

# PORTFOLIO INVESTING USING AN ENSEMBLE OF DATA SCIENCE TECHNIQUES

[UNSUPERVISED MACHINE LEARNING AND LONG SHORT TERM  
MEMORY (LSTM) DEEP LEARNING TECHNIQUES]

capstone project by Jeffrey Sim  
GA, DSI-18

*disclaimer: all material presented herein are for academic illustration purposes and should not be considered as investment advice.*

## Agenda

00

Introduction

11

Problem Statement

22

Data Science  
Methodology  
Overview

33

Data Collection  
EDA

44

Portfolio Building

55

Repeatable  
Prediction Model

66

Performance Conclusion  
Recommendation



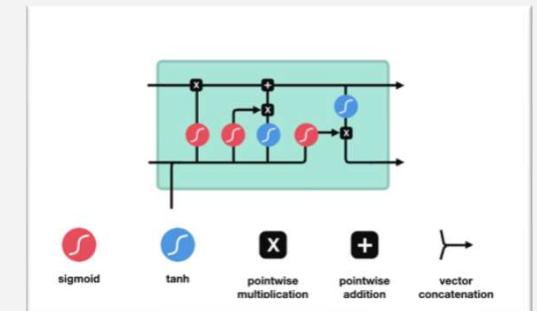
**“every crisis brings  
opportunities, and this is no  
lesser in the world of investing”**





## Traditional Way

“by the time you hear it  
it is too late”



## Data Science Techniques

# 11

## Problem Statements

1. **Too time consuming to collect and analyze information to build a viable portfolio**
2. **Not having a reliable and repeatable way to predict stock price movement**
3. **Not having an automated trading strategy that wins consistently**

Metrics of success

**Our modeled portfolios should achieve:**

1. **Above 50% of right calls**
2. **Positive returns higher than bank interest (@ 1.5% per yr)**



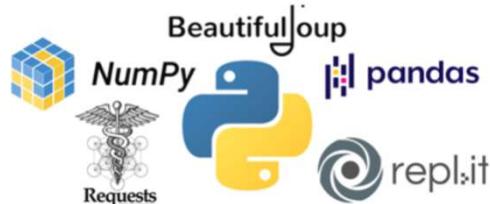
Solution  
Methodology  
Overview

22

1

## Data Pipeline & Collection

- **Web scrapping** with BeautifulSoup, Regex, JSON, yfinance API, etc.
- Fundamentals indicators, Analyst's Recommendations, News articles, Historical prices
- Valence Aware Dictionary and Sentiment Reasoner (**VADER**)



## Unsupervised Machine Learning

EDA, Domain Knowledge & Kmeans Clustering

- **EDA** to create "**Human portfolio**"
- **Standard Scaler pre-processing**
- **Kmeans** clustering to create "**AI portfolio**"

2

## Ensemble of techniques

4

## Trading Strategy Back Testing

- **Simple Buy low, Sell High, No Short**
- **Trading Account Class & Psuedo Trading Bot**
- Long-Term & Short-Term Performance



## LSTM Modeling & Tuning

- **MinMax Scaler** pre-processing
- **Generic LSTM model** with fixed layer and nodes for all (i.e. un-tune baseline)
- Keras RandomSearch **Tuned LSTM models** for each company (i.e. tuned portfolios)
- Trained with 20yrs data

## Deep Machine Learning Neural Networks

# 33

## Data Collection & EDA

yfinance

finviz

**yahoo!**  
finance

**MORNINGSTAR**

**MarketWatch**

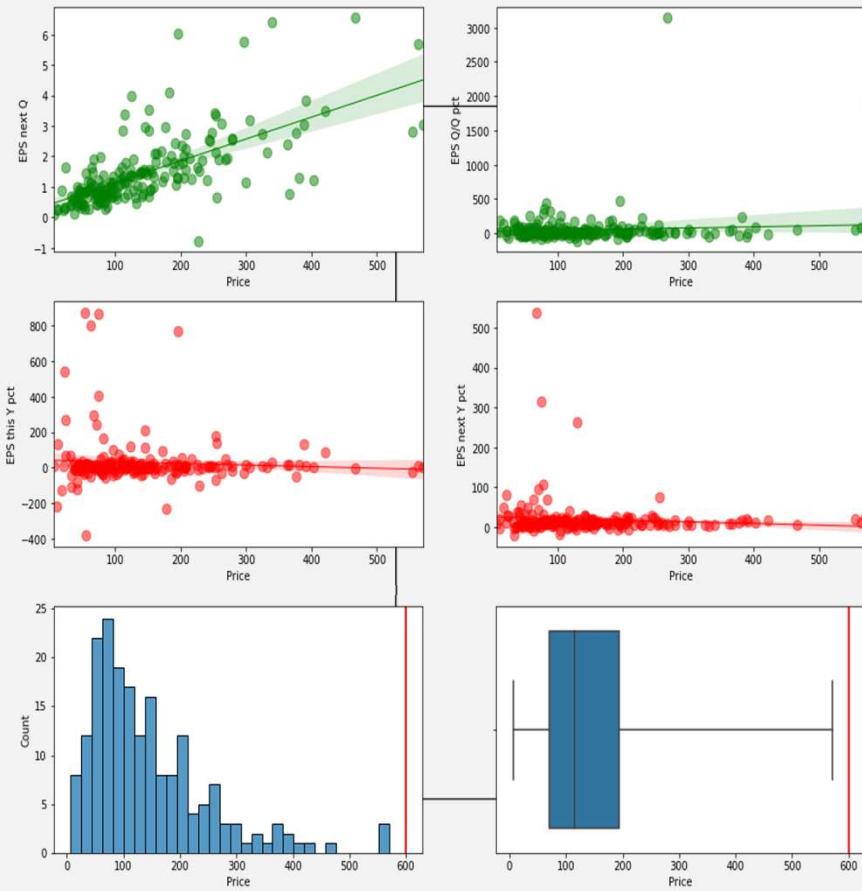
**► ALPHA VANTAGE**

**Google Finance**



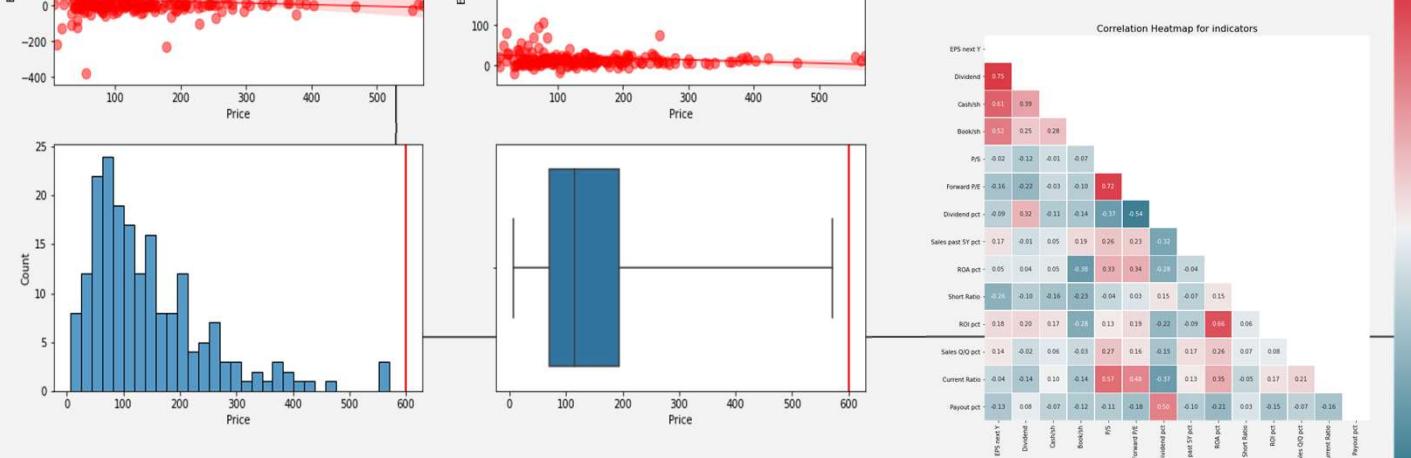
*“snapshots of company’s financial performance and how the market views its performance”*

## fundamental indicators



**Market Capitalization, Volume and Price**  
**Price-to-something (P/E, P/B, P/FCF)**  
**Earning per Share**  
**Insider, Institutional and shares float/outstd**  
**Technical Indicators (SMA, RSI)**  
**Dividends**  
**Profits, Earnings, Debts**  
**Financial ratios (ROE, ROI, ROA)**

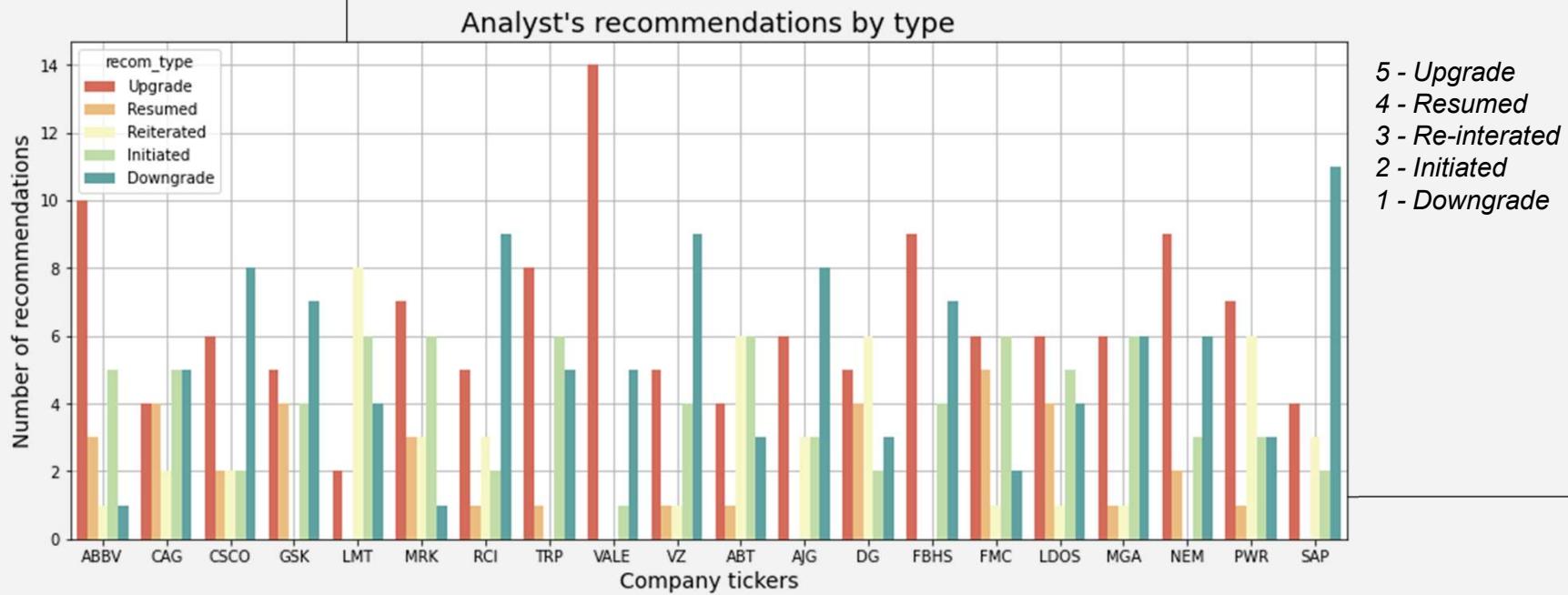
**DATASET ONE**  
**for Unsupervised Machine**  
**Learning Automated**  
**Portfolio Building**

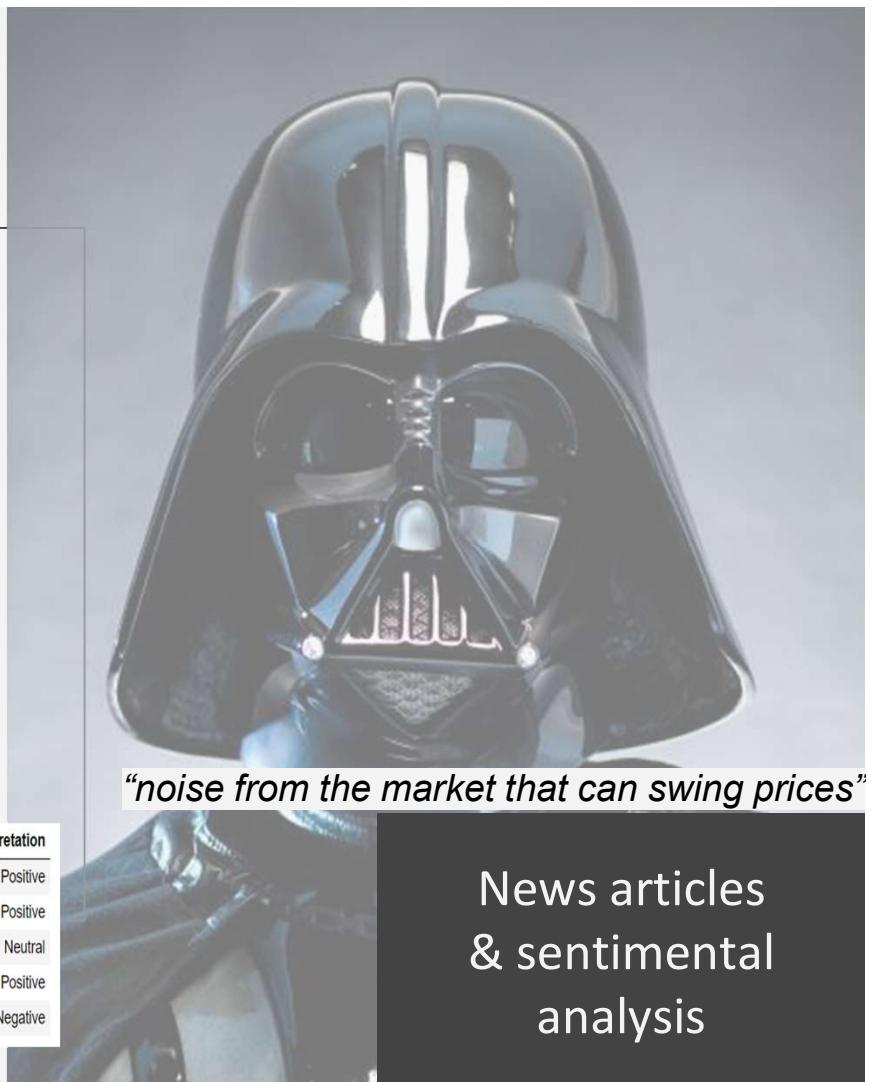
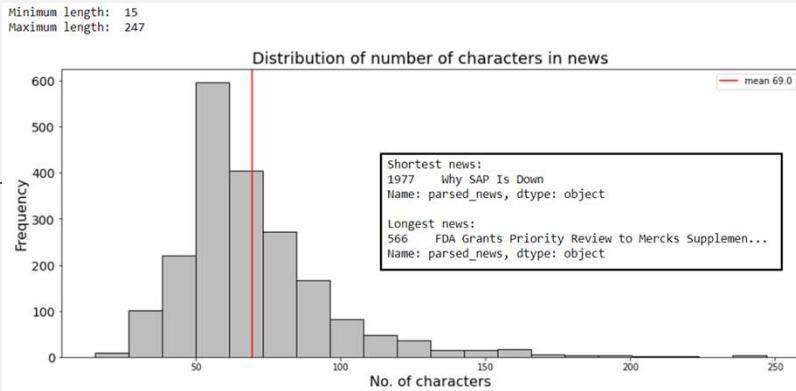


## analyst's recommendations

*"opinions from financial professionals who had spent a lot of time studying the financial reports a company"*

400 recommendations for 20 companies spanning from **2015-10-23 to 2021-01-21**





## VADER (Valence Aware Dictionary and Sentiment Reasoner)

2000 articles for 20 companies spanning from 2020-01-03 to 2021-01-22

positive sentiment: compound score  $\geq 0.05$

neutral sentiment: (compound score  $> -0.05$  and  $< 0.05$ )

negative sentiment: compound score  $\leq -0.05$

	ticker	date	time		parsed_news	neg	neu	pos	compound
494	LMT	2020-12-08	04:03PM	Boeing, Kratos To Build Skyborg Drones For Air...	0.00	1.000	0.000	0.0000	
1707	NEM	2021-01-05	08:56AM	Newmont to Provide Reserves and Exploration Up...	0.00	0.759	0.241	0.2263	
1122	AJG	2020-12-09	09:24AM	Arthur J. Gallagher (AJG) Announces Cool Insur...	0.00	0.753	0.247	0.3182	
996	VZ	2020-12-18	05:58PM	5G Auction Soars to \$34 Billion With Verizon L...	0.00	1.000	0.000	0.0000	
415	LMT	2021-01-12	03:03PM	More companies pause political donations follo...	0.32	0.680	0.000	-0.5106	

	ticker	date	compound	Interpretation
49	NEM	2021-01-05	0.226300	Positive
17	ABBV	2020-12-15	0.051033	Positive
17	VZ	2021-01-09	0.000000	Neutral
29	SAP	2020-12-31	0.102700	Positive
26	ABT	2020-12-05	-0.381800	Negative

historical prices  
& volumes

TA-Lib

#### Momentum indicators

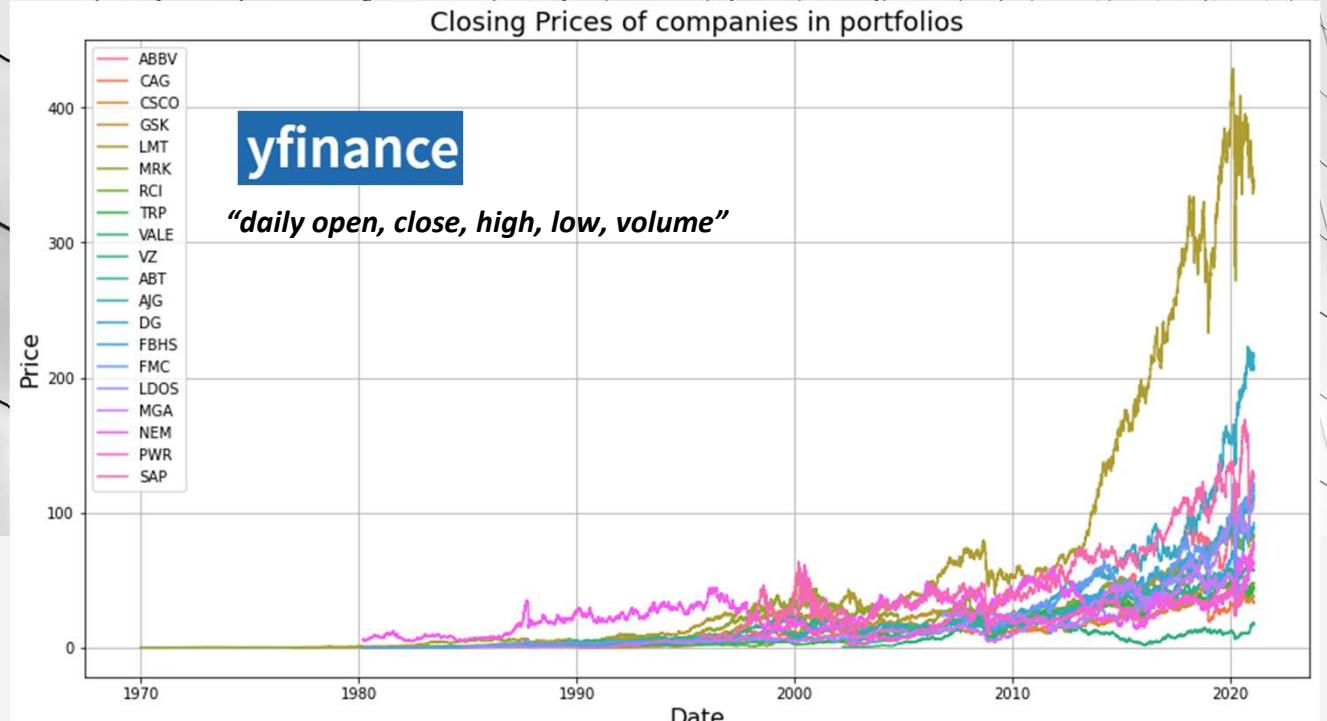
- SMA20 –Simple Moving Average
- RSI14 Relative Strength Index
- MOM5 – Momentum

#### Volume indicators

- OBV – On-balance Volume

#### Volatility indicators

- ATR14 – Average True Range



**DATASET TWO**  
for Deep Machine Learning (LSTM)  
Modeling for a repeatable price  
movement prediction model

#### Final merged data:

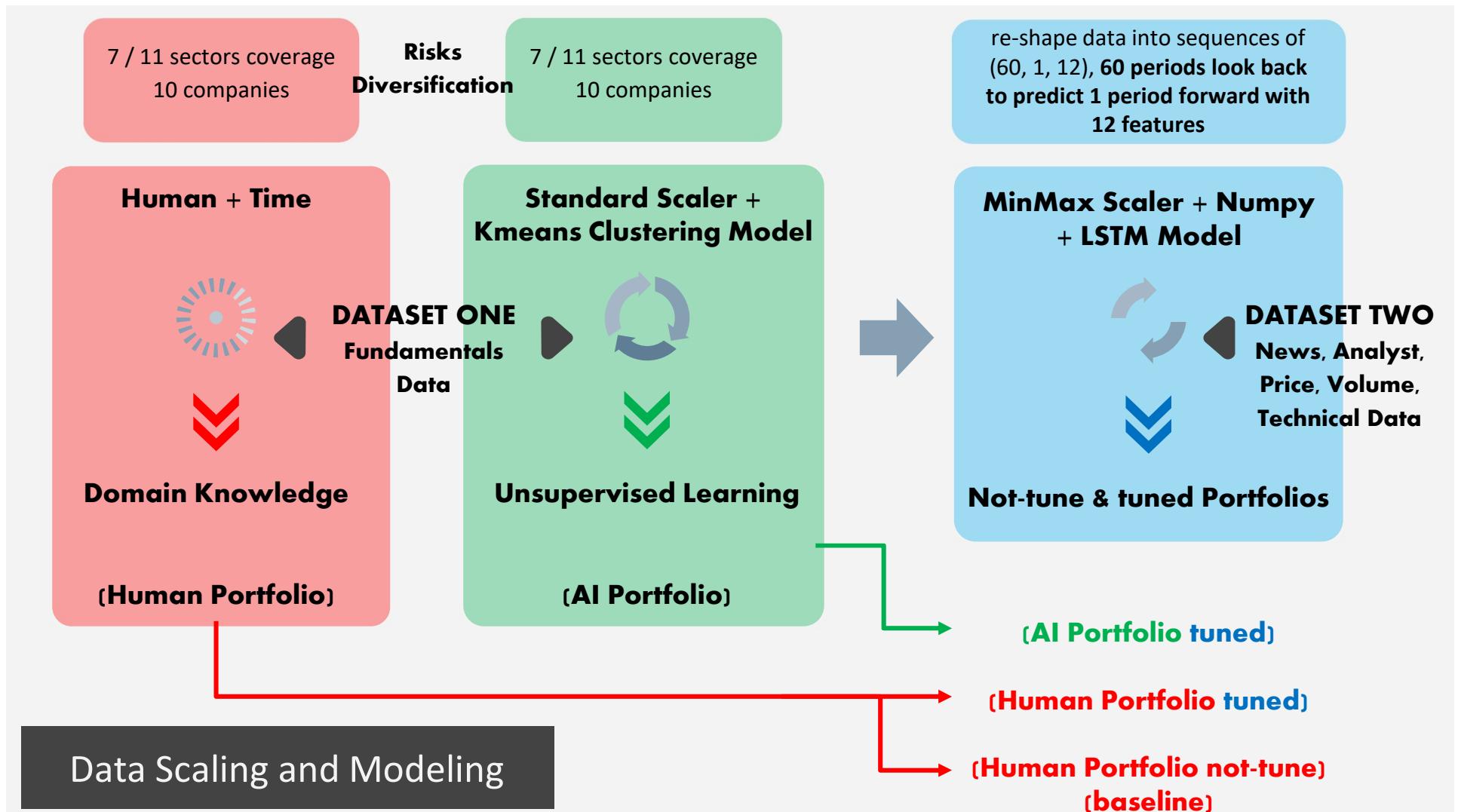
- News, Analysts, Historical Prices, Volumes, Technical Indicators
- 20 companies
- more than 150K+ rows
- 13 features



# Portfolio Building

Clustering model  
with Dataset One and  
Data Pre-processing

44



# 55

## Modeling Tuning Evaluation

Price movement prediction  
with Dataset Two

# Modeling

```
1 # instantiation
2 lstm_model = Sequential()
3
4 # set activation
5 activ = 'tanh'
6
7 # adding the input and first LSTM layer
8 lstm_model.add(LSTM(100, activation=activ, return_sequences=True,
9                     input_shape=(n_per_in, n_features)))
10
11 # # final hidden layer w/o return_sequences
12 lstm_model.add(LSTM(100, activation=activ))
13
14 # output layer
15 lstm_model.add(Dense(n_per_out))
executed in 492ms, finished 11:39:34 2021-01-24
```

```
1 lstm_model.summary()
executed in 15ms, finished 11:39:34 2021-01-24
```

Model: "sequential"

Layer (type)	Output Shape	Param #
<hr/>		
lstm (LSTM)	(None, 60, 100)	45200
lstm_1 (LSTM)	(None, 100)	80400
dense (Dense)	(None, 1)	101
<hr/>		
Total params:	125,701	
Trainable params:	125,701	
Non-trainable params:	0	

## Generic LSTM (not-tune)

```
3 def build_model(hp):
4
5     # instantiation
6     lstm_model = Sequential()
7
8     # set activation
9     activ = 'tanh'
10
11    # adding the input and first LSTM layer
12    lstm_model.add(LSTM(units=hp.Int("input_units", min_value=64, max_value=256, step=32),
13                         activation=activ, return_sequences=True,
14                         input_shape=(n_per_in, n_features)))
15
16    # hidden layers with Dropout regularisation
17    for i in range(hp.Int("n_layers", 1, 2)):
18        lstm_model.add(LSTM(units=hp.Int("hidden_{i}_units", min_value=64, max_value=256, step=32),
19                            activation=activ, return_sequences=True))
20        lstm_model.add(Dropout(0.2))
21
22    # final hidden layer w/o return_sequences
23    lstm_model.add(LSTM(units=hp.Int("final_units", min_value=64, max_value=256, step=32),
24                        activation=activ))
25
26    # output layer
27    lstm_model.add(Dense(n_per_out))
28
29    # Compiling the RNN
30    lstm_model.compile(optimizer=Adam(lr=0.01), loss='mse')
31
32    return lstm_model
33
34 #initializing randomsearch
35 tuner = RandomSearch(build_model,
36                       objective = 'val_loss',
37                       max_trials = 3, # how many random permutations do you want
38                       executions_per_trial = 2, # how many time to train for each permutations
39                       directory = LOG_DIR,
40                       project_name = 'SEEDED(42) RandomSearch')
```

trials results

```
1 # display the search results and params
2 tuner.results_summary()
executed in 15ms, finished 00:59:18 2021-01-25

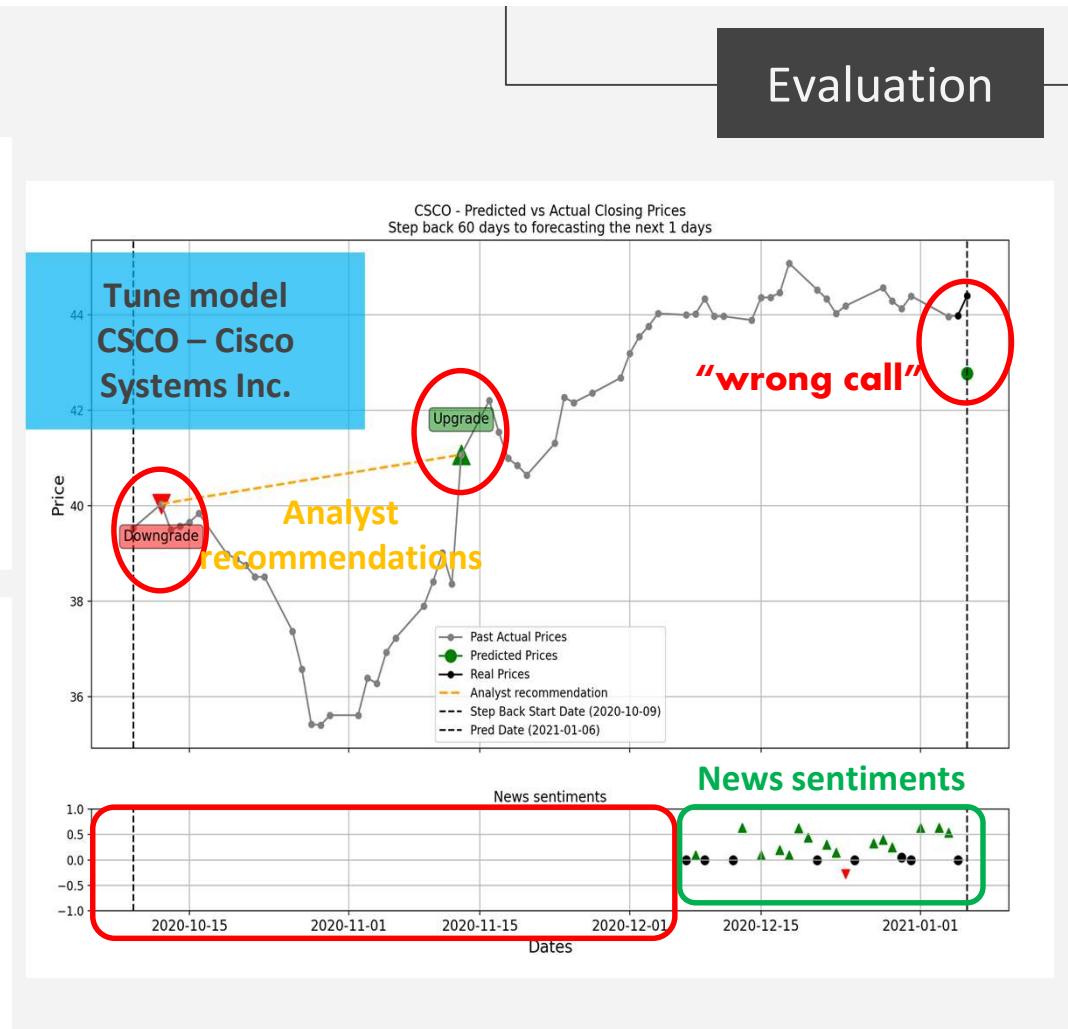
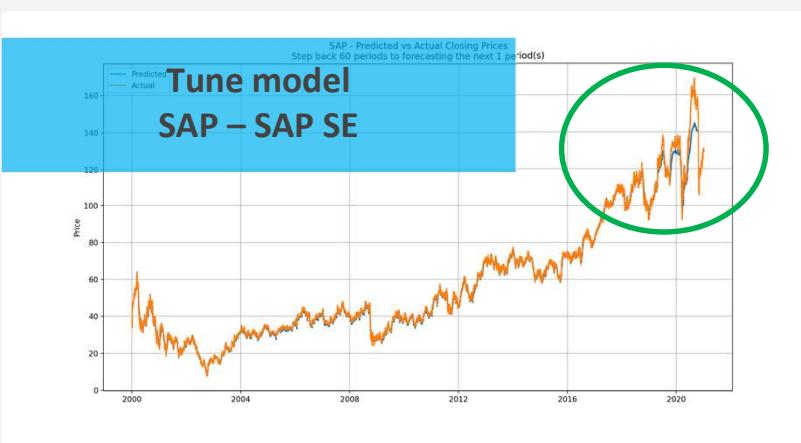
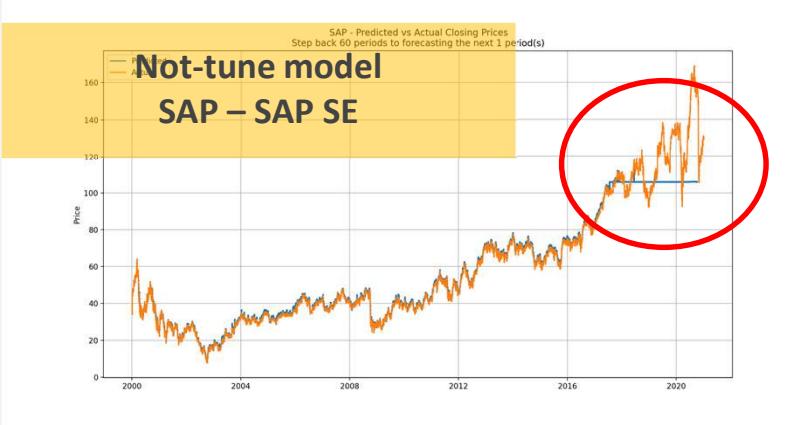
Results summary
Results in keras_tuner_logs/RCI_1611507447\S
Showing 10 best trials
Objective(name='val_loss', direction='min')
Trial summary
Hyperparameters:
input_units: 160
n_layers: 1
hidden_0_units: 160
final_units: 64
hidden_1_units: 160
Score: 0.00032069699955172837
Trial summary
Hyperparameters:
input_units: 192
n_layers: 2
hidden_0_units: 64
final_units: 224
hidden_1_units: 64
Score: 0.003214719356037673
Trial summary
Hyperparameters:
input_units: 128
n_layers: 1
hidden_0_units: 160
final_units: 128
hidden_1_units: 192
Score: 0.17871791124343872
```

## Dynamic LSTM (tuned)



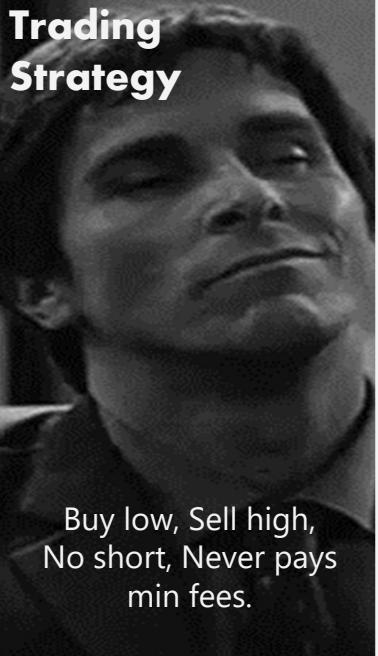
## Examples of prediction results

## Evaluation



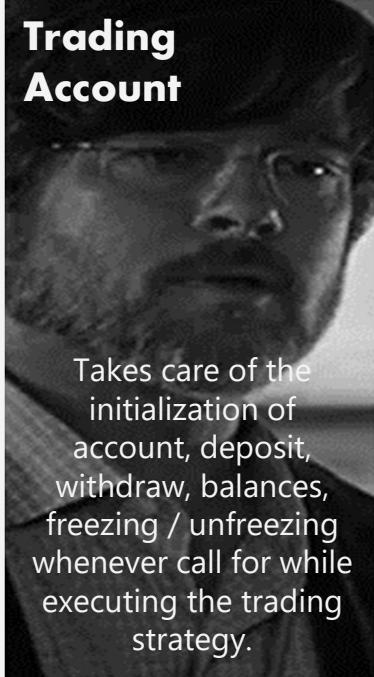
**Based on DBS OET US fees**

Commission = 0.15% on transaction amount  
Min. Fees = US 18 per transaction



**Trading  
Strategy**

Buy low, Sell high,  
No short, Never pays  
min fees.



**Trading  
Account**

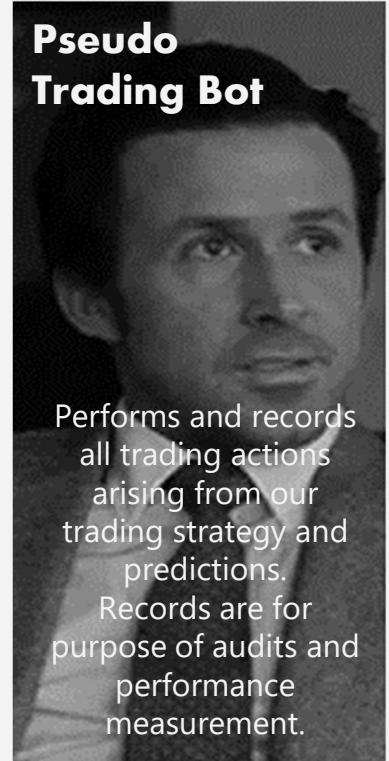
Takes care of the initialization of account, deposit, withdraw, balances, freezing / unfreezing whenever call for while executing the trading strategy.

**Account Starting Balance = 50K each company (i.e 500K each portfolio)**

**Min. predicted percentage gain = 5 %**

*Short Selling – selling what you do not have at a high price, hoping to buy the require quantity at a lower price (very high risk)*

**Meet Our Team**



**Pseudo  
Trading Bot**

Performs and records all trading actions arising from our trading strategy and predictions.

Records are for purpose of audits and performance measurement.

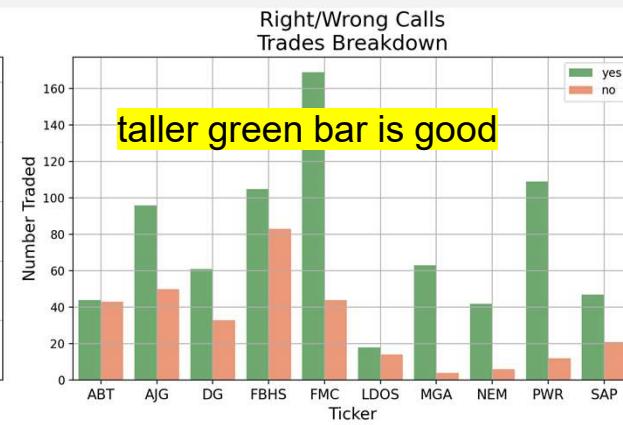
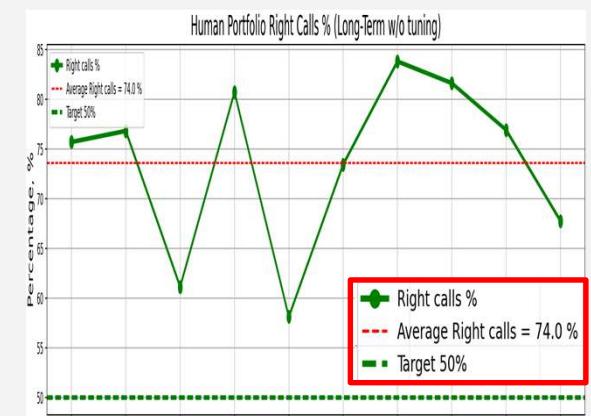
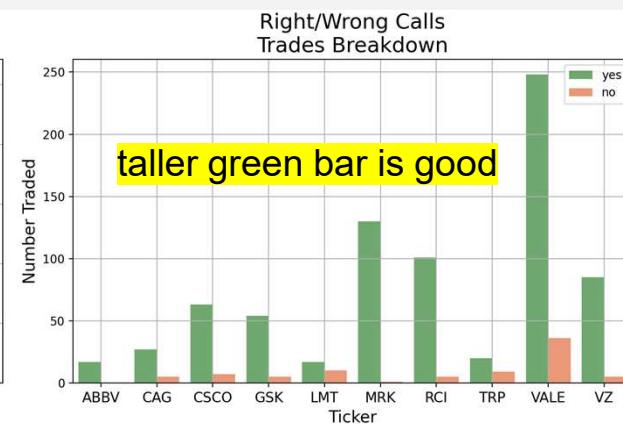
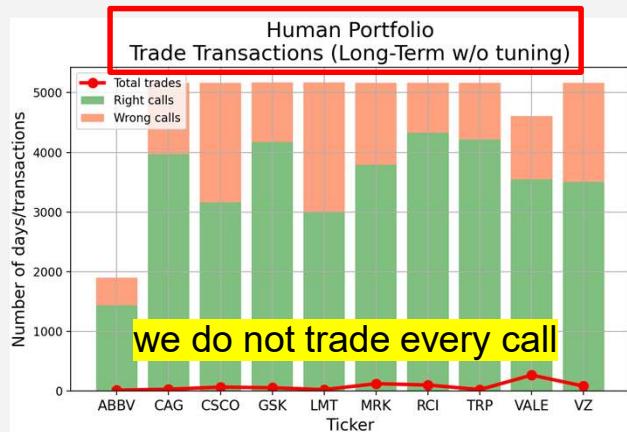


Conclusion  
Recommendation

66

## Conclusion

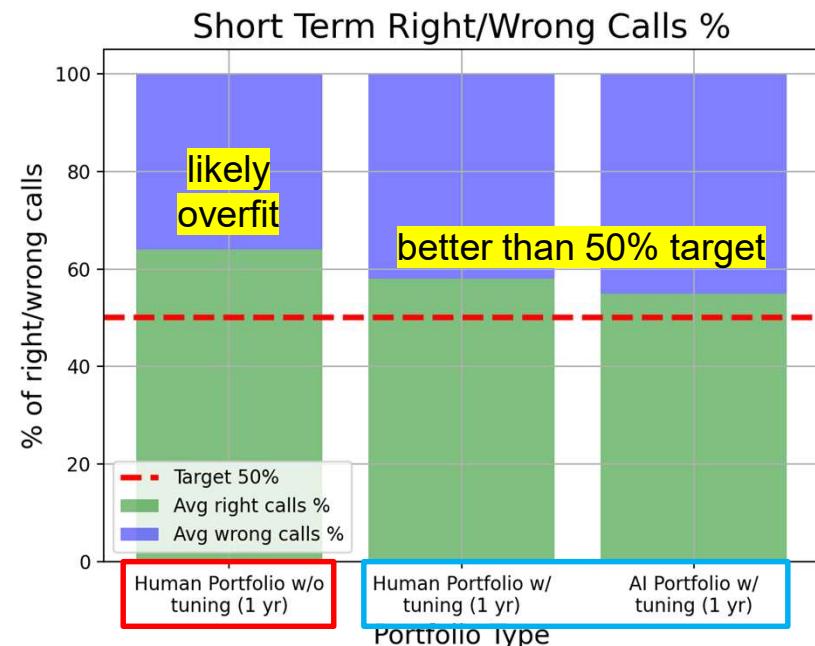
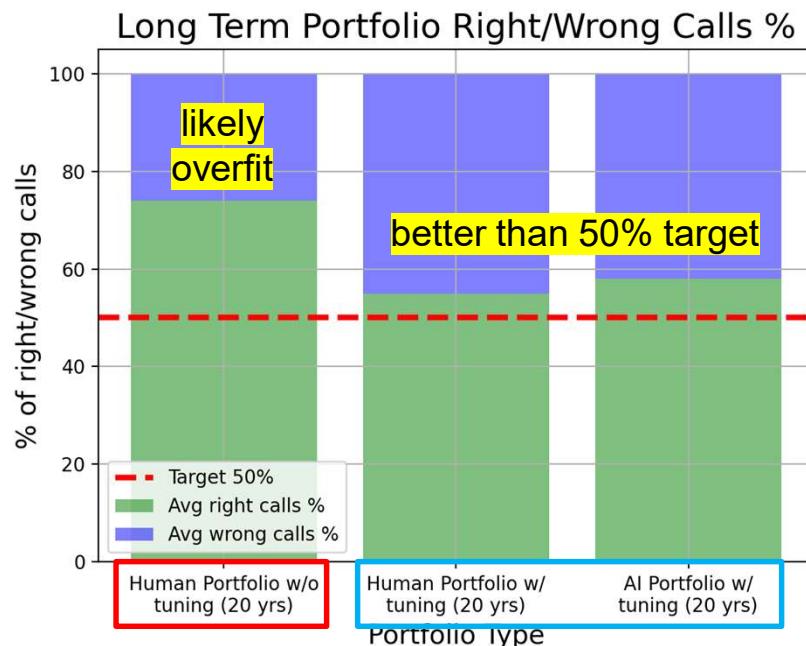
How often did we get it right?  
@ success metric - above 50% of right calls



## Conclusion

How often did we get it right?

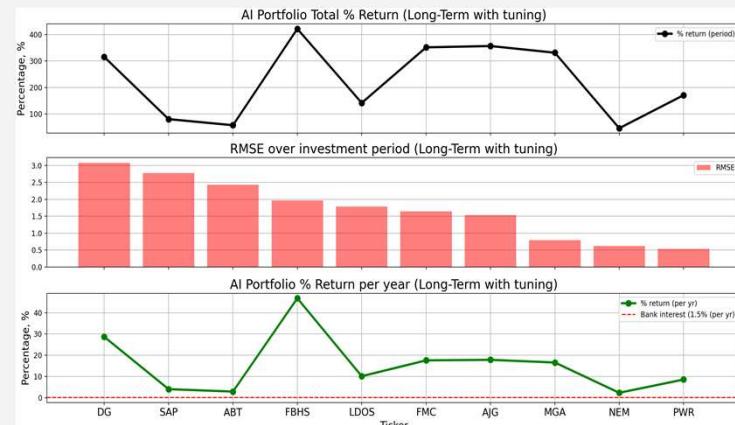
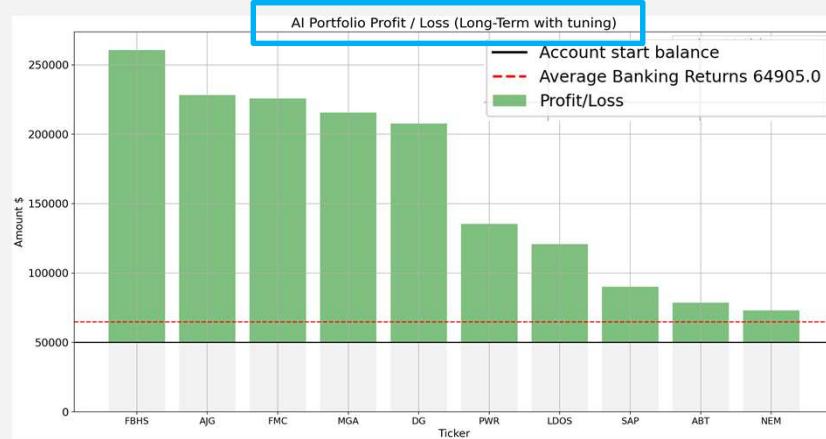
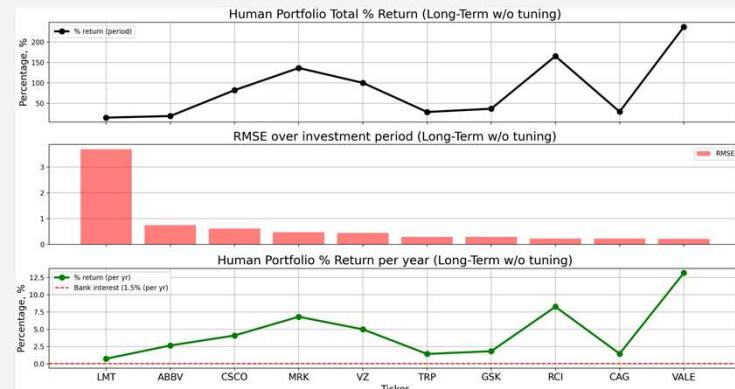
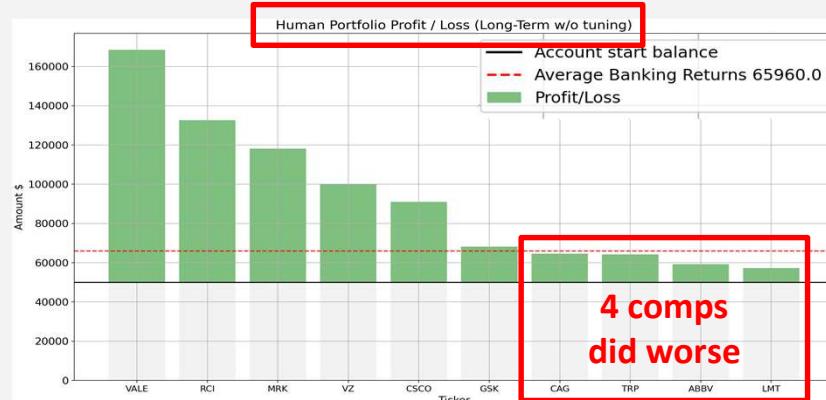
@ success metric - above 50% of right calls



## Conclusion

How much did we make or lose?

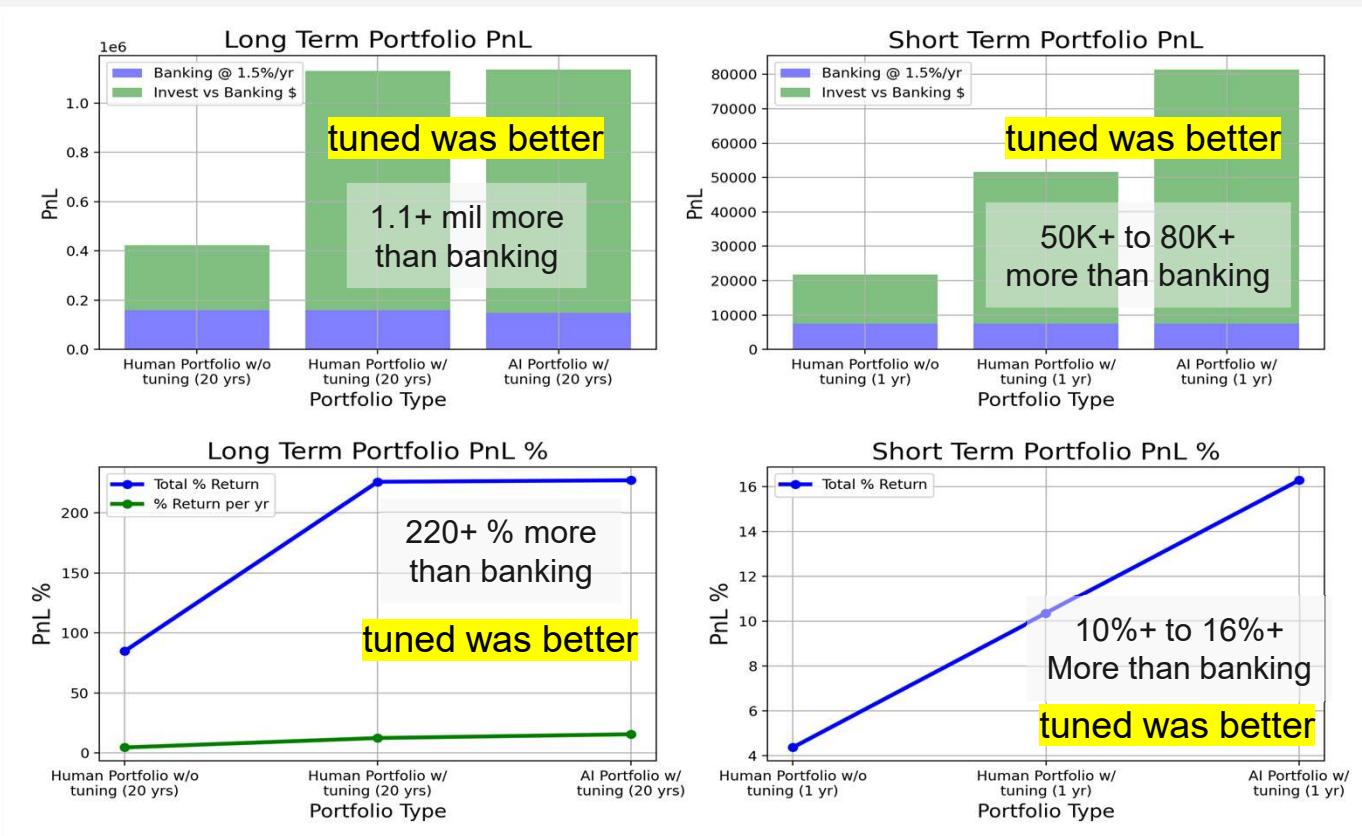
@ success metric - positive returns higher than bank interest (1.5% per yr)



## Conclusion

How much did we make or lose?

@ success metric - positive returns higher than bank interest (1.5% per yr)



# 11

## Problem Statements

BUSTED  
BUSTED  
BUSTED

1. Too time consuming to collect and analyze information to build a viable portfolio.
2. Not having a reliable and repeatable way to predict stock price movement.
3. Not having an automated trading strategy that wins consistently.

### Metrics of success

Our modeled portfolios should achieve:

1. Above 50% of right calls
2. Positive returns higher than bank interest (@ 1.5% per yr)



## Recommendations

- Getting more diversified data sources into the data pipe
- Making a few AI portfolios with more companies in each to further diversified the risk
- Building an image classification model to identify technical trends on candlestick charts
- Re-tune the dynamic model according to your investment horizons and time-frame
- Expanding on the trading strategy to include short selling and integrating it with trading brokerage APIs
- Deploying the solution on a platform with more computational resources
- Designing a dashboard to easily track performance regularly

## Credits

Our sincere and heartfelt thanks to Divya, Ryan, Alexis and the rest of GA Team for their un-reserving dedication to help us and for sharing their knowledge and experience.

Kudos to all classmates. We did it !!  
See you all in the data science world.



# The End

*disclaimer: all material presented herein are for academic illustration purposes and should not be considered as investment advice.*

