

# Stock Prediction

With Tweet sentiment and financial data

Group: Fight Potatoes



Wenxin Gu | Siwen Chen | Carolyn Liu



# Project Objective

**Output:**

**Healthcare Sector Stock Price Prediction**

**Input:**

**Financial Data**

**Generate financial ratios from quarter report**

- Profitability
- Liquidity
- Growth
- Credibility

**Generate stock trend indicators using stock price**

- RSI
- MACD

**Tweets Data**

**Generate sentiment indicators using Covid-related Tweets**

- Polarity
- Subjectivity

**Overview**

**Phase: 1**

**Phase 2**

**Phase 3**

**Summary**

# Phase 1: Data Collection

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes	
Collect Data	Independent Variables	Dependent variables	<div>24 sheets (per stock)</div> <div>*</div> <div>785 rows (2020.3.1 – 2022.4.24)</div> <div>*</div> <div>21 features</div>	
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li></ul>	<ul style="list-style-type: none"><li>Dummy stock price: -1/1</li><li>Stock price</li></ul>		
Data processing	Transfer all data into daily data	Fill weekend data		
	<ul style="list-style-type: none"><li>Transfer the quarterly data by filling the value by the last valid observation</li></ul>	<ul style="list-style-type: none"><li>Stock Market only has data on weekdays</li><li>Fill by taking average</li></ul>		
Feature selection	Split into training and testing	Information Gain Method		
	<ul style="list-style-type: none"><li>Use 2020.3.1 – 2021.12.31 as training</li><li>Use 2022.1.1 – 2022.4.24 as testing</li></ul>	<ul style="list-style-type: none"><li>Select the feature that has an importance larger than 0 using training data</li></ul>		
Modeling	Classifier	SGD Regression		
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li></ul>	<ul style="list-style-type: none"><li>Tried SGD regression to predict the stock price instead of classifying</li></ul>		

	G	H	I	J	K	L	M	N	O	P	Q	
/ (C Current Ri Debt Ratic Dividends Dummy_ C Gross Prof Interest C Inventory Log Reven Net Acqui Net Profit Quick Rati	737	1.533173	0.254655	0.179577	1	0.539137	16.66667	2.013255	9.135769	0.122725	0.141917	1.111538
338	1.639535	0.259131	-0.60215	1	0.532793	12.64706	1.922096	9.13258			0.145173	1.158562
211	1.630334	0.263141	0.196429	1	0.530695	6	1.650667	9.092721			0.081583	1.136247
388	2.469559	0.243348	0.210526	1	0.530531	12.77778	1.690349	9.100715			0.157811	1.741248
263	2.32788	0.24504	0.157143	1	0.531355	16.61111	2.059722	9.171141			0.149697	1.689843
474	2.064612	0.258321	0.299492	1	0.541344	17.26316	2.050331	9.189771			0.186047	1.432128
397	1.998862	0.281977	0.133787	1	0.535738	15.15789	1.927939	9.18327	0.052606		0.141639	1.396473
357	2.106729	0.272424	0.210714	1	0.537201	16	1.938875	9.200303	-9.53E-05		0.166456	1.479118
519	2.224239		0.155263	1	0.542169	19.75	2	9.220108	0		0.266265	1.608314
334	2.193182		0.35	1	0.543608	18.8	1.904437	9.223755			0.169056	1.491793
236	3.17735	0.748785	0.512532	1	0.769531	9	4.800883	9.939719	0.012736		0.321806	2.909978
514	3.142371	0.735162	0.477778	1	0.774684	8.418224	4.674078	9.935457	0.000132		0.349228	2.884099
326	0.862452	0.584712	0.601648	1	0.644029	2.614007	2.568367	10.01808	0.256256	-0.07079		0.58403
577	0.951532	0.581864	0.374131	1	0.608588	5.322581	3.713874	10.11066	0.005882	-0.178887		0.676812
109	0.843411	0.571554	0.459235	1	0.662	6.558252	4.186707	10.1417	0.002657	0.002598		0.603643
376	0.831555	0.568302	0.495202	1	0.676172	6.596463	3.976161	10.11428	0.001316	0.273098		0.606084
336	0.906394	0.555092	0.492013	1	0.67598	7.60396	4.122563	10.14485	0.000993	0.054875		0.645447
222	1.01486	0.542268	0.299948	1	0.693906	8.02735	4.635423	10.15661	0.003305	0.221657		0.754565
403	0.70344	0.533333	0.473404	1	0.70344	0.508443	4.750851	10.17378	0.007368	0.271655		0.652705

Overview

Phase: 1

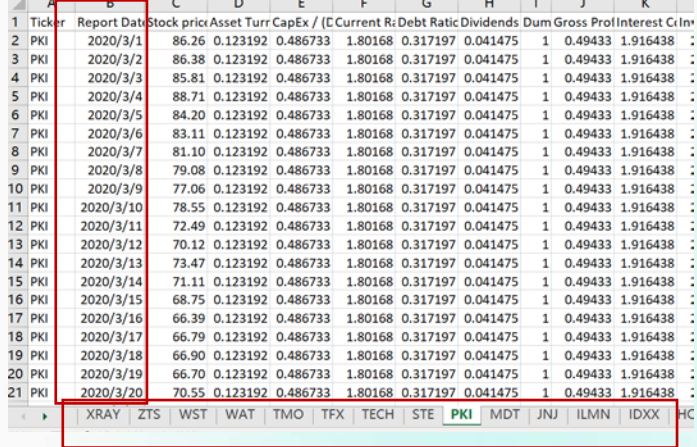
Phase 2

Phase 3

Summary

# Phase 1: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes
Collect Data	Independent Variables	Dependent variables	<b>In total:</b> <b>24 stocks (sheets) *</b> <b>785 rows * 21 features</b> 
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li></ul>	<ul style="list-style-type: none"><li>Dummy stock price: -1/1</li><li>Stock price</li></ul>	
Data processing	Transfer all data into daily data	Fill weekend data	
	<ul style="list-style-type: none"><li>Transfer the quarterly data by filling the value by the last valid observation</li></ul>	<ul style="list-style-type: none"><li>Stock Market only has data on weekdays</li><li>Fill by taking average</li></ul>	
Feature selection	Split into training and testing	Information Gain Method	
	<ul style="list-style-type: none"><li>Use 2020.3.1 – 2021.12.31 as training</li><li>Use 2022.1.1 – 2022.4.24 as testing</li></ul>	<ul style="list-style-type: none"><li>Select the feature that has an importance larger than 0 using training data</li></ul>	
Modeling	Classifier	SGD Regression	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li></ul>	<ul style="list-style-type: none"><li>Tried SGD regression to predict the stock price instead of classifying</li></ul>	

Overview

Phase: 1

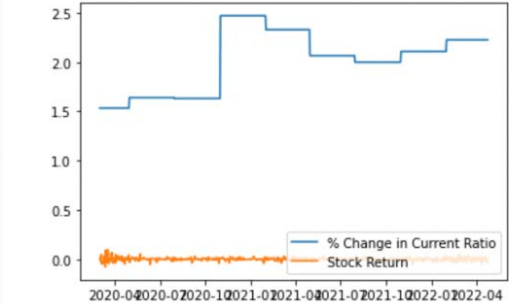
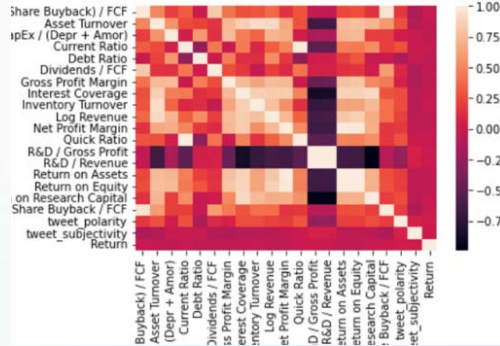
Phase 2

Phase 3

Summary

# Phase 1: Data Processing (Continue)

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes
Collect Data	<b>Independent Variables</b> <ul style="list-style-type: none"> <li>Quarterly data: Financial ratios</li> <li>Daily data: Tweets sentiment indicators</li> </ul>	<b>Dependent variables</b> <ul style="list-style-type: none"> <li>Dummy stock price: -1/1</li> <li>Stock price</li> </ul>	<p><b>Irrelevant features</b></p>  <p><b>Redundant features</b></p> 
Data processing	<b>Transfer all data into daily data</b> <ul style="list-style-type: none"> <li>Transfer the quarterly data by filling the value by the last valid observation</li> </ul>	<b>Fill weekend data</b> <ul style="list-style-type: none"> <li>Stock Market only has data on weekdays</li> <li>Fill by taking average</li> </ul>	
Feature selection	<b>Split into training and testing</b> <ul style="list-style-type: none"> <li>Use 2020.3.1 – 2021.12.31 as training</li> <li>Use 2022.1.1 – 2022.4.24 as testing</li> </ul>	<b>Information Gain Method</b> <ul style="list-style-type: none"> <li>Select the feature that has an importance larger than 0 using training data</li> </ul>	
Modeling	<b>Classifier</b> <ul style="list-style-type: none"> <li>Ran logistic regression for each of the 24 stocks</li> </ul>	<b>SGD Regression</b> <ul style="list-style-type: none"> <li>Tried SGD regression to predict the stock price instead of classifying</li> </ul>	

Overview

Phase: 1

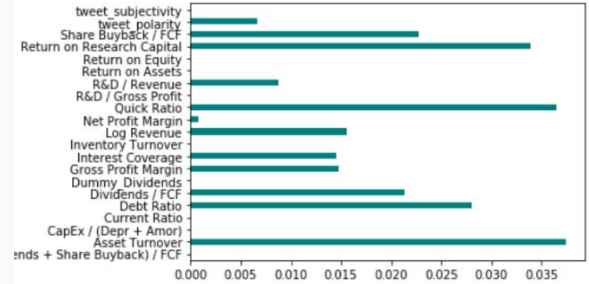
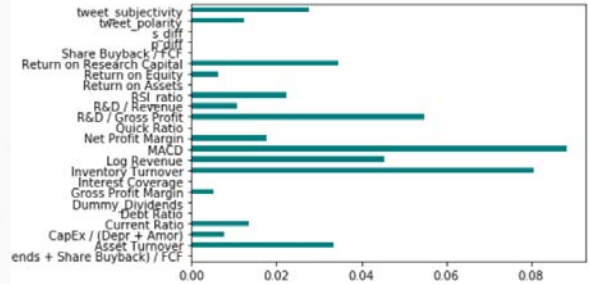
Phase 2

Phase 3

Summary

# Phase 1: Feature Selection

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes	
Collect Data	Independent Variables	Dependent variables	<b>Information Gain Method showing the importance of different features</b> Example: Stock 1 	
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li></ul>	<ul style="list-style-type: none"><li>Dummy stock price: -1/1</li><li>Stock price</li></ul>		
Data processing	Transfer all data into daily data	Fill weekend data	<b>Example: Stock 2</b> 	
	<ul style="list-style-type: none"><li>Transfer the quarterly data by filling the value by the last valid observation</li></ul>	<ul style="list-style-type: none"><li>Stock Market only has data on weekdays</li><li>Fill by taking average</li></ul>		
Feature selection	Split into training and testing	Information Gain Method		
	<ul style="list-style-type: none"><li>Use 2020.3.1 – 2021.12.31 as training</li><li>Use 2022.1.1 – 2022.4.24 as testing</li></ul>	<ul style="list-style-type: none"><li>Select the feature that has an importance larger than 0 using training data</li></ul>		
Modeling	Classifier	SGD Regression		
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li></ul>	<ul style="list-style-type: none"><li>Tried SGD regression to predict the stock price instead of classifying</li></ul>		

Overview

Phase: 1

Phase 2

Phase 3

Summary

# Phase 1: Modeling

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes
Collect Data	Independent Variables	Dependent variables	Outcome for the classifier model <pre>logit_accuracy_results  {'XRAY': 0.43859649122807015, 'ZTS': 0.6842105263157895, 'WST': 0.42105263157894735, 'WAT': 0.3333333333333333, 'TMO': 0.4473684210526316, 'TFX': 0.5701754385964912, 'TECH': 0.38596491228070173, 'STE': 0.43859649122807015, 'PKI': 0.49122807017543857, 'MDT': 0.47368421052631576, 'JNJ': 0.5350877192982456, 'ILMN': 0.5877192982456141, 'IDXX': 0.35964912280701755, 'HOLX': 0.45614035087719296, 'DXCM': 0.5263157894736842, 'A': 0.37719298245614036, 'ABBV': 0.6140350877192983, 'ABT': 0.4473684210526316, 'AMGN': 0.4473684210526316, 'BAX': 0.4649122807017544, 'BSX': 0.4649122807017544, 'BIO': 0.6666666666666666, 'COO': 0.5175438596491229, 'DHR': 0.3508771929824561}</pre>
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li></ul>	<ul style="list-style-type: none"><li>Dummy stock price: -1/1</li><li>Stock price</li></ul>	
Data processing	Transfer all data into daily data	Fill weekend data	
	<ul style="list-style-type: none"><li>Transfer the quarterly data by filling the value by the last valid observation</li></ul>	<ul style="list-style-type: none"><li>Stock Market only has data on weekdays</li><li>Fill by taking average</li></ul>	
Feature selection	Split into training and testing	Information Gain Method	
	<ul style="list-style-type: none"><li>Use 2020.3.1 – 2021.12.31 as training</li><li>Use 2022.1.1 – 2022.4.24 as testing</li></ul>	<ul style="list-style-type: none"><li>Select the feature that has an importance larger than 0 using training data</li></ul>	
Modeling	Classifier	SGD Regression	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li></ul>	<ul style="list-style-type: none"><li>Tried SGD regression to predict the stock price instead of classifying</li></ul>	

Overview

Phase: 1

Phase 2

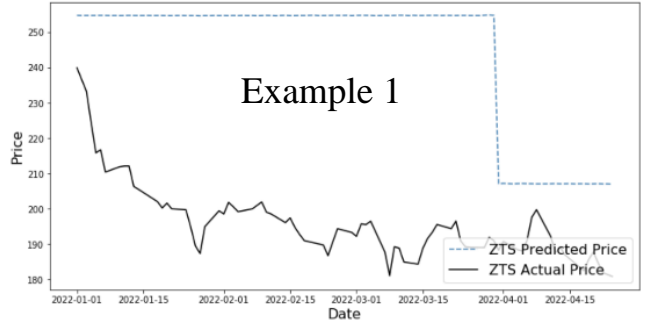
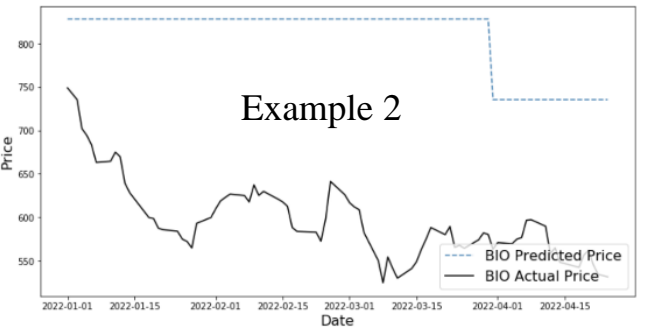
Phase 3

Summary



# Phase 1: Modeling (Continue)

Goal: Increase the Stock Price Prediction Accuracy

Phase One			Outcomes	
Collect Data	Independent Variables	Dependent variables	<div>Outcome for the SGD regression</div> <div>Example 1</div>  <div>Example 2</div> 	
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li></ul>	<ul style="list-style-type: none"><li>Dummy stock price: -1/1</li><li>Stock price</li></ul>		
Data processing	Transfer all data into daily data	Fill weekend data		
	<ul style="list-style-type: none"><li>Transfer the quarterly data by filling the value by the last valid observation</li></ul>	<ul style="list-style-type: none"><li>Stock Market only has data on weekdays</li><li>Fill by taking average</li></ul>		
Feature selection	Split into training and testing	Information Gain Method		
	<ul style="list-style-type: none"><li>Use 2020.3.1 – 2021.12.31 as training</li><li>Use 2022.1.1 – 2022.4.24 as testing</li></ul>	<ul style="list-style-type: none"><li>Select the feature that has an importance larger than 0 using training data</li></ul>		
Modeling	Classifier	SGD Regression		
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li></ul>	<ul style="list-style-type: none"><li>Tried SGD regression to predict the stock price instead of classifying</li></ul>		

Overview

Phase: 1

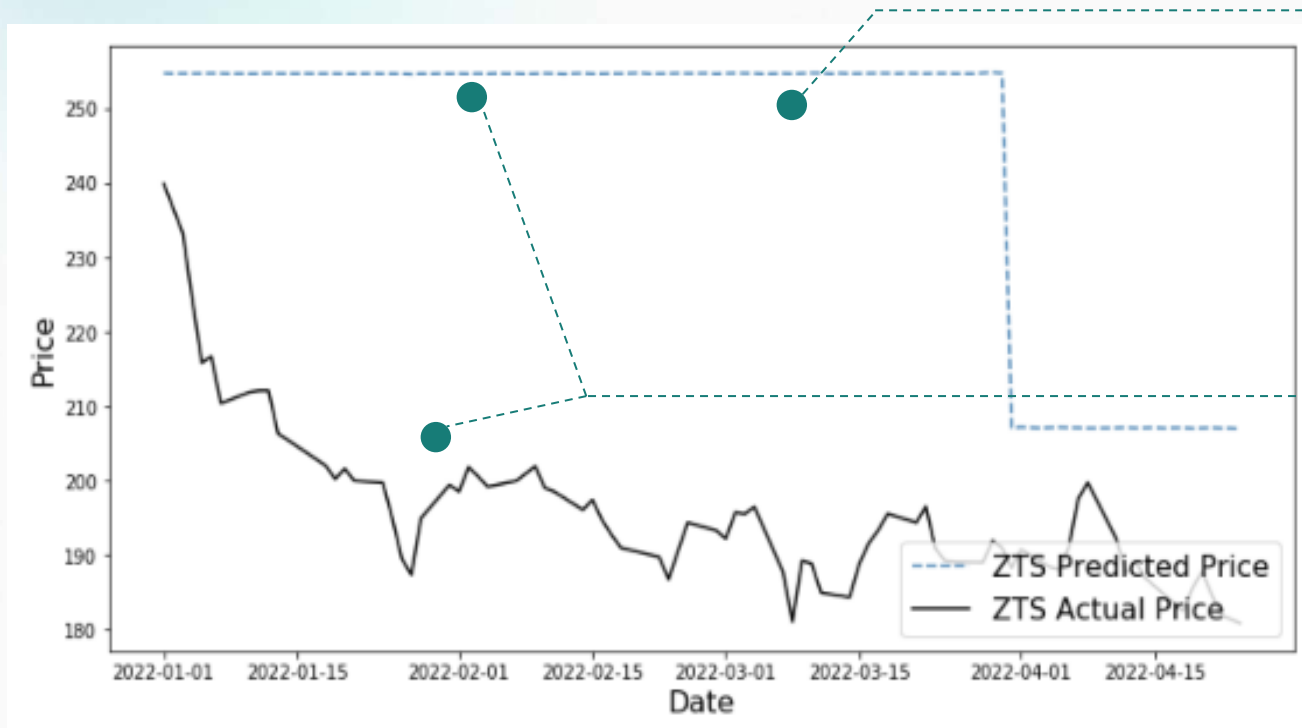
Phase 2

Phase 3

Summary



# Phase 1: Problems



## Fluctuation too small?

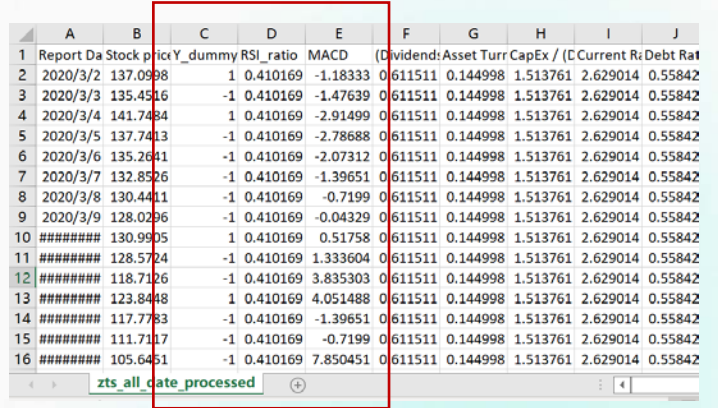
- Potential Reason:  
tweets data not so informative?
- Possible Solution:
  - Normalization?
  - Add more daily data?

## The predicted value is too high comparing to the actual value

- Potential Reason:  
Using the past two years data is too old for the current situation?
- Possible Solution:  
Splitting the data into smaller windows?

# Phase 2: Data collection

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes
Collect Data	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li><li>Dummy stock price: 1, 0</li></ul>	<ul style="list-style-type: none"><li>Stock trend indicator: RSI, MACD</li><li>Stock price: 1, 0, -1</li></ul>	
Data processing	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Transfer all data into daily data</li><li>Fill weekend and missing data</li></ul>	<ul style="list-style-type: none"><li>Normalization</li><li>Change in variable</li></ul>	
Feature selection	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Split into training and testing</li><li>Select feature using Information Gain Method</li></ul>	<ul style="list-style-type: none"><li>Add Fisher's score to select feature together with information gain method</li></ul>	
Modeling	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li><li>Tried SGD regression</li></ul>	<ul style="list-style-type: none"><li>Focus on 2 stocks instead of 24 stocks</li><li>Divide into 5 smaller windows</li></ul>	

Overview

Phase: 1

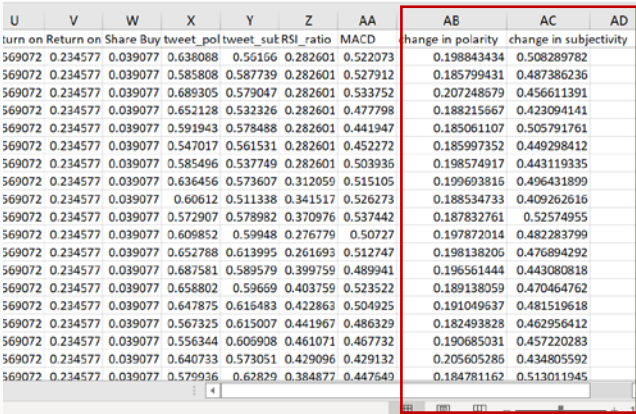
Phase 2

Phase 3

Summary

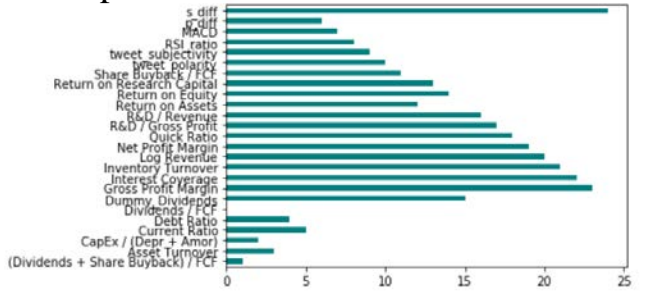
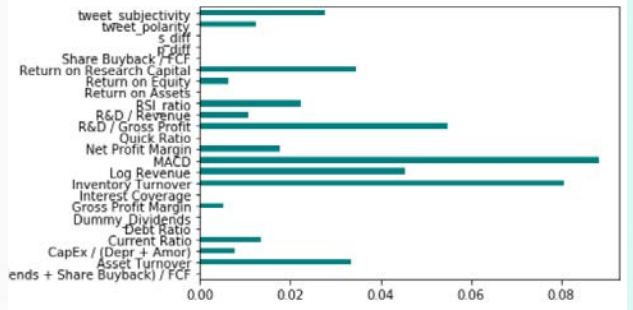
# Phase 2: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes			
Collect Data	Phase 1	Based on Phase 1, We add:				
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li><li>Dummy stock price: 1, 0</li></ul>	<ul style="list-style-type: none"><li>Stock trend indicator: RSI, MACD</li><li>Stock price: 1, 0, -1</li></ul>				
Data processing	Phase 1	Based on Phase 1, We add:				
	<ul style="list-style-type: none"><li>Transfer all data into daily data</li><li>Fill weekend and missing data</li></ul>	<ul style="list-style-type: none"><li>Normalization</li><li>Change in variable</li></ul>				
Feature selection	Phase 1	Based on Phase 1, We add:				
	<ul style="list-style-type: none"><li>Split into training and testing</li><li>Select feature using Information Gain Method</li></ul>	<ul style="list-style-type: none"><li>Add Fisher's score to select feature together with information gain method</li></ul>				
Modeling	Phase 1	Based on Phase 1, We add:				
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li><li>Tried SGD regression</li></ul>	<ul style="list-style-type: none"><li>Focus on 2 stocks instead of 24 stocks</li><li>Divide into 5 smaller windows</li></ul>				

# Phase 2: Feature Selection

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes
Collect Data	Phase 1	Based on Phase 1, We add:	<p>Use both the method to select features</p> <p>Example: Fisher's score</p> 
	<ul style="list-style-type: none"> <li>Quarterly data: Financial ratios</li> <li>Daily data: Tweets sentiment indicators</li> <li>Dummy stock price: 1, 0</li> </ul>	<ul style="list-style-type: none"> <li>Stock trend indicator: RSI, MACD</li> <li>Stock price: 1, 0, -1</li> </ul>	
Data processing	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"> <li>Transfer all data into daily data</li> <li>Fill weekend and missing data</li> </ul>	<ul style="list-style-type: none"> <li>Normalization</li> <li>Change in variable</li> </ul>	
Feature selection	Phase 1	Based on Phase 1, We add:	<p>Example: Information gain</p> 
	<ul style="list-style-type: none"> <li>Split into training and testing</li> <li>Select feature using Information Gain Method</li> </ul>	<ul style="list-style-type: none"> <li>Add Fisher's score to select feature together with information gain method</li> </ul>	
Modeling	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"> <li>Ran logistic regression for each of the 24 stocks</li> <li>Tried SGD regression</li> </ul>	<ul style="list-style-type: none"> <li>Focus on 2 stocks instead of 24 stocks</li> <li>Divide into 5 smaller windows</li> </ul>	

Overview

Phase: 1







Phase 2

Phase 3

Summary

# Phase 2: Modeling

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes
Collect Data	Phase 1	Based on Phase 1, We add:	<b>5 windows:</b>  bio_test_group0_2021-01-01_2021-06-30  bio_test_group1_2021-04-01_2021-09-30  bio_test_group2_2021-07-01_2021-12-31  bio_test_group3_2021-10-01_2022-04-24  bio_test_group4_2022-01-01_2022-04-24  bio_train_group0_2020-03-02_2020-12-31  bio_train_group1_2020-06-01_2021-03-31  bio_train_group2_2020-09-01_2021-06-30  bio_train_group3_2021-01-01_2021-09-30  bio_train_group4_2021-03-01_2021-12-31  zts_test_group0_2021-01-01_2021-06-30  zts_test_group1_2021-04-01_2021-09-30  zts_test_group2_2021-07-01_2021-12-31  zts_test_group3_2021-10-01_2022-04-24  zts_test_group4_2022-01-01_2022-04-24  zts_train_group0_2020-03-02_2020-12-31  zts_train_group1_2020-06-01_2021-03-31  zts_train_group2_2020-09-01_2021-06-30  zts_train_group3_2021-01-01_2021-09-30  zts_train_group4_2021-03-01_2021-12-31
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li><li>Dummy stock price: 1, 0</li></ul>	<ul style="list-style-type: none"><li>Stock trend indicator: RSI, MACD</li><li>Stock price: 1, 0, -1</li></ul>	
Data processing	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Transfer all data into daily data</li><li>Fill weekend and missing data</li></ul>	<ul style="list-style-type: none"><li>Normalization</li><li>Change in variable</li></ul>	
Feature selection	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Split into training and testing</li><li>Select feature using Information Gain Method</li></ul>	<ul style="list-style-type: none"><li>Add Fisher's score to select feature together with information gain method</li></ul>	
Modeling	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li><li>Tried SGD regression</li></ul>	<ul style="list-style-type: none"><li>Focus on 2 stocks instead of 24 stocks</li><li>Divide into 5 smaller windows</li><li>Other models: Decision Tree, SVM, etc.</li></ul>	

Overview

Phase: 1

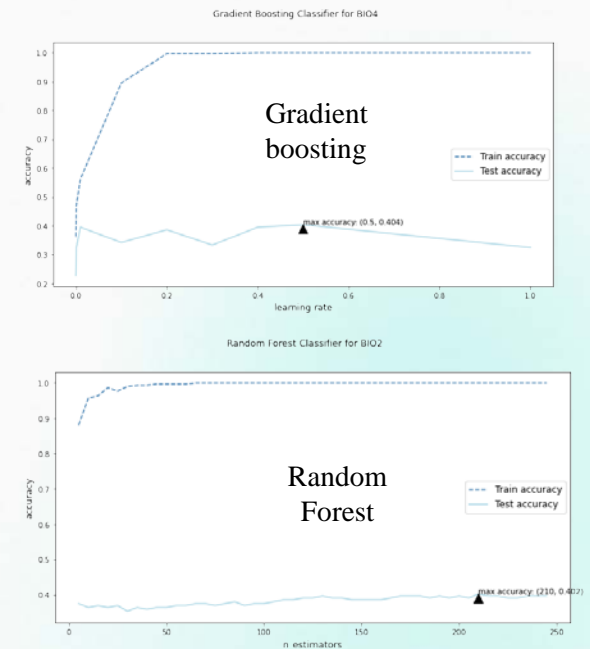
Phase 2

Phase 3

Summary

## Phase 2: Modeling (Continue)

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes
Collect Data	Phase 1	Based on Phase 1, We add:	<b>A glimpse of different models: accuracy slightly increase</b> 
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li><li>Dummy stock price: 1, 0</li></ul>	<ul style="list-style-type: none"><li>Stock trend indicator: RSI, MACD</li><li>Stock price: 1, 0, -1</li></ul>	
Data processing	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Transfer all data into daily data</li><li>Fill weekend and missing data</li></ul>	<ul style="list-style-type: none"><li>Normalization</li><li>Change in variable</li></ul>	
Feature selection	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Split into training and testing</li><li>Select feature using Information Gain Method</li></ul>	<ul style="list-style-type: none"><li>Add Fisher's score to select feature together with information gain method</li></ul>	
Modeling	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li><li>Tried SGD regression</li></ul>	<ul style="list-style-type: none"><li>Focus on 2 stocks instead of 24 stocks</li><li>Divide into 5 smaller windows</li><li>6 other models: Decision Tree, SVM, etc.</li></ul>	

Overview

Phase: 1

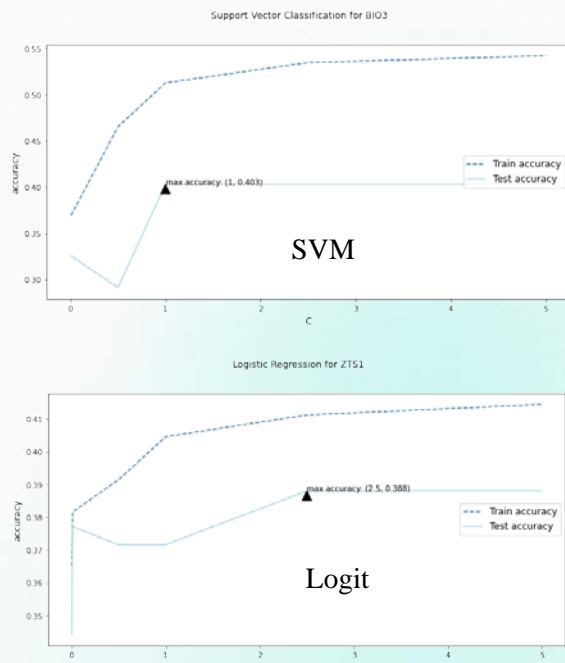
Phase 2

Phase 3

Summary

# Phase 2: Modeling (Continue)

Goal: Increase the Stock Price Prediction Accuracy

Phase Two			Outcomes
Collect Data	Phase 1	Based on Phase 1, We add:	<b>A glimpse of different models: accuracy slightly increase</b> 
	<ul style="list-style-type: none"><li>Quarterly data: Financial ratios</li><li>Daily data: Tweets sentiment indicators</li><li>Dummy stock price: 1, 0</li></ul>	<ul style="list-style-type: none"><li>Stock trend indicator: RSI, MACD</li><li>Stock price: 1, 0, -1</li></ul>	
Data processing	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Transfer all data into daily data</li><li>Fill weekend and missing data</li></ul>	<ul style="list-style-type: none"><li>Normalization</li><li>Change in variable</li></ul>	
Feature selection	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Split into training and testing</li><li>Select feature using Information Gain Method</li></ul>	<ul style="list-style-type: none"><li>Add Fisher's score to select feature together with information gain method</li></ul>	
Modeling	Phase 1	Based on Phase 1, We add:	
	<ul style="list-style-type: none"><li>Ran logistic regression for each of the 24 stocks</li><li>Tried SGD regression</li></ul>	<ul style="list-style-type: none"><li>Focus on 2 stocks instead of 24 stocks</li><li>Divide into 5 smaller windows</li><li>6 other models: Decision Tree, SVM, etc.</li></ul>	

Overview

Phase: 1

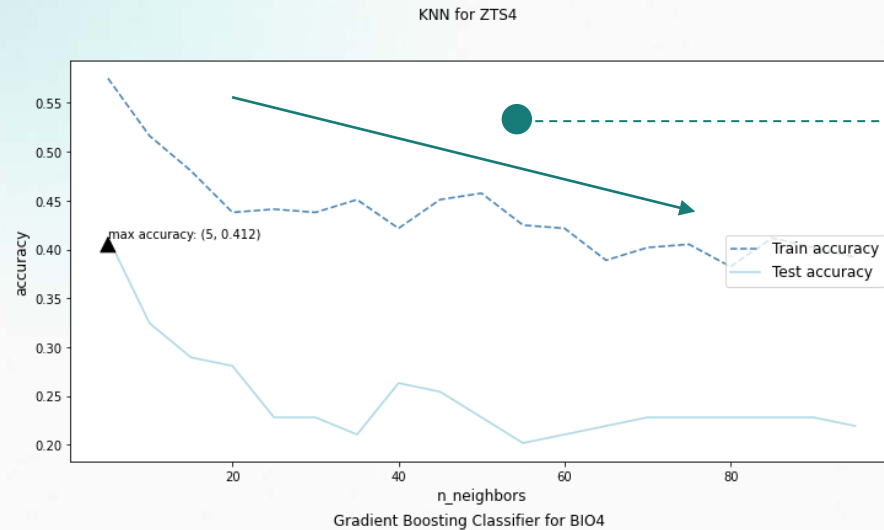
Phase 2

Phase 3

Summary

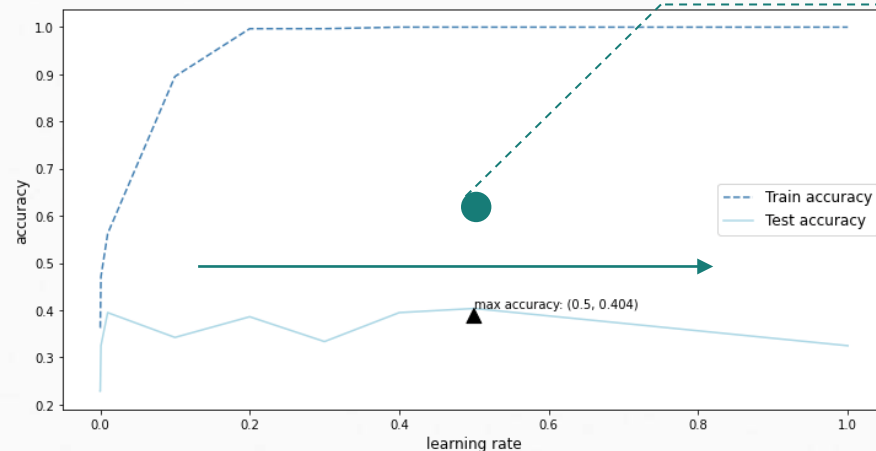


## Phase 2: some interesting findings



### KNN method accuracy decrease as $k$ increasing

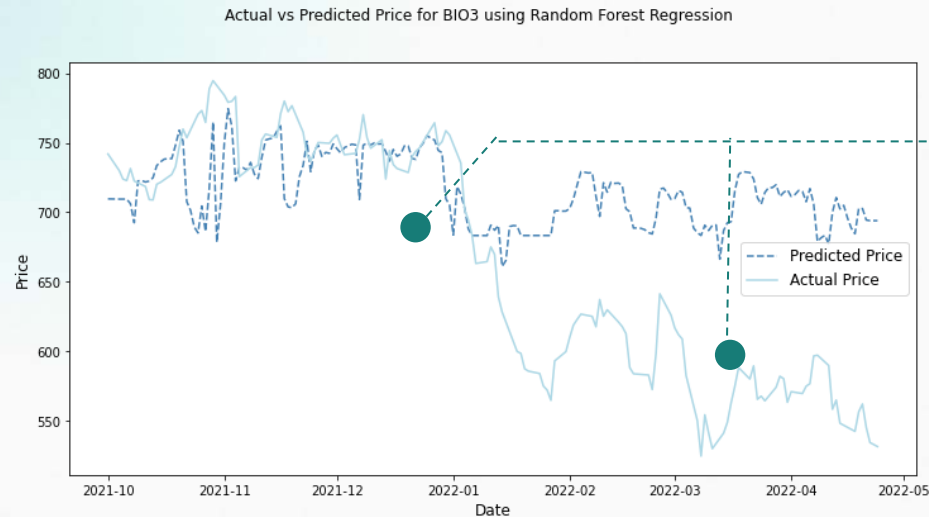
- Potential Reason:
  - similar datapoint indicates the similar time period?
  - similar time period datapoint has similar trend?
  - As  $k$  increase, it tends to differ a lot



### Gradient Boosting accuracy remain unchanged as not matter what learning rate to choose

- Potential Reason:
  - the data itself might not be so informative

## Phase 2: Problems



### Predict well at first but then worsen

- Potential Reason:  
the window method works?
- Possible Solution:  
use smaller window?

Shall we use the current day information for the current day stock price prediction?

- Possible Solution:  
Shift Y and use previous X to predict current Y?

Only use training data to check accuracy?

- Possible Solution:  
Validation data?

Manually selecting parameters causing not choosing the optimal parameters?

- Possible Solution:  
Use Loop to select the optimal parameters

# Phase 3: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

## Phase Three

## Outcomes

### Data Processing

Phase 1 + 2

- Transfer all data into daily data
- Fill weekend and missing data
- Normalization
- Change in variable

Based on Phase 1 + 2, We add:

- 20 windows are applied to use the closer date to predict the data

### Feature selection

Phase 1 + 2

- Split into training and testing
- Select feature using Information Gain Method
- Add Fisher's score to select feature together with information gain method

Based on Phase 1 + 2, We add:

- Use validation data to select feature instead of training data
- Loop to select optimal bar for feature selection

### Modeling

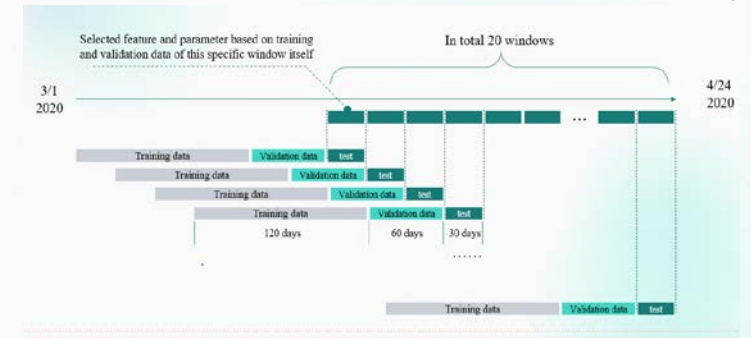
Phase 1 + 2

- Focus on 2 stocks instead of 24 stocks
- Divide into 5 smaller windows
- tried 7 models: Decision Tree, SVM, Random Forest, etc.

Based on Phase 1 + 2, We add:

- Loop to select optimal number for model parameter
- Shift Y for one day, use the Xs of day t-1 to predict Y of day t

## Window illustration



Overview

Phase: 1

Phase 2

Phase 3

Summary

# Phase 3: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

## Phase Three

## Outcomes

### Data Processing

Phase 1 + 2

- Transfer all data into daily data
- Fill weekend and missing data
- Normalization
- Change in variable

Based on Phase 1 + 2, We add:

- 20 windows are applied to use the closer date to predict the data

### Feature selection

Phase 1 + 2

- Split into training and testing
- Select feature using Information Gain Method
- Add Fisher's score to select feature together with information gain method

Based on Phase 1 + 2, We add:

- Use validation data to select feature instead of training data
- Loop to select optimal bar for feature selection

### Modeling

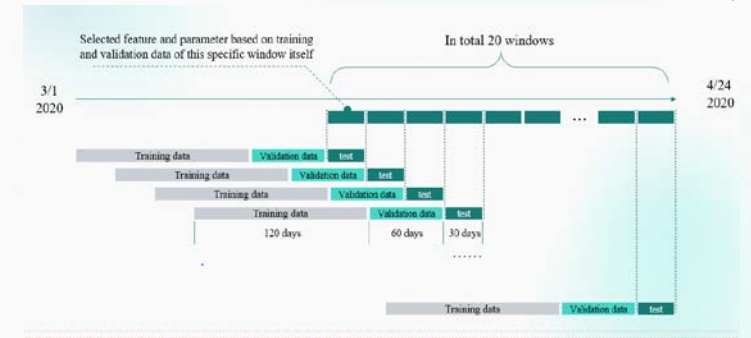
Phase 1 + 2

- Focus on 2 stocks instead of 24 stocks
- Divide into 5 smaller windows
- tried 7 models: Decision Tree, SVM, Random Forest, etc.

Based on Phase 1 + 2, We add:

- Loop to select optimal number for model parameter
- Shift Y for one day, use the Xs of day t-1 to predict Y of day t

## Window illustration



Overview

Phase: 1

Phase 2

Phase 3

Summary

# Phase 3: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

## Phase Three

## Outcomes

### Data Processing

Phase 1 + 2

- Transfer all data into daily data
- Fill weekend and missing data
- Normalization
- Change in variable

Based on Phase 1 + 2, We add:

- 20 windows are applied to use the closer date to predict the data

### Feature selection

Phase 1 + 2

- Split into training and testing
- Select feature using Information Gain Method
- Add Fisher's score to select feature together with information gain method

Based on Phase 1 + 2, We add:

- Use validation data to select feature instead of training data
- Loop to select optimal bar for feature selection

### Modeling

Phase 1 + 2

- Focus on 2 stocks instead of 24 stocks
- Divide into 5 smaller windows
- tried 7 models: Decision Tree, SVM, Random Forest, etc.

Based on Phase 1 + 2, We add:

- Loop to select optimal number for model parameter
- Shift Y for one day, use the Xs of day t-1 to predict Y of day t

## Window illustration



Overview

Phase: 1

Phase 2

Phase 3

Summary

# Phase 3: Data Processing

Goal: Increase the Stock Price Prediction Accuracy

## Phase Three

## Outcomes

### Data Processing

Phase 1 + 2

- Transfer all data into daily data
- Fill weekend and missing data
- Normalization
- Change in variable

Based on Phase 1 + 2, We add:

- 20 windows are applied to use the closer date to predict the data

### Feature selection

Phase 1 + 2

- Split into training and testing
- Select feature using Information Gain Method
- Add Fisher's score to select feature together with information gain method

Based on Phase 1 + 2, We add:

- Use validation data to select feature instead of training data
- Loop to select optimal bar for feature selection

### Modeling

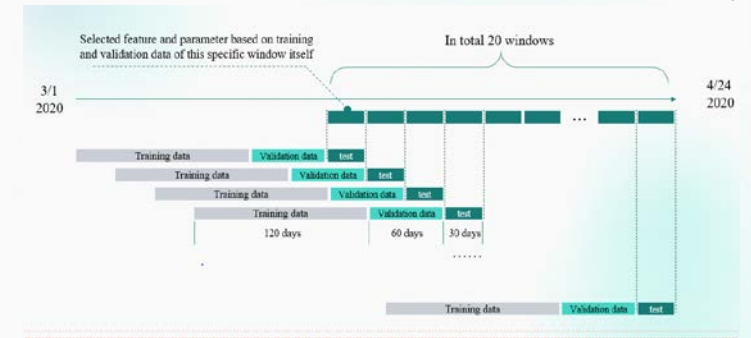
Phase 1 + 2

- Focus on 2 stocks instead of 24 stocks
- Divide into 5 smaller windows
- tried 7 models: Decision Tree, SVM, Random Forest, etc.

Based on Phase 1 + 2, We add:

- Loop to select optimal number for model parameter
- Shift Y for one day, use the Xs of day t-1 to predict Y of day t

## Window illustration



Overview

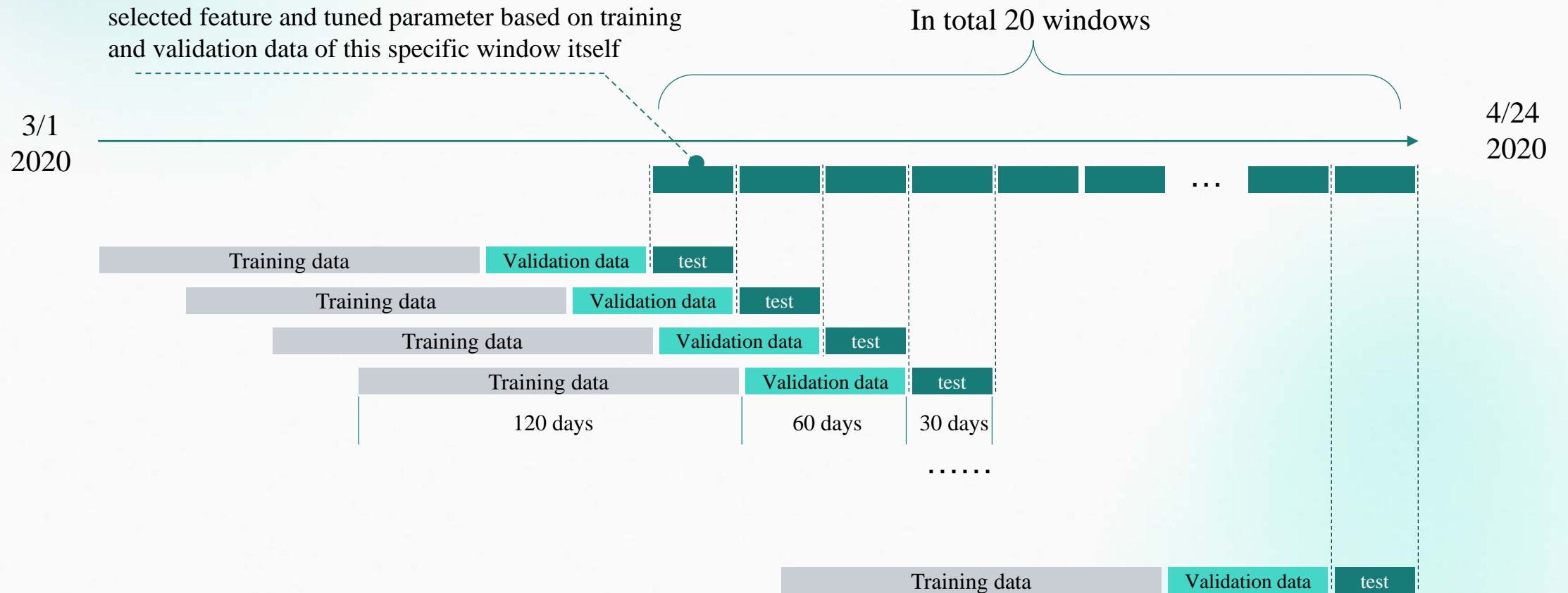
Phase: 1

Phase 2

Phase 3

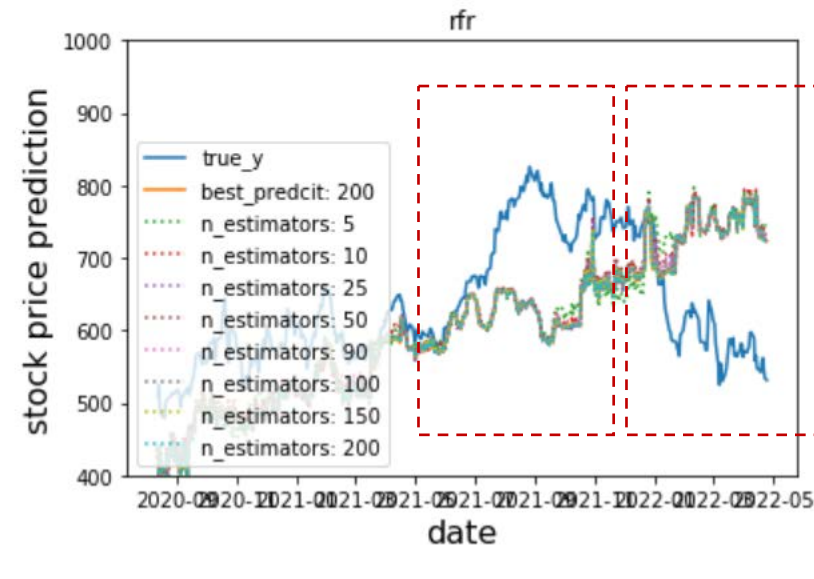
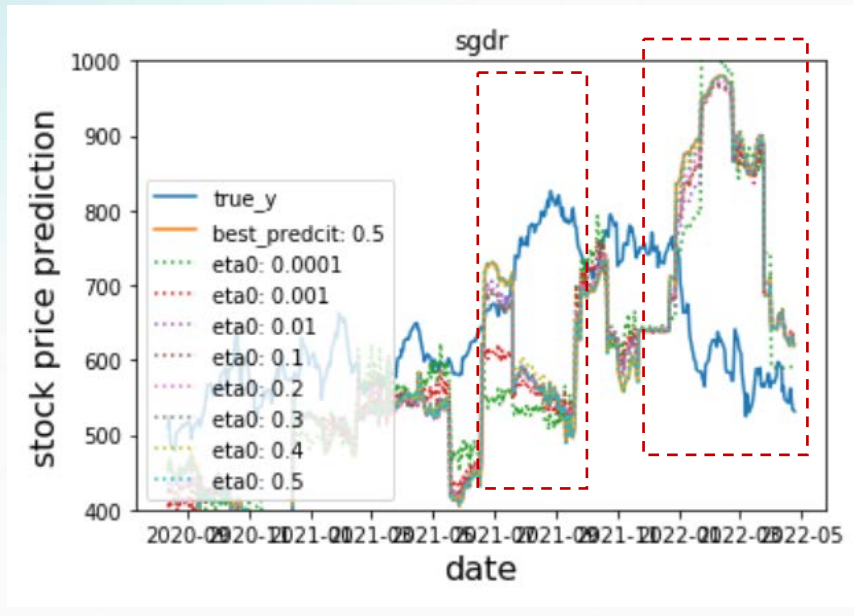
Summary

## Phase 3: Modeling





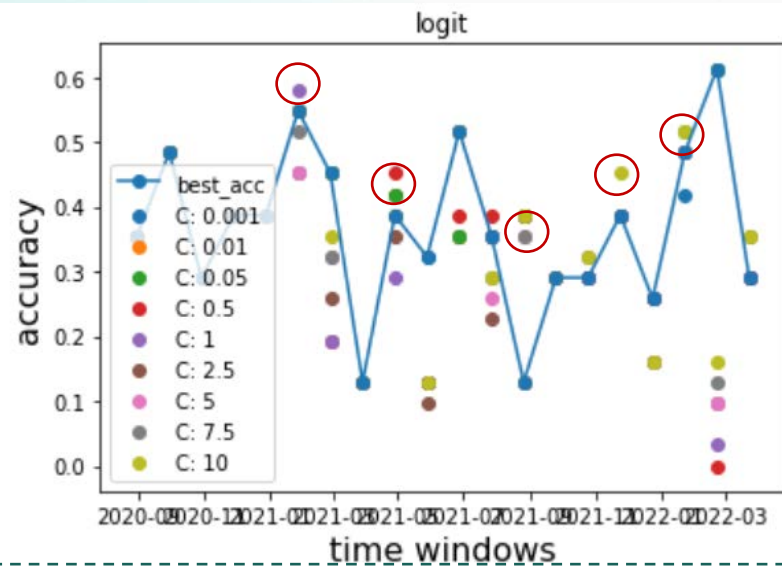
## Phase 3: a better result



**Did capture the smaller trend, but did not capture the larger trend**

- Potential Reason:  
the factor capturing the general trend should have higher weight?
- Potential Solution:  
Ensemble methods?

## Phase 3: some interesting findings



Validation accuracy	Test accuracy
other: 0.39344262295081966	test acc: 0.41935483870967744
current best: 0.39344262295081966	test acc: 0.41935483870967744
other: 0.39344262295081966	test acc: 0.41935483870967744
other: 0.39344262295081966	test acc: 0.41935483870967744
other: 0.39344262295081966	test acc: 0.41935483870967744
other: 0.4426229508196721	test acc: 0.41935483870967744
current best: 0.4426229508196721	test acc: 0.41935483870967744
other: 0.4426229508196721	test acc: 0.3548387096774194
other: 0.4262295081967213	test acc: 0.45161290322580644
other: 0.4098360655737705	test acc: 0.5161290322580645
other: 0.4098360655737705	test acc: 0.5806451612903226

VS

The parameter having higher validation accuracy might not always mean a better accuracy in testing

# Summary: What do we learned?



Real-life accuracy is not always as expected.



Be careful with data.



Ethical issues