

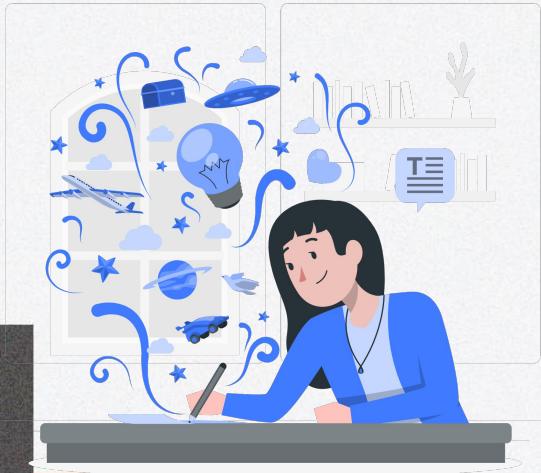
# SC1015 MINI PROJECT: STOCKS

Done by:

1. Yeoh Wei Yang (U2121112A)
2. Chantharojwong Kasidis (U2020731L)
3. Harvey Zhang Tianren (U2122097A)



# Agenda



1

Motivation

2

Exploratory  
data analysis

3

Model  
building

4

Applying  
our model

1

# Motivation

# Why stocks?



## Wealth

A method of generating income passively



## Inflation

Growth of money can beat out inflation



## Potential

Potential for high returns

# However!





Find the top factors that affect  
the performance of a stock, so  
as to develop a safe stock  
investing strategy to ensure  
positive rates of return for  
investors while minimizing  
losses

2

# Exploratory data analysis

# Understanding the dataset

	# R&D Expense Gro...	# SG&A Expenses ...	▲ Sector	# 2015 PRICE VAR [...	# Class
730		-1	Financial Services Healthcare Other (2566)	17% 15% 67% 	
0.0	-0.1746		Consumer Defensive	-9.323275997445537	0
1.6484	1.7313		Consumer Defensive	-25.512192888957696	0
0.0	0.0234		Consumer Defensive	33.11829671550496	1
0.0	-0.006		Consumer Defensive	2.7522914680574364	1
0.0	-0.02200000000000002		Consumer Defensive	12.897715165910551	1
0.0	0.0161		Consumer Defensive	13.980936777937483	1
0.0	-0.0053		Consumer Defensive	5.339412837835392	1
0.0	0.0307		Consumer Defensive	-26.65370178517993	0
0.0	-0.0256		Consumer Defensive	23.809817736521	1

## Summary

- ▶ 5 files
- ▶ 1125 columns

# Understanding the dataset



**In the dataset:**

**1) 200+ financial indicators of a year**

**2) Price Var [%] :**

- Measure of percent price variation of the next year

**3) Class :**

- Binary classification
- "0" for stocks which value will drop (don't buy)
- "1" for stocks which value will increase (buy)



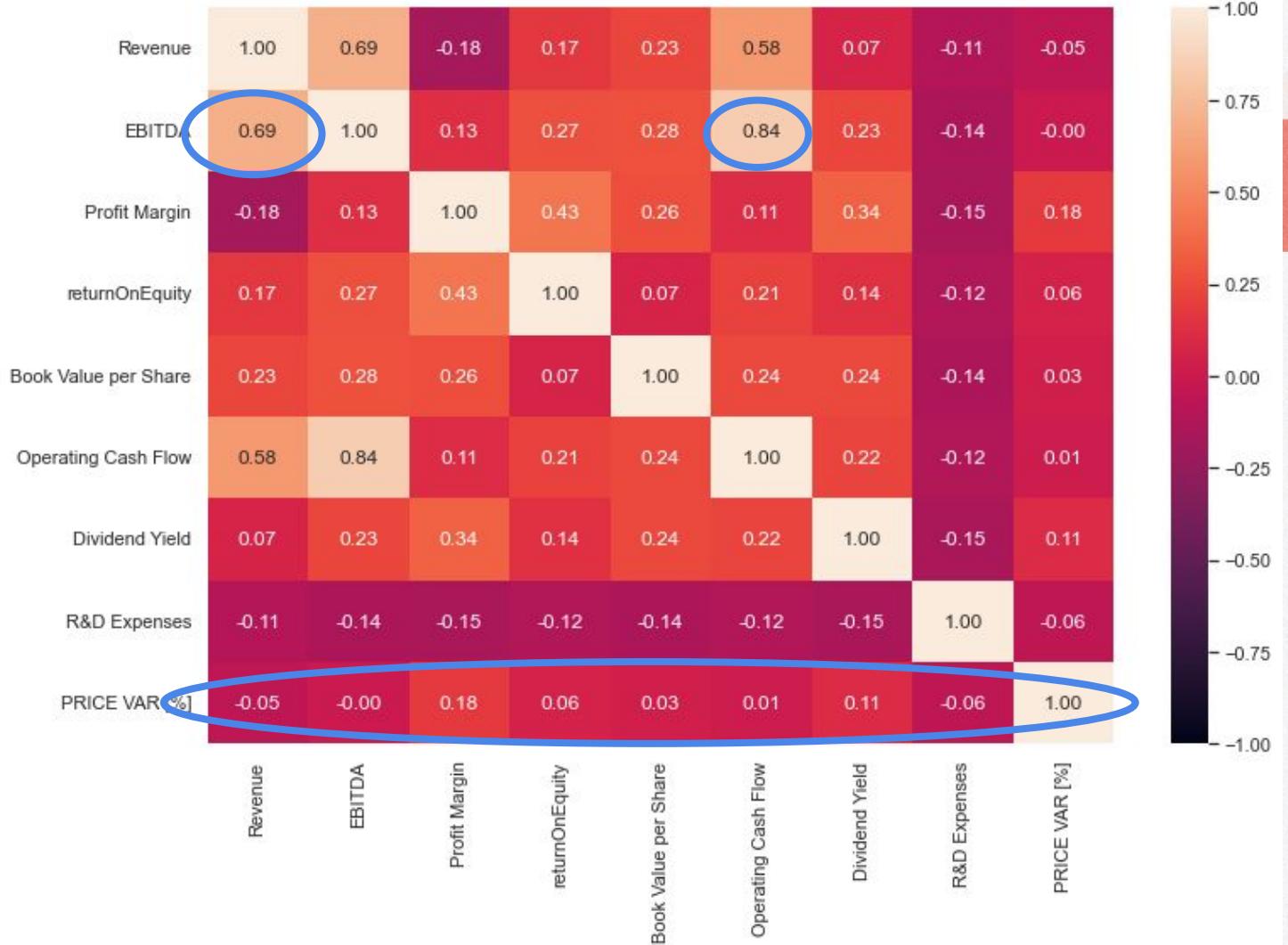
# Data preparation and cleaning

- ◊ Check for missing values
- ◊ Fill missing values with median
- ◊ Remove outliers
- ◊ Select a few indicators to work with

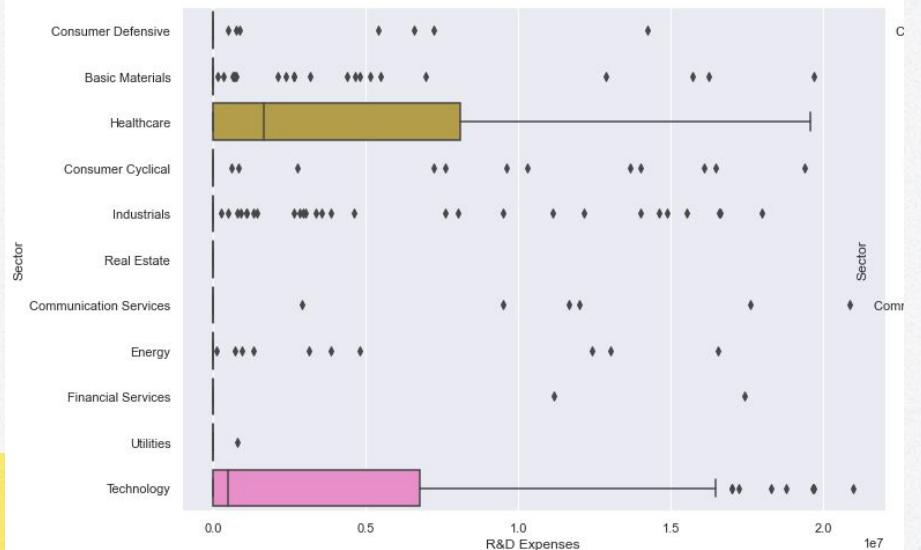


Out[13]:		Revenue	EBITDA	Profit Margin	returnOnEquity	Book Value per Share	Operating Cash Flow	Dividend Yield	R&D Expenses	PRICE VAR [%]	Class	Sector
14	5.727000e+09	683400000.0	0.066	0.2041		6.706	634100000.0	0.0173	0.0	18.603845	1	Consumer Defensive
22	4.551600e+09	241900000.0	-0.021	-0.1154		2.211	536500000.0	0.0117	0.0	26.681241	1	Consumer Defensive
23	2.464867e+09	771439000.0	0.196	0.3189		3.020	597491000.0	0.0000	0.0	37.721889	1	Consumer Defensive
30	3.297600e+09	743500000.0	0.126	0.1969		7.779	540300000.0	0.0157	0.0	10.792636	1	Consumer Defensive
32	5.973810e+08	183876000.0	0.122	0.1289		10.909	111582000.0	0.0000	0.0	49.607672	1	Consumer Defensive
...	...	...	...	...	...	...	...	...	...	...	...	...
3802	1.185080e+08	9650000.0	0.034	0.1324		4.249	7612000.0	0.0000	0.0	-2.453386	0	Technology
3803	4.952987e+07	-53213.0	-0.002	-0.0097		4.505	523987.0	0.0000	0.0	29.362884	1	Technology
3804	1.532400e+08	20887000.0	0.085	0.3646		2.426	-1587000.0	0.0000	11326000.0	-31.167763	0	Technology
3806	3.407580e+08	8512000.0	0.017	0.1456		8.489	5745000.0	0.0395	0.0	7.779579	1	Technology
3807	4.033737e+07	4959141.0	0.060	0.0721		1.645	4012331.0	0.0000	3379920.0	-34.099613	0	Technology

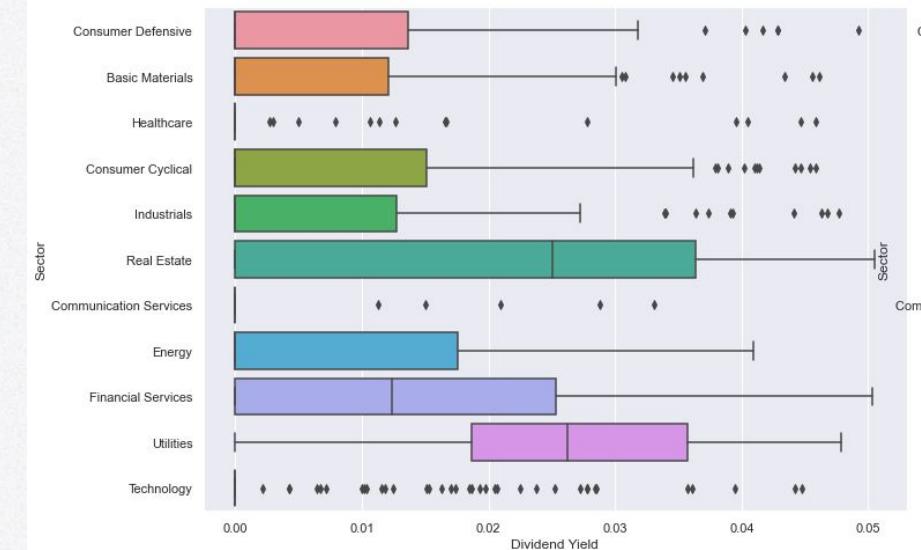
1521 rows × 11 columns

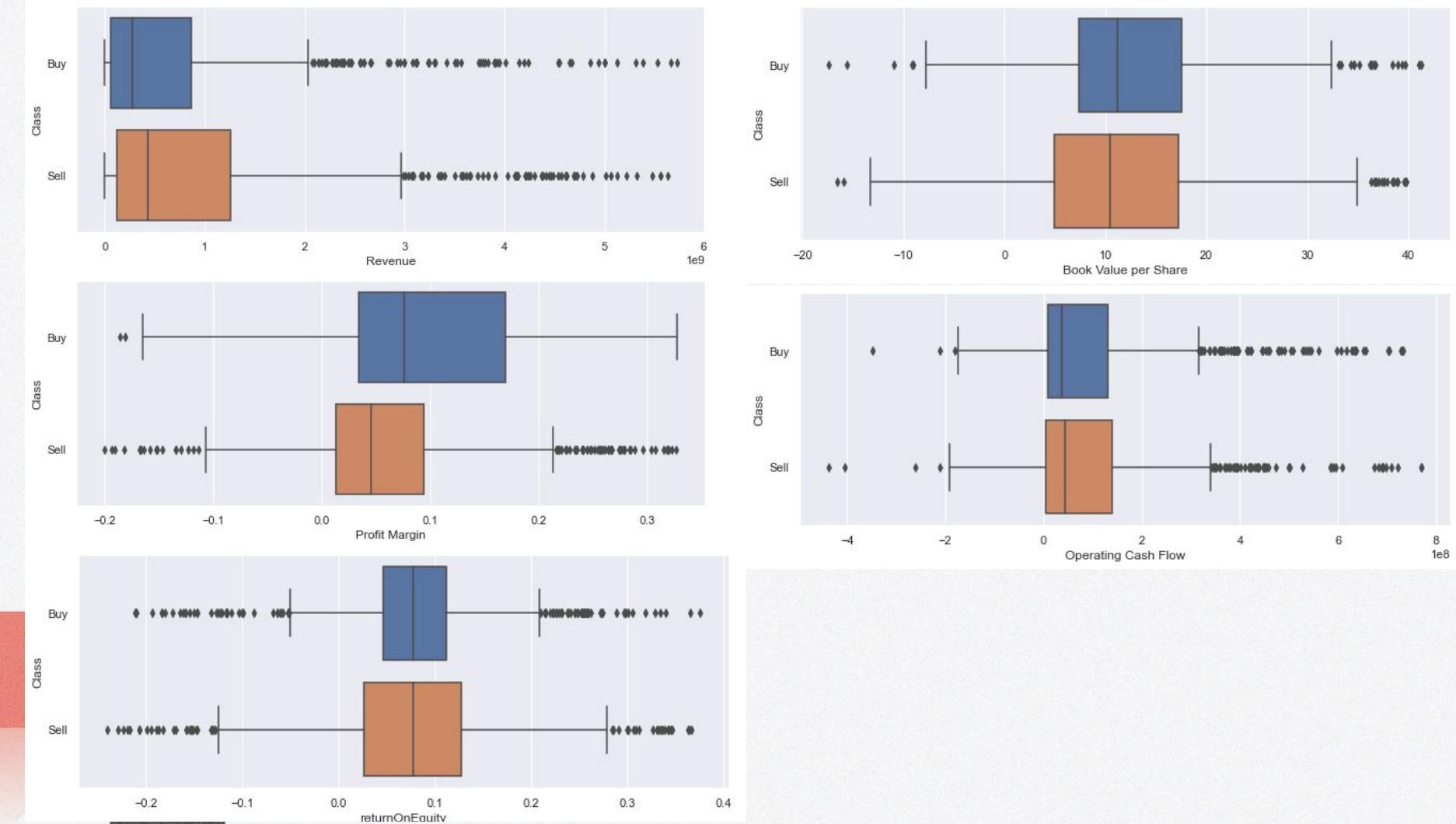


## R&D expenses against sector



## Dividend yield against sector



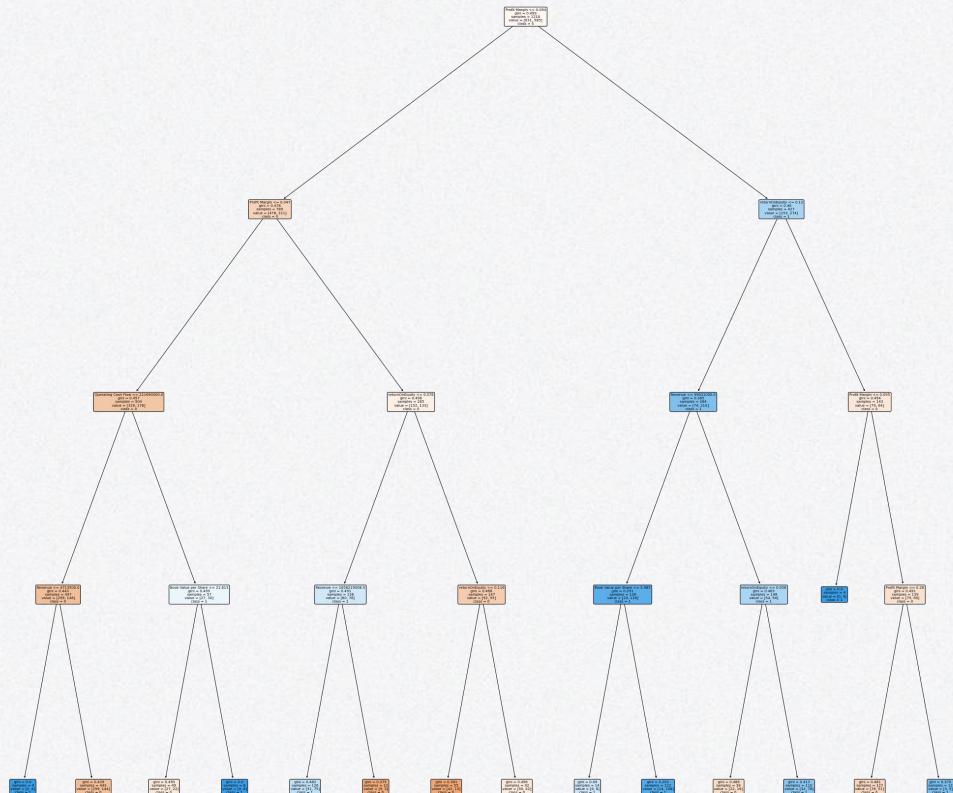


# 3

# Machine learning

- MultiVariate Decision Tree
  - Random Forest
  - GridsearchCV

# Decision Tree



- A Decision Tree is used to classify the response variables against Class
- A classification accuracy of 0.67 is obtained on the train set and 0.57 for the test set
- (A more detailed image will be in our GitHub repository)

# Random Forest

## Train Dataset

Classification accuracy: 0.691

True Positive Rate: 0.561

True Negative Rate: 0.808

False Positive Rate: 0.438

False Negative Rate: 0.191



## Test Dataset

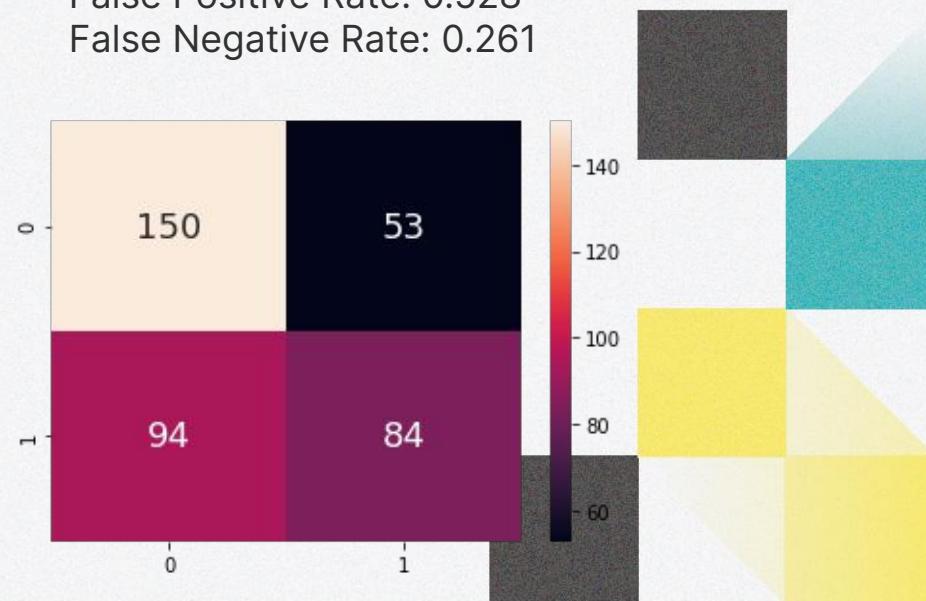
Classification accuracy: 0.614

True Positive Rate: 0.472

True Negative Rate: 0.738

False Positive Rate: 0.528

False Negative Rate: 0.261



# Cross Validation GridSearch

- Find the best hyperparameters for our random forest model
- GridSearch was run with `n_estimators` over a range of (100,400) and depth over a range of (2,4)
- K-fold Cross Validation, k is set to 5
- Best parameters: `n_estimators` = 128, depth = 4



# After GridSearchCV

## Train Dataset

Classification accuracy: 0.691 -> 0.695

True Positive Rate: 0.561 -> 0.569

True Negative Rate: 0.808 -> 0.809

False Positive Rate: 0.438 -> 0.430

False Negative Rate: 0.191 -> 0.190



## Test Dataset

Classification accuracy: 0.614 → 0.634

True Positive Rate: 0.472 → 0.488

True Negative Rate: 0.738 → 0.763

False Positive Rate: 0.528 → 0.511

False Negative Rate: 0.261 → 0.236



4

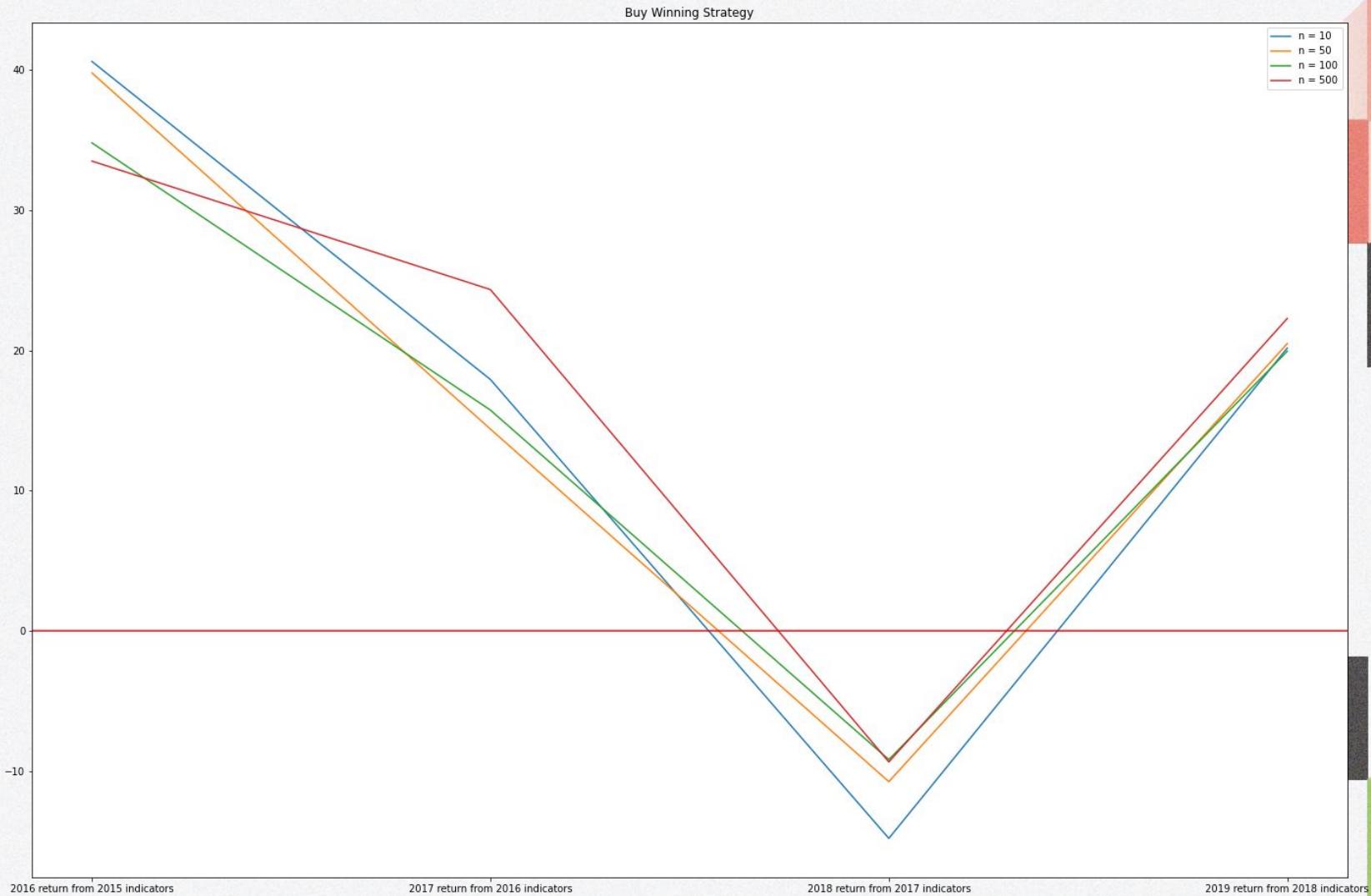
# Applying our model

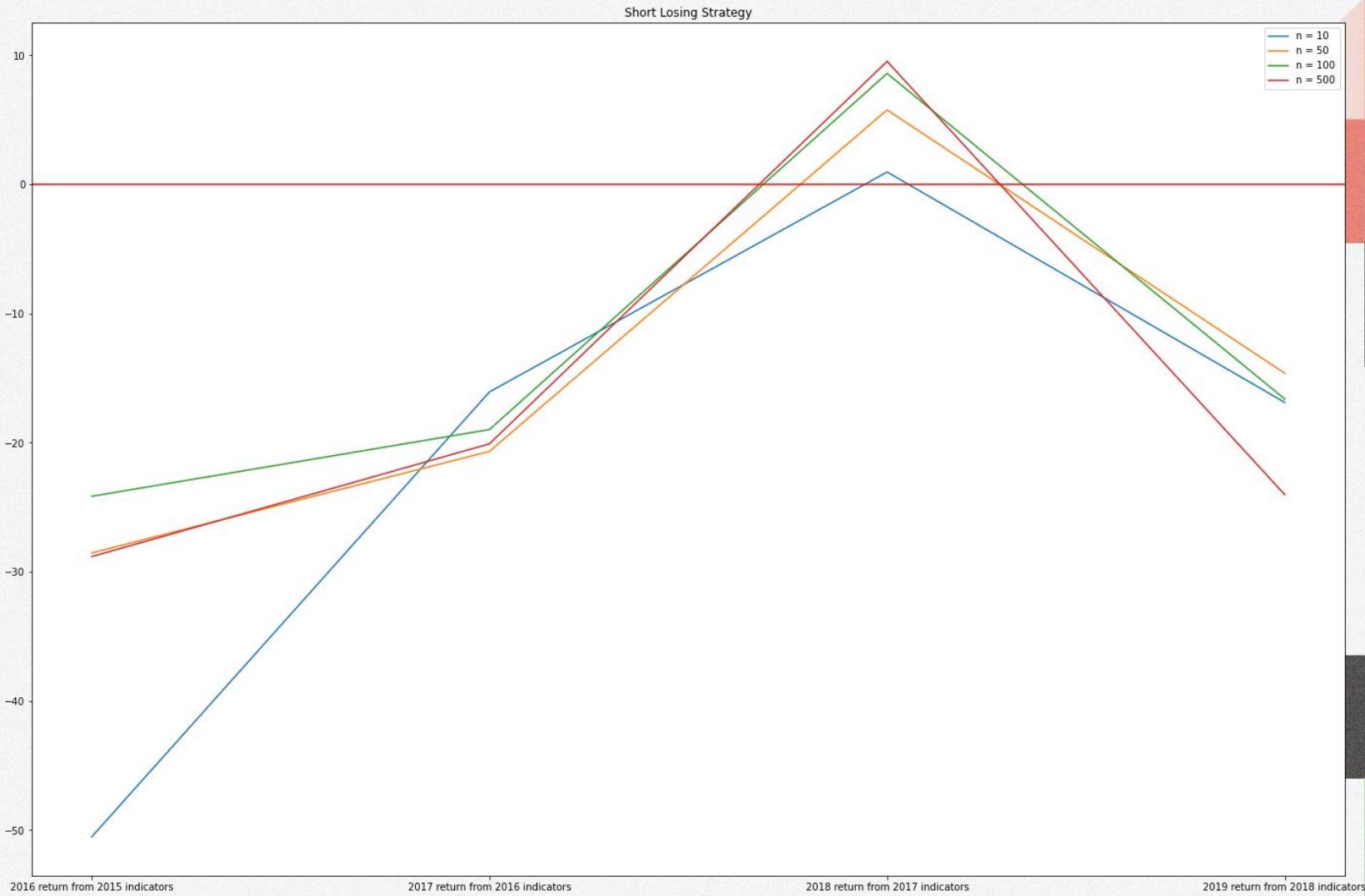
# Research topic

Find the top factors that affect the performance of a stock, so as to develop a safe stock investing strategy to ensure positive rates of return for investors while minimizing losses

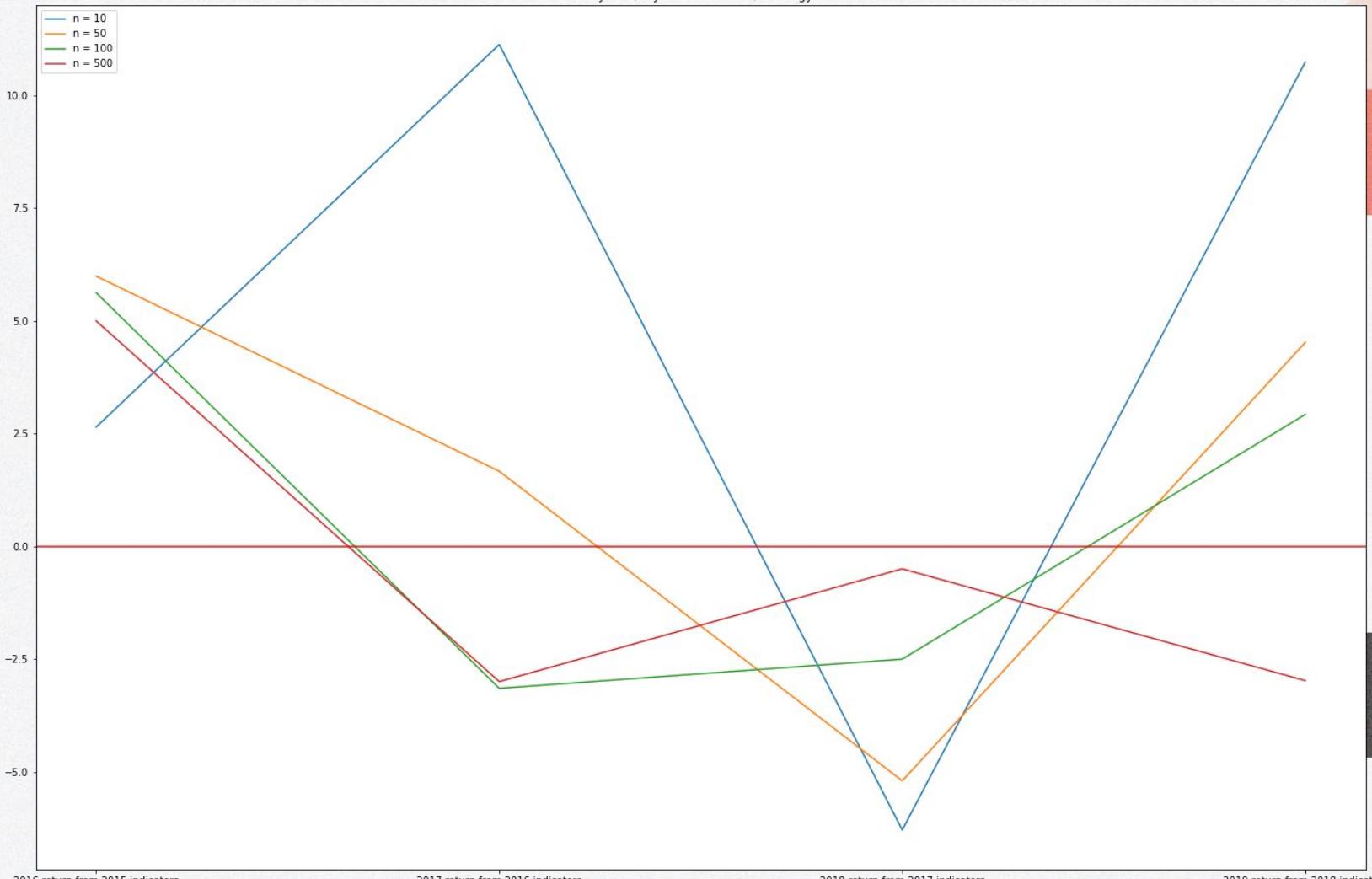
# Investing strategies

- 01** Buy top n stocks with highest Class “1” probability
- 02** Short sell top n stocks with lowest Class “1” probability
- 03** Buy top  $n/2$  stocks with highest Class “1” probability and Short sell top  $n/2$  stocks with lowest Class “1” probability

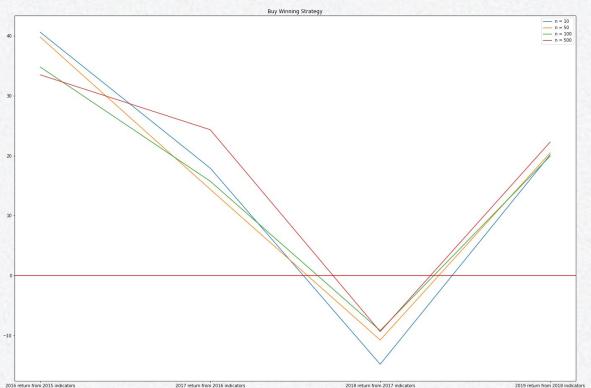




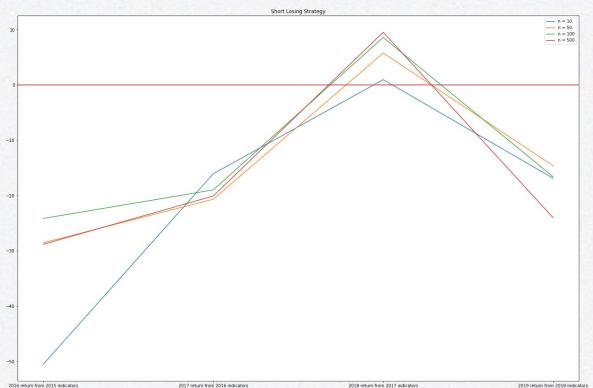
### Hybrid (Buy Win Short Lose) Strategy



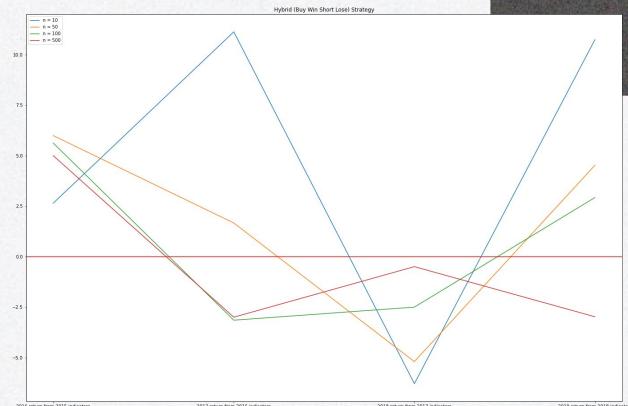
## Strategy 1: Buy winning



## Strategy 2: Short sell



## Strategy 3: Combination



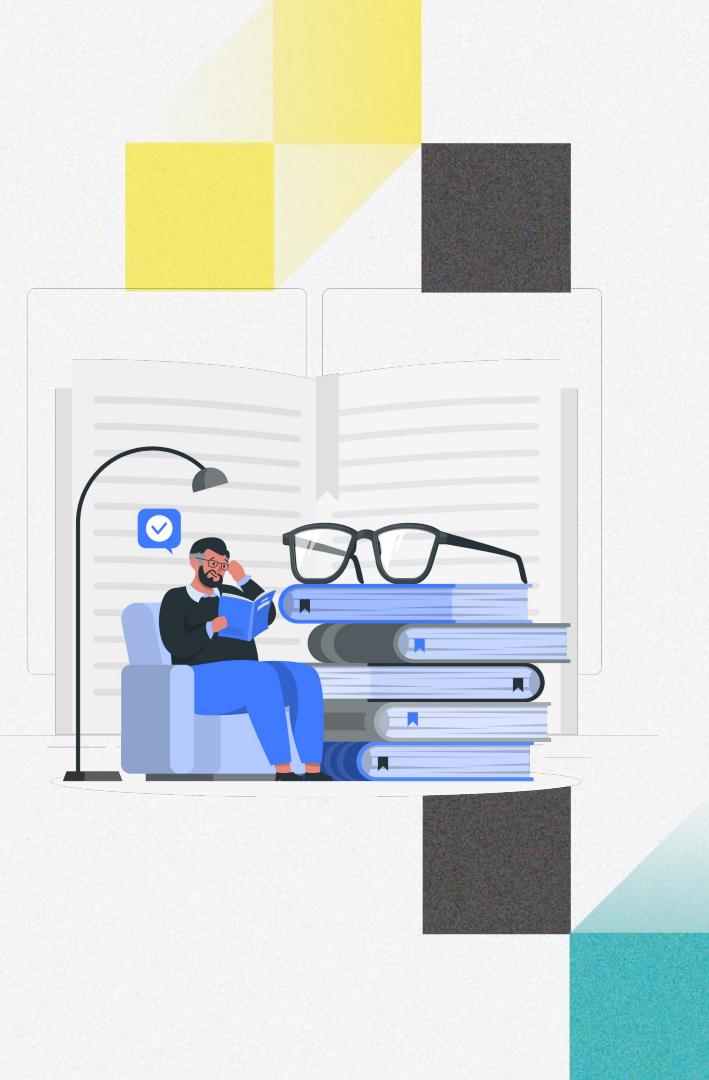
# Learning something new

- Cross Validation GridSearch
- Time series

## Outcomes

With our machine learning model + strategy investors can:

- Make informed decisions on which stocks to invest in
- Maximize profit whilst minimizing losses



# Assumptions

- Based on US stocks
- Long term investments

# Limitations

- Data used is from 2014-2018 (model might change if there is a drastic change in the market)
- Dataset was too large to pick best predictors
- Focus on maximizing stock price change and neglect stock dividends



# Thank you!

CREDITS: This presentation template was created by **Slidesgo**,  
including icons by **Flaticon**, and infographics & images by **Freepik**

Please keep this slide for attribution