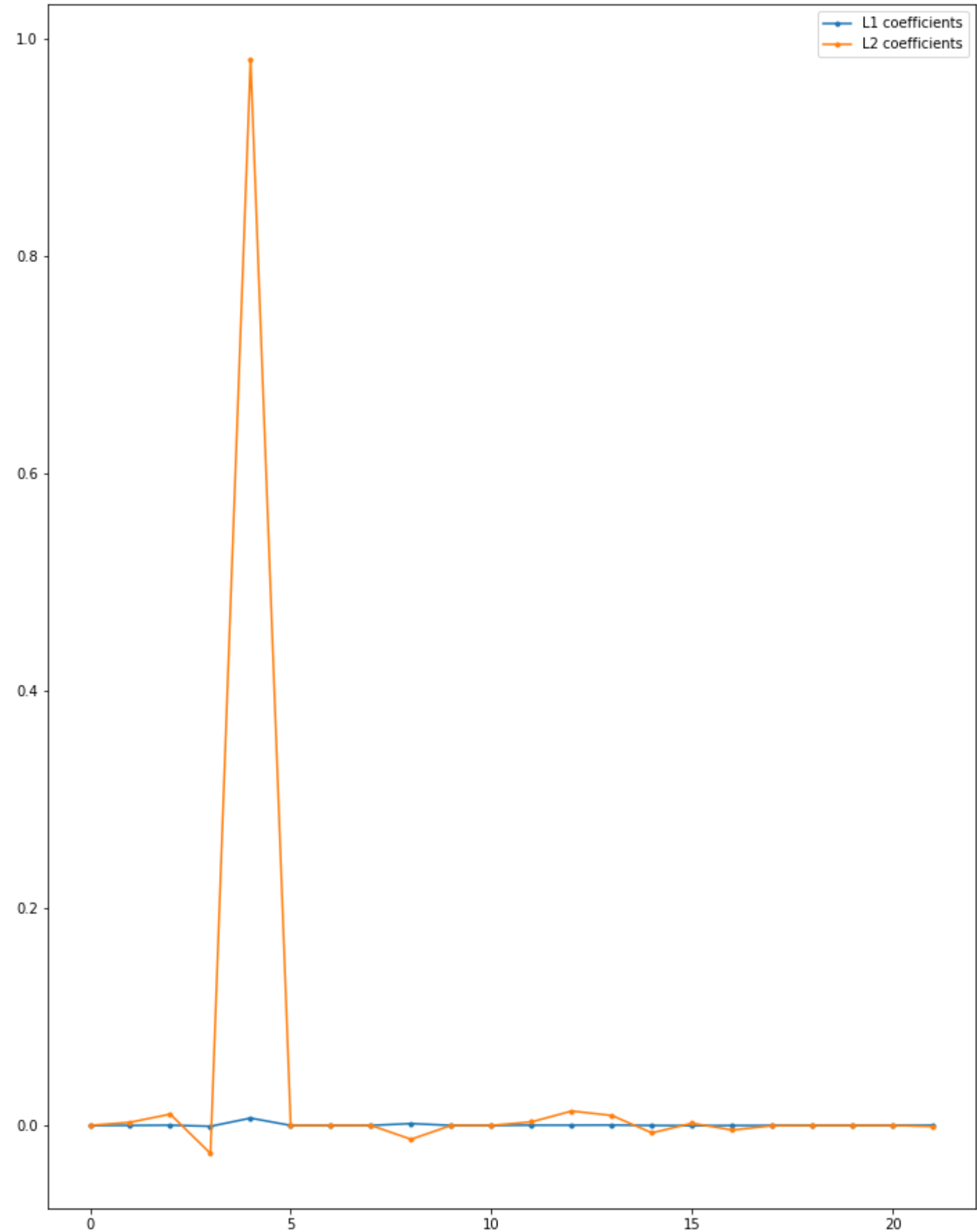
```
In [14]: 1 from sklearn.linear_model import LogisticRegression
2
3 y = group_a['sign_1'] # sign of daily return
4 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
5 assert 'sign_1' not in X.columns
6 # L1 model
7 l1_model = LogisticRegression(penalty='l1', solver='saga')
8 l1_model.fit(X, y)
9 l1_prediction = l1_model.predict(X) # predicting
10
11 # L2 model
12 l2_model = LogisticRegression(penalty='l2')
13 l2_model.fit(X, y)
14 l2_prediction = l2_model.predict(X) # predicting
15
16 l1_model_coefs = pd.DataFrame(l1_model.coef_, columns = X.columns)
17 l2_model_coefs = pd.DataFrame(l2_model.coef_, columns = X.columns)
18 table = pd.concat([l1_model_coefs, l2_model_coefs])
19 table.index=['L1', 'L2']
20 table
```

Out[14]:

	return_4	O-C	close	ewma	momentum_1	sign_5	return_1	sign_4	momentum_2	sign_3	...	open	momentum_3
L1	0.000000	0.000013	0.000222	-0.000826	0.006719	1.001688e-07	3.228137e-07	6.636860e-07	0.001691	0.000001	...	0.000235	0.000368
L2	0.000004	0.002875	0.010372	-0.025796	0.981807	-3.660943e-05	5.167622e-05	6.667411e-06	-0.012800	0.000052	...	0.013247	0.009137

2 rows × 22 columns

```
In [15]: 1 import matplotlib.pyplot as plt
2
3 fig = plt.figure(figsize=(12,16))
4 ax = plt.axes()
5 ax.plot(l1_model.coef_[0], marker ='.')
6 ax.plot(l2_model.coef_[0], marker ='.')
7 ax.legend(['L1 coefficients', 'L2 coefficients'])
8 fig.show()
```



Group B: Chicago Fed Relative Midwest Economy Index

before training a model, I'd first bucket the target variable with 10 bins based on quantile of the periodic(monthly) return of data. Then run L1 and L2 logistic regression to predict the bucketed periodic return, table below this cell shows a comparison between L1 and L2 coeffcients (index of table show if the row is L1 or L2), L1 and L2 model each has 10 rows, each of which shows coefficients to predict a specific bucket level of the periodic return.


```
In [16]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import KBinsDiscretizer
3
4 ##### bucketing dependent variable #####
5 y = group_b['return_1'] # sign of
6 X = group_b[set(group_b.columns) - set(['date', 'return_1'])]
7
8 est = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='quantile')
9 y_bucket = est.fit_transform([[i] for i in y.to_list()]) # arrange it as a 2-D array
10 y_bucket = [i[0] for i in y_bucket] # flats it back to 1D
11 ##### bucketing dependent variable #####
12
13 # L1 model
14 l1_model = LogisticRegression(multi_class='multinomial', solver='saga', penalty='l1')
15 l1_model.fit(X, y_bucket)
16 l1_model_coefs = pd.DataFrame(l1_model.coef_, columns = X.columns)
17 l1_prediction = l1_model.predict(X) # predicting
18 l1_model_coefs.index = ['L1']*len(l1_model_coefs)
19
20 # L2 model
21 l2_model = LogisticRegression(multi_class='multinomial', solver='lbfgs', penalty='l2')
22 l2_model.fit(X, y_bucket)
23 l2_model_coefs = pd.DataFrame(l2_model.coef_, columns = X.columns)
24 l2_model_coefs.index = ['L2']*len(l2_model_coefs)
25 l2_prediction = l2_model.predict(X) # predicting
26
27 table = pd.concat([l1_model_coefs, l2_model_coefs])
28 table
```

```
Out[16]:
```

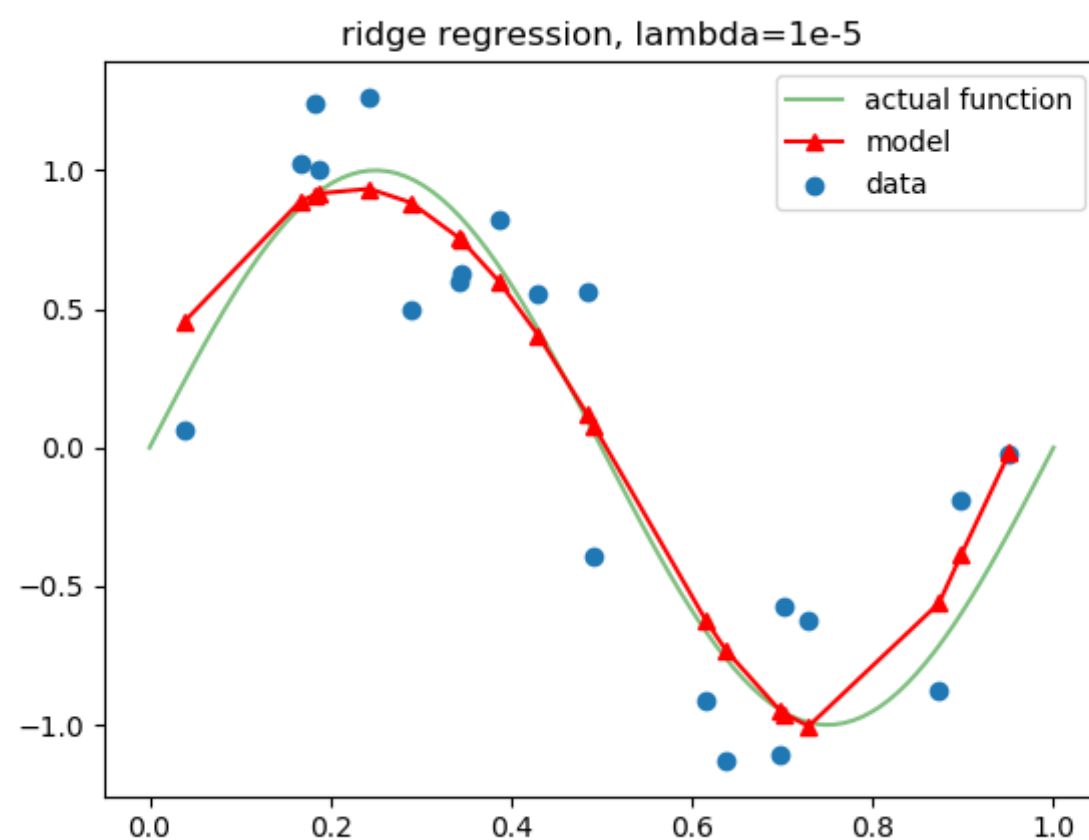
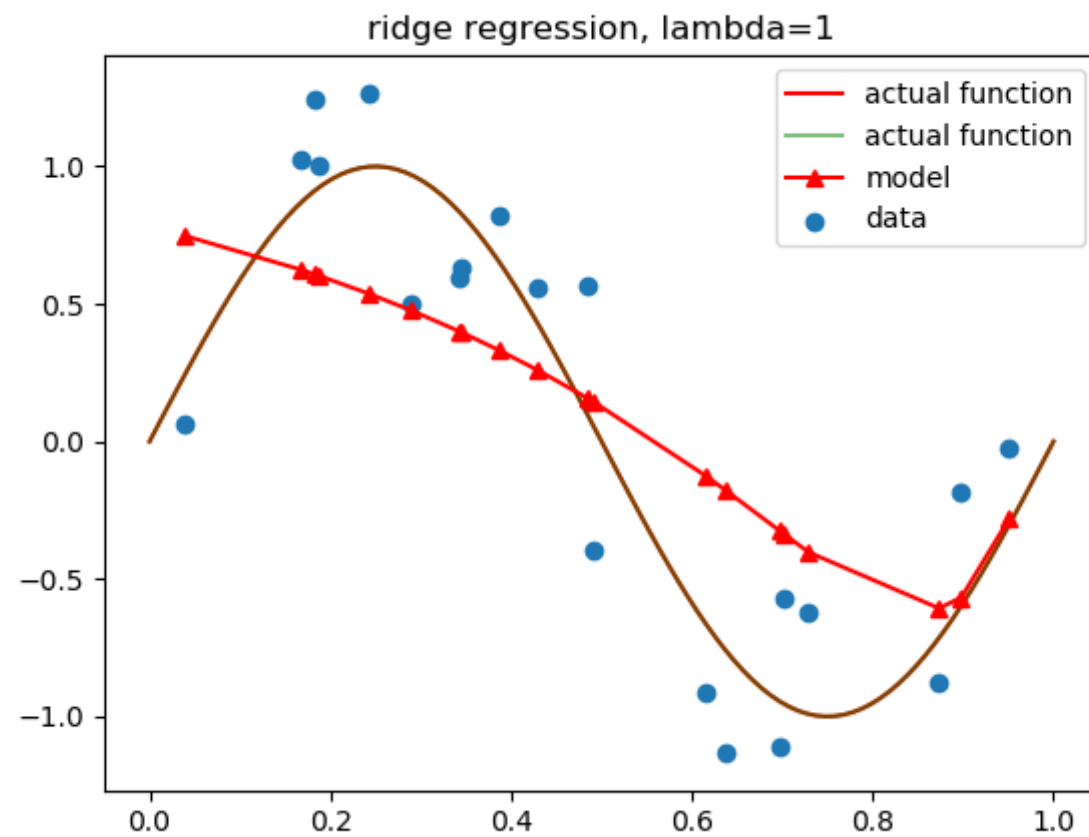
	ma_4	value	ma_2	return_4	sign_1	ma_5	momentum_1	sign_5	sign_4	momentum_2	...	momentum_3	momentum_4
L1	0.001745	-0.011476	-0.001104	-0.031650	-0.208707	0.002427	-0.015166	-0.052909	-0.023051	-0.029793	...	-0.035706	-0.0227
L1	-0.014424	0.003330	-0.011308	0.012162	-0.201096	-0.008873	0.045649	-0.102978	-0.081923	0.015801	...	0.037196	-0.0004
L1	-0.000710	-0.008362	-0.018033	0.031490	-0.209872	-0.003902	0.013985	-0.041588	-0.050190	-0.019554	...	-0.020068	0.0035
L1	0.007702	0.004366	0.013293	0.000384	-0.216683	0.004486	-0.012278	0.006485	-0.016479	-0.014940	...	0.019475	0.0072
L1	0.019868	0.012840	0.016119	-0.027327	-0.230787	0.017501	-0.001061	0.034434	0.003097	-0.012717	...	0.000000	0.0000
L1	0.000000	0.000381	0.000006	0.000381	-0.001013	0.000097	0.000000	0.062873	0.060278	0.000806	...	0.002873	-0.0157
L1	-0.020249	-0.015550	-0.017196	-0.010458	0.265452	-0.021952	0.000000	0.069701	0.056249	-0.000142	...	0.014842	0.0082
L1	0.004594	0.004509	0.003623	0.006815	0.254427	0.005948	0.000000	0.059651	0.084326	0.017475	...	-0.019797	-0.0011
L1	0.000000	0.008490	0.013625	0.036798	0.265885	-0.002221	-0.004977	-0.004659	0.018380	0.043379	...	0.014577	0.0435
L1	0.000000	-0.007892	0.000000	-0.023775	0.284843	0.000000	-0.017751	-0.034147	-0.050284	0.000574	...	-0.015026	-0.0290
L2	0.077699	-0.125609	0.013171	-0.039504	-2.072282	-0.023634	-0.277562	-0.158103	0.146634	-0.178627	...	-0.357046	0.3035
L2	0.032220	0.147961	-0.153822	0.043905	-1.536727	0.084102	0.603565	-0.536774	-0.130749	-0.248710	...	0.108109	-0.1436
L2	0.061495	0.008949	-0.218703	0.030546	-1.818749	0.067040	0.455305	0.242778	0.035671	-0.232516	...	-0.432973	-0.0802
L2	-0.056827	-0.045722	0.090191	-0.011895	-1.884748	0.031800	-0.271826	0.159532	-0.007216	-0.251420	...	0.567668	-0.4320
L2	0.067691	0.023612	-0.013568	-0.041750	-2.057006	0.015705	0.074360	-0.103710	-0.013803	-0.195486	...	-0.055191	0.2158
L2	-0.071845	0.028571	0.168442	-0.009101	0.017591	0.003061	-0.279742	0.043012	-0.019694	0.349052	...	0.332352	-0.2745
L2	-0.071472	-0.036979	-0.080523	-0.014068	2.339189	-0.077610	0.087088	0.316466	-0.052581	-0.466922	...	0.517808	0.0657
L2	0.061315	0.049942	0.005899	0.002455	2.213265	0.054591	0.088086	0.013855	0.028205	0.672703	...	-0.806283	0.0222
L2	0.016829	-0.003814	0.084988	0.056064	2.330196	-0.122895	-0.177604	-0.076858	-0.108366	0.079628	...	0.015404	0.6775
L2	-0.117105	-0.046910	0.103925	-0.016652	2.469269	-0.032160	-0.301671	0.099804	0.121899	0.472299	...	0.110153	-0.3545

20 rows × 24 columns

1.2.(b) clearly explain whether L1 or L2 regression likely to have a high bias and low variance.

looks like model from group a has an accuracy score of 99.8% _, and model from group b has an accuracy score of 10% (almost like randomly guess from a set of 10 targets, since it has 10 buckets), no matter what regularization factor I use, apparently both need to be tweaked (preprocessing, scaling, things like that) to show correctly the L1/L2's effect on bias/variance. So for now I'm going to borrow pictures from the internet to answer this question:

It is likely if we use an improper λ value of lasso/ridge regression. as the following images show (high λ versus low λ), apparently 1 is a overfit and the other being a underfit by tweaking λ



Section 2 - Model Selection

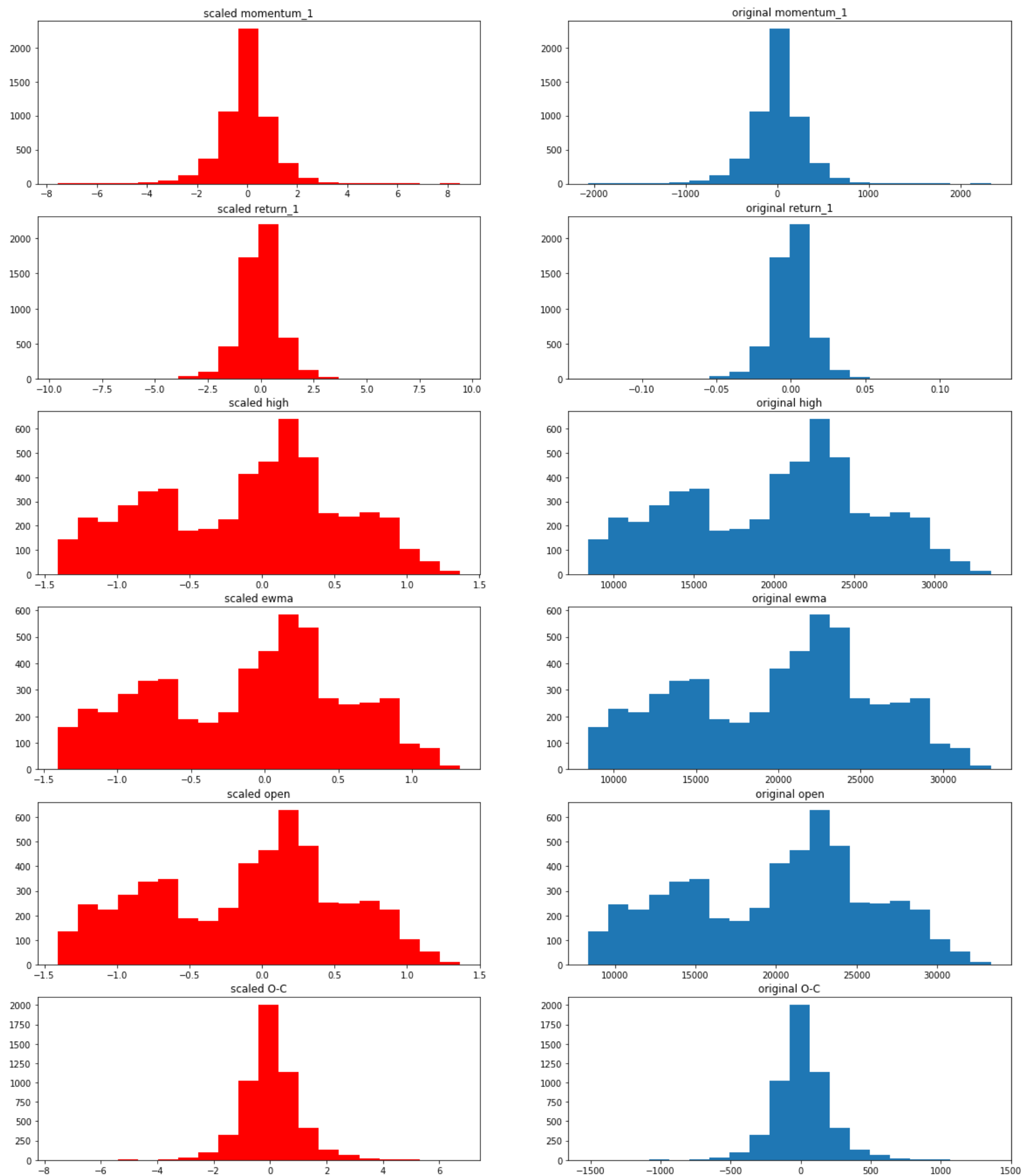
21. Return to preprocessing and fit one model using 2 types of scaling on features data. Produce histograms of changed data columns. Report on changes to coefficients.

2.1.(a) min-max scaling or robust quantiles-based scaling;

below is a trial of scaling on group A with RobustScaler, to best show the histograms, here I'm going to only plot a subset ('momentum_1', 'return_1', 'high', 'ewma', 'open', 'O-C') of all scaled columns rather than all (signed columns are not changed since they are bool). left column of histograms (in red) are scaled values, and right ones are original. One thing to note is shapes from left to right are not curved out (not changed), only the units are different.

In [17]:

```
1 import numpy as np
2 from sklearn import preprocessing
3
4 # X and y
5 y = group_a['sign_1'] # sign of daily return
6 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
7 # Create scaler
8 minmax_scale = preprocessing.RobustScaler()
9 # Scale or Transform the feature
10 scaled_X = minmax_scale.fit_transform(X)
11 scaled_X = pd.DataFrame(scaled_X, columns =X.columns)
12
13 ## subplots scaled X
14 plot_cols = ['momentum_1', 'return_1', 'high', 'ewma', 'open', 'O-C']
15 fig, ((ax1,ax2),(ax3,ax4),(ax5,ax6),(ax7,ax8),(ax9,ax10),(ax11,ax12)) = plt.subplots(6,2, figsize=(20,24))
16 fig.suptitle('Sharing x per column, y per row')
17
18 ax1.hist(scaled_X['momentum_1'], bins=20, color="red")
19 ax1.set_title("scaled momentum_1")
20 ax2.hist(X['momentum_1'], bins=20)
21 ax2.set_title("original momentum_1")
22
23 ax3.hist(scaled_X['return_1'], bins=20, color="red")
24 ax3.set_title("scaled return_1")
25 ax4.hist(X['return_1'], bins=20)
26 ax4.set_title("original return_1")
27
28 ax5.hist(scaled_X['high'], bins=20, color="red")
29 ax5.set_title("scaled high")
30 ax6.hist(X['high'], bins=20)
31 ax6.set_title("original high")
32
33 ax7.hist(scaled_X['ewma'], bins=20, color="red")
34 ax7.set_title("scaled ewma")
35 ax8.hist(X['ewma'], bins=20)
36 ax8.set_title("original ewma")
37
38 ax9.hist(scaled_X['open'], bins=20, color="red")
39 ax9.set_title("scaled open")
40 ax10.hist(X['open'], bins=20)
41 ax10.set_title("original open")
42
43 ax11.hist(scaled_X['O-C'], bins=20, color="red")
44 ax11.set_title("scaled O-C")
45 ax12.hist(X['O-C'], bins=20)
46 ax12.set_title("original O-C")
47
48 fig.show()
```



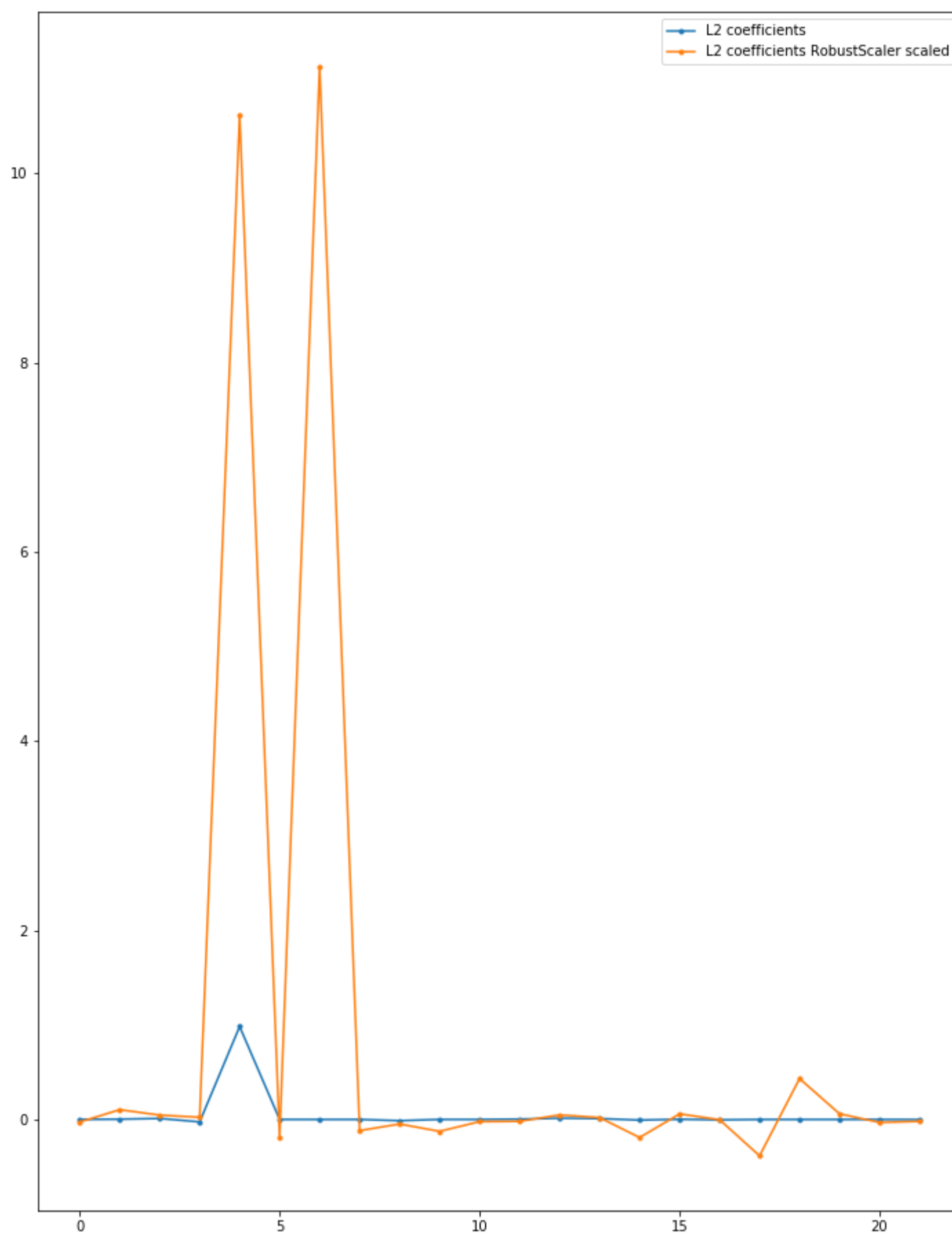
changes to coefficients of robust scaler:

below we gonna repeat L2 logistic regression modelling from section 1.2.1, and looks like from the plotting(each dot on both lines represents a coefficient value) below scaling the feature columns would make the coefficients boosted,

```

In [18]: 1 from sklearn.linear_model import LogisticRegression
2
3 y = group_a['sign_1'] # sign of daily return
4 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
5 assert 'sign_1' not in X.columns
6 # L2 model
7 l2_model = LogisticRegression(penalty='l2')
8 l2_model.fit(X, y)
9 l2_prediction = l2_model.predict(X) # predicting
10
11 # L2 model
12 l2_model_scaled = LogisticRegression(penalty='l2')
13 l2_model_scaled.fit(scaled_X, y)
14 l2_prediction_scaled = l2_model_scaled.predict(scaled_X) # predicting
15
16 fig = plt.figure(figsize=(12,16))
17 ax = plt.axes()
18 ax.plot(l2_model.coef_[0], marker = '.')
19 ax.plot(l2_model_scaled.coef_[0], marker = '.')
20 ax.legend(['L2 coefficients', 'L2 coefficients RobustScaler scaled'])
21 fig.show()

```

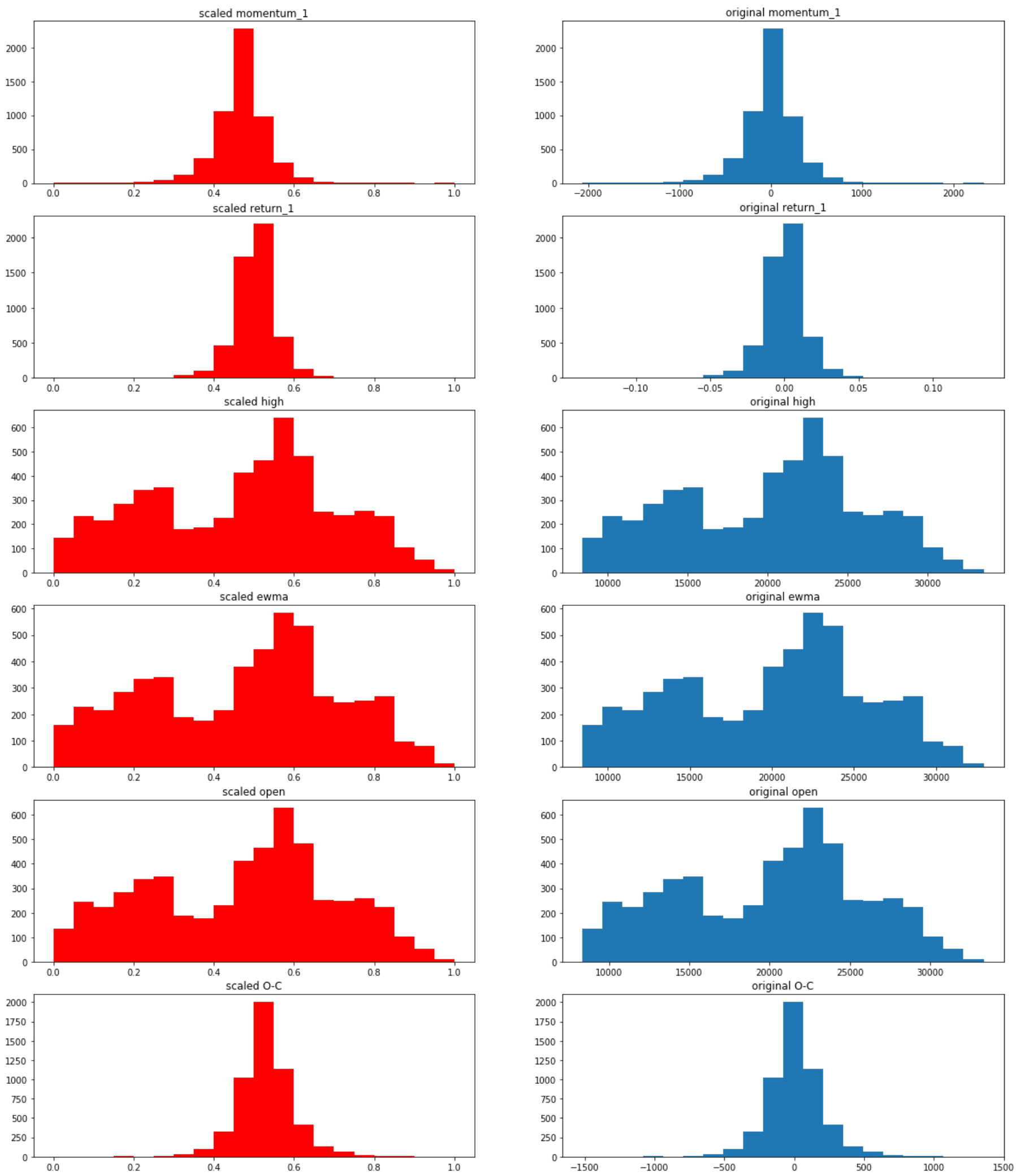


2.1.(b) scaling to uniform [0, 1] range:

below is a trial of scaling on group A with `MinMaxScaler` , to best show the histograms, here I'm going to only plot a subset ('momentum_1', 'return_1', 'high', 'ewma', 'open', 'O-C') of all scaled columns rather than all (signed columns are not changed since they are bool). left column of histograms (in red) are scaled values, and right ones are original. One thing to note is all scaled columns are now ranging from 0 - 1. Also, coefficients are boosted as well compared to L2 model with original unscaled features data

In [19]:

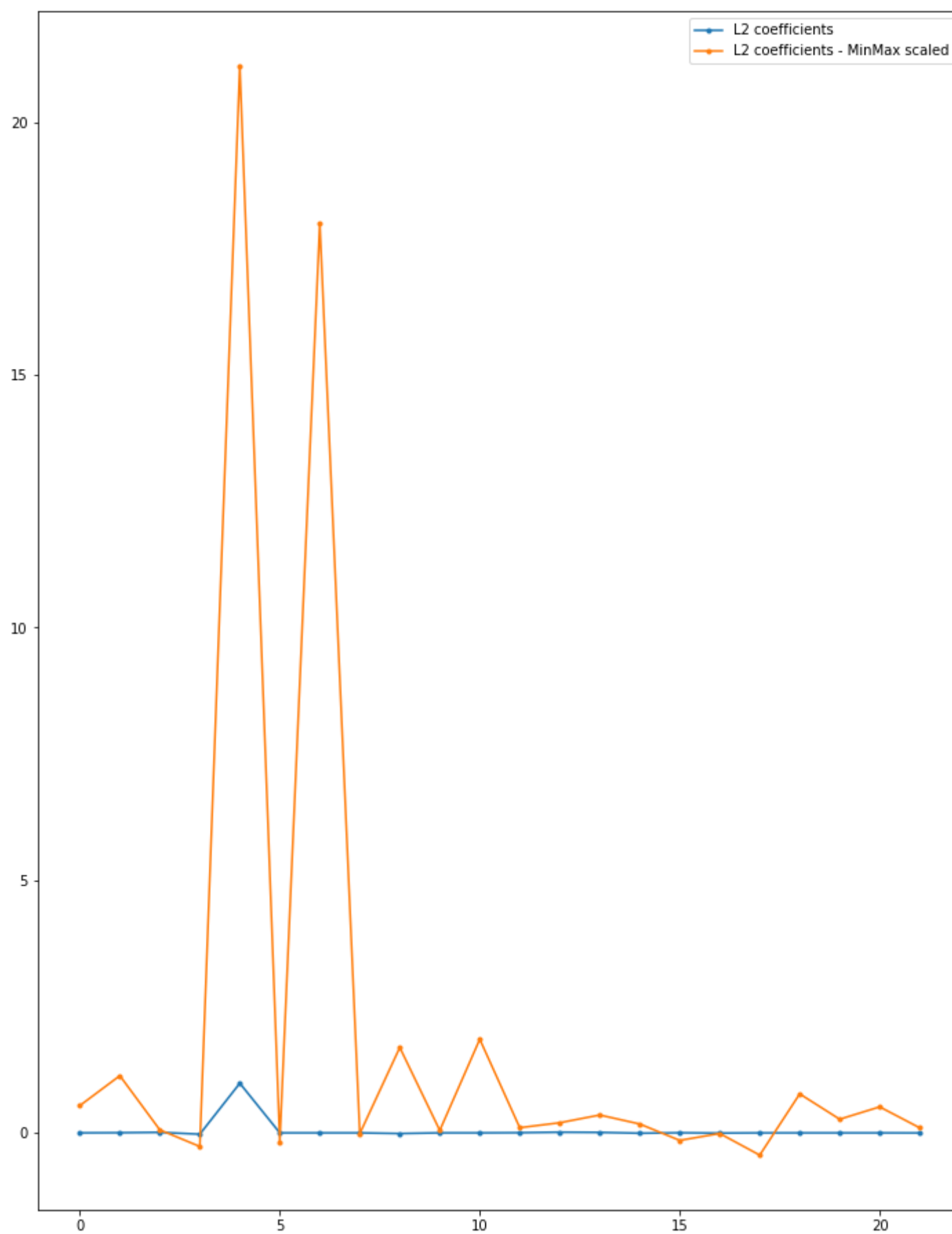
```
1 import numpy as np
2 from sklearn import preprocessing
3
4 # X and y
5 y = group_a['sign_1'] # sign of daily return
6 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
7 # Create scaler
8 minmax_scale = preprocessing.MinMaxScaler()
9 # Scale or Transform the feature
10 scaled_X = minmax_scale.fit_transform(X)
11 scaled_X = pd.DataFrame(scaled_X, columns =X.columns)
12
13 ## subplots scaled X
14 plot_cols = ['momentum_1', 'return_1', 'high', 'ewma', 'open', 'O-C']
15 fig, ((ax1,ax2),(ax3,ax4),(ax5,ax6),(ax7,ax8),(ax9,ax10),(ax11,ax12)) = plt.subplots(6,2, figsize=(20,24))
16 fig.suptitle('Sharing x per column, y per row')
17
18 ax1.hist(scaled_X['momentum_1'], bins=20, color="red")
19 ax1.set_title("scaled momentum_1")
20 ax2.hist(X['momentum_1'], bins=20)
21 ax2.set_title("original momentum_1")
22
23 ax3.hist(scaled_X['return_1'], bins=20, color="red")
24 ax3.set_title("scaled return_1")
25 ax4.hist(X['return_1'], bins=20)
26 ax4.set_title("original return_1")
27
28 ax5.hist(scaled_X['high'], bins=20, color="red")
29 ax5.set_title("scaled high")
30 ax6.hist(X['high'], bins=20)
31 ax6.set_title("original high")
32
33 ax7.hist(scaled_X['ewma'], bins=20, color="red")
34 ax7.set_title("scaled ewma")
35 ax8.hist(X['ewma'], bins=20)
36 ax8.set_title("original ewma")
37
38 ax9.hist(scaled_X['open'], bins=20, color="red")
39 ax9.set_title("scaled open")
40 ax10.hist(X['open'], bins=20)
41 ax10.set_title("original open")
42
43 ax11.hist(scaled_X['O-C'], bins=20, color="red")
44 ax11.set_title("scaled O-C")
45 ax12.hist(X['O-C'], bins=20)
46 ax12.set_title("original O-C")
47
48 fig.show()
```




```

In [20]: 1 from sklearn.linear_model import LogisticRegression
2
3 y = group_a['sign_1'] # sign of daily return
4 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
5 assert 'sign_1' not in X.columns
6 # L2 model
7 l2_model = LogisticRegression(penalty='l2')
8 l2_model.fit(X, y)
9 l2_prediction = l2_model.predict(X) # predicting
10
11 # L2 model
12 l2_model_scaled = LogisticRegression(penalty='l2')
13 l2_model_scaled.fit(scaled_X, y)
14 l2_prediction_scaled = l2_model_scaled.predict(scaled_X) # predicting
15
16 fig = plt.figure(figsize=(12,16))
17 ax = plt.axes()
18 ax.plot(l2_model.coef_[0], marker = '.')
19 ax.plot(l2_model_scaled.coef_[0], marker = '.')
20 ax.legend(['L2 coefficients', 'L2 coefficients - MinMax scaled'])
21 fig.show()

```



2.2. Choose a restricted set of features, use penalisation and other (below), and produce a model that makes predictions with reduced variance (MSE). Accuracy score is not estimator variance

2.2(a) formally explain your approach to feature selection or elimination. For example, VIF for elimination of interdependent (colinear) features, or selection according to a score from the F-test.

Gonna use group b data (St Louis Fed Relative Midwest Economy Index) and VIF as an approach to select best limited set of features, and only keep features with VIF value between 5 and 0. The features being kept are ['return_2', 'return_5', 'sign_2', 'sign_3', 'return_3', 'return_6', 'sign_4', 'sign_6', 'sign_1', 'sign_5', 'return_4'] which are sign of T-1 thru T-6 returns and sign of them

```
In [21]: 1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 from sklearn.preprocessing import StandardScaler
3
4 def vif(X):
5     # perform feature scaling
6     scaler = StandardScaler()
7     xs = scaler.fit_transform(X)
8     # subsume into a dataframe
9     vif = pd.DataFrame()
10    vif["Features"] = X.columns
11    vif["VIF Factor"] = [variance_inflation_factor(xs, i) for i in range(xs.shape[1])]
12    return vif
13
14 y = group_b['return_1'] # sign of
15 X = group_b[set(group_b.columns) - set(['date', 'return_1'])]
16
17 result = vif(X)
18 select_features = result[(result['VIF Factor'] <= 5) & (result['VIF Factor'] > 0)]['Features'].tolist()
19 print(select_features)
```

```
['return_4', 'sign_1', 'sign_5', 'sign_4', 'return_6', 'sign_3', 'return_2', 'sign_2', 'return_3', 'sign_6', 'return_5']
```

below is an experiment using group b data, L2 model(C=1) with dependent variable (periodic return) bucketed, 10 bins based on quantile, between full set of features and a subset of features output by VIF metrics. we use 400/532 observations for training, 132/532 as out of sample test. As the result shows, MSE metrics actually decreases a bit if we were to use a subset of features by $(3.45 - 2.96)/3.45 = 14\%$.

```
In [22]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import KBinsDiscretizer
3 from sklearn.metrics import mean_squared_error
4
5 # create training and test dataset with full set of features
6 y = group_b['return_1'] # sign of
7 X = group_b[set(group_b.columns) - set(['date', 'return_1'])]
8 ##### bucketing dependent variable #####
9 est = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='quantile')
10 y_bucket = est.fit_transform([i for i in y.to_list()]) # arrange it as a 2-D array
11 y_bucket = [i[0] for i in y_bucket] # flats it back to 1D
12
13 y_train, y_test = y_bucket[:400], y_bucket[400:]
14 X_train_all_features, X_test_features = X[:400], X[400:]
15
16 # create training and test dataset with subset of features from VIF experiment
17 X_subset = group_b[select_features]
18 X_train_subset_features, X_test_subset_features = X_subset[:400], X_subset[400:]
19
20 ##### bucketing dependent variable #####
21 def calc_mse(X_train, y_train, X_test, y_test):
22     # L2 model
23     l2_model = LogisticRegression(multi_class='multinomial', solver='lbfgs', penalty='l2', C=1)
24     #l2_model = LogisticRegression(penalty='l2')
25     l2_model.fit(X_train, y_train)
26     l2_prediction = l2_model.predict(X_test) # predicting
27     return mean_squared_error(y_test, l2_prediction)
28
29 print(f"MSE when using full set of features: {calc_mse(X_train_all_features, y_train, X_test_features, y_test)}")
30 print(f"MSE when using features selected by VIF: {calc_mse(X_train_subset_features, y_train, X_test_subset_features,
```

```
MSE when using full set of features: 3.4545454545454546
```

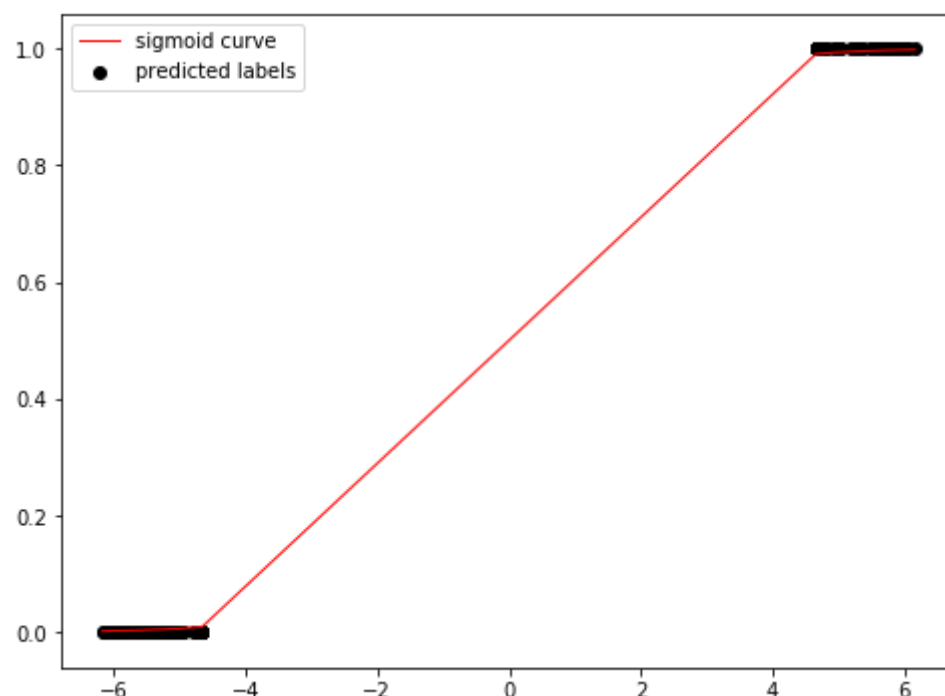
```
MSE when using features selected by VIF: 2.962121212121212
```

2.2.(b) if prediction task allows (Group A, restricted model) graphically plot a logistic sigmoid for each feature. If prediction is not binomial (Group B), sigmoid only possible between two classes.

not so understand what's plot a logistic sigmoid for each feature mean, but here I'm gonna use group a data to train a L2 logistic regression model with features coming from VIF selection, and plot a sigmoid function of subset of features with scatter plotting the predicted labels.

```
In [23]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import train_test_split
3 from scipy.special import expit
4
5 ### select subset of features by VIF ###
6 y = group_a['sign_1'] # sign of daily return
7 X = group_a[set(group_a.columns) - set(['date', 'return_1'])]
8 result = vif(X)
9 select_features = result[(result['VIF Factor'] <= 5) & (result['VIF Factor'] > 0)][['Features']].tolist()
10 X = group_a[select_features]
11
12 # Scale or Transform the feature
13 minmax_scale = preprocessing.MinMaxScaler()
14 scaled_X = minmax_scale.fit_transform(X)
15 scaled_X = pd.DataFrame(scaled_X, columns=X.columns)
16
17 X_train, X_test, y_train, y_test = train_test_split(scaled_X, y, test_size=0.05, random_state=42)
18
19 # L2 model
20 l2_model = LogisticRegression(penalty='l2')
21 l2_model.fit(X_train, y_train)
22 l2_prediction = l2_model.predict(X_test) # predicting
23
24 # plot sigmod curve
25 plt.figure(1, figsize=(8, 6))
26 plt.clf()
27 # sort data
28 X_test['row_sum'] = ((X_test * l2_model.coef_).sum(axis=1) + l2_model.intercept_)
29 X_test.sort_values('row_sum', inplace=True)
30 row_sum = X_test['row_sum']
31 X_test.drop('row_sum', axis=1, inplace=True)
32
33 plt.scatter(row_sum, l2_model.predict(X_test), color='black', zorder=1)
34
35 # sigmoid curve somehow it looks like a straightline below --
36 loss = expit(row_sum).ravel()
37
38 plt.plot(row_sum, loss, color='red', linewidth=1)
39 plt.legend(["sigmoid curve", 'predicted labels'])
```

Out[23]: <matplotlib.legend.Legend at 0x216e8940508>



2.3. For the winning restricted model produce a suitable set of Evaluation metrics, such as area under ROC curve (each class) and confusion matrix. Give expressions for precision/recall.

Using group a data for this question as a demo, below is code to do a grid search with cross validation to find the best model to predict sign of T-1 daily return of HSI data. And cells following it would be some evaluation metrics to gauge the best model.

The best model is:

best params: {'logistic__C': 0.01, 'logistic__penalty': 'l2'}

best score: 0.7835389627765361

```
In [31]: 1 from sklearn.pipeline import Pipeline
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.metrics import balanced_accuracy_score, mean_squared_error, roc_curve
4 from sklearn import metrics
5
6 ### select subset of features by VIF ###
7 y = group_a['sign_1'] # sign of daily return
8 est = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='quantile')
9 #y_bucket = est.fit_transform([[i] for i in y.to_list()]) # arrange it as a 2-D array
10 #y_bucket = [i[0] for i in y_bucket] # flats it back to 1D
11
12 X = group_a[set(group_a.columns) - set(['date', 'sign_1'])]
13 result = vif(X)
14 select_features = result[(result['VIF Factor'] <= 5) & (result['VIF Factor'] > 0)][['Features']].tolist()
15 X = group_a[select_features]
16
17 pipe = Pipeline([("scaler", StandardScaler()), ("logistic", LogisticRegression(solver='liblinear'))])
18 pipe.fit(X, y)
19
20 # set up a
21 penalty = ['l1', 'l2']
22 C = np.linspace(0.01, 10, 10)
23 param_grid = dict(logistic__C=C, logistic__penalty=penalty)
24
25 scoring = {'AUC': 'roc_auc', 'Accuracy': 'accuracy'}
26 gridsearch = GridSearchCV(pipe, param_grid, n_jobs=1, cv=5, verbose=1,
27                           scoring=scoring, refit='AUC')
28 # fit grid search
29 best_model = gridsearch.fit(X, y)
30 print("best params: ", best_model.best_params_)
31 print("best score: ", best_model.best_score_)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

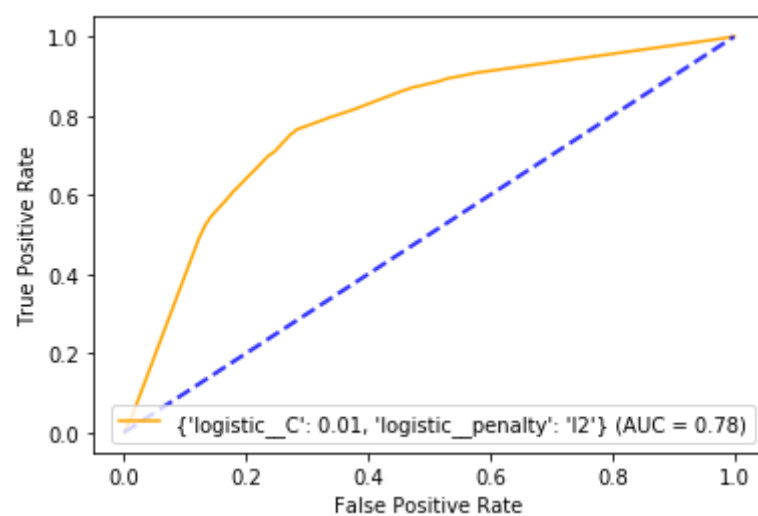
best params: {'logistic__C': 0.01, 'logistic__penalty': 'l2'}
best score: 0.7835389627765361

[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 1.3s finished

Evaluation metrics - ROC

Below is ROC Curve metrics (in sample), the orange curve is above the benchmark line (blue dotted), and area under curve (AUC) is 0.78, so the model is not too bad.

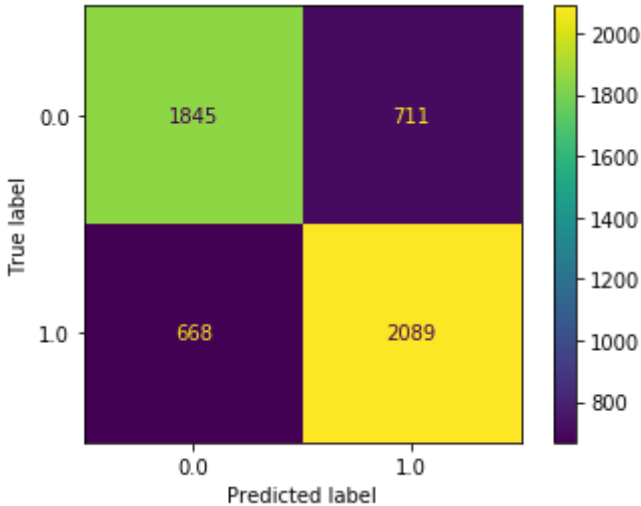
```
In [34]: 1 import matplotlib.pyplot as plt
2 from sklearn import datasets, metrics, model_selection, svm
3
4 metrics.plot_roc_curve(pipe, X, y, color='orange', name=best_model.best_params_)
5 plt.plot([0, 1], [0, 1], linestyle='--', lw=2, color='b',
6          label='Chance', alpha=.8)
7
8 plt.show()
```



Evaluation metrics - confusion matrix

Below is confusion matrix metrics (in sample), next cell would be a confusion matrix heat map:

```
In [46]: 1 from sklearn.metrics import plot_confusion_matrix
2 from sklearn.metrics import confusion_matrix
3
4 result = confusion_matrix(y, best_model.predict(X))
5 metrics.plot_confusion_matrix(pipe, X, y)
6 plt.show()
7
8 true_negative, false_positive, false_negative ,true_positive = result[0][0], result[0][1], result[1][0], result[1][1]
9 print(f"True Negative: {true_negative}")
10 print(f"False Positive: {false_positive}")
11 print(f"False Negative: {false_negative}")
12 print(f"True Positive: {true_positive}")
```



True Negative: 1845
False Positive: 711
False Negative: 668
True Positive: 2089

and the expression for precision/recall would be :

$$precision = \frac{truepositive}{truepositive+falsepositive} = 2089/(2089 + 711) = 0.746$$
$$recall = \frac{truepositive}{truepositive+falsenegative} = 2089/(2089 + 668) = 0.7577$$
$$accuracy = \frac{truepositive+truenegative}{total} = (2089 + 1845)/(2089 + 1845 + 711 + 668) = 0.74$$

Section 3 - Mathematical bases of supervised learning

$$MSE(\hat{\beta}) = E[(\hat{\beta} - \beta)^2] = Var[\hat{\beta}] + (E[\hat{\beta}] - \beta)^2$$

3.1.(a) can there exist an estimator with the smaller MSE than minimal least squares?

Yes and No, "No" as a first impression of this question, of course not! since it gives the minial squared error according to its definition given model type and number of training data are fixed. BUT "YES" at risk of being over meticulous, estimator of a Lass/ridge regression has an MLS larger than MLS of ordinary linear regression. For example,

Lasso's cost function $\beta = (X^t X + \mu * I) - 1 X^t Y$

ordinary linear regression's cost function $\beta = (X^t X) - 1 X^t Y$

so of course Lasso's cost function is minimal when $\mu = 0$, which is ordinary Linear regression's cost function

3.1(b) for a prediction, does the MSE measure an irreducible error or model error?

Yes. MSE could be described as the addition of model variance, model bias, and irreducible uncertainty, which comes from noise in our data set.

```
In [ ]: 1
```