
Feature Selection in Finance, It is Delicious

Benjamin A. Schiffman

Justin J. Siekmann

Department of Electrical and Computer Engineering

University of Arizona

Tucson, AZ 85719

bschifman@email.arizona.edu

jsiekman@email.arizona.edu

Abstract

The abstract paragraph should be indented ½ inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Introduction

There exist technical analysis indicators traditionally used by analysts to evaluate and predict market and equity performance, as they “can provide a unique perspective on the strength and direction of the underlying price action” (<–what is this quote from??).

Feature selection is used to determine relevant indicators while identifying those that are irrelevant and redundant. Different implementations of algorithms based on these indicators could be used to predict performance of individual equities, sectors, or overall markets. They could also be used to classify and identify the correlation and interdependencies between equities, sectors, and markets. Our goal is to implement various approaches to determine efficacy of technical indicators as enablers to financial analysis. This project presented a couple challenges during implementation including developing an accurate testing method as well as handling and computing such large volumes of data.

Possibly reword and keep/move:

From this project we hope to deepen our understanding of the usage cases for applying specific machine learning algorithms as well as expanding upon our technical analysis of the stock market and which indicators play a role in successful market analysis.

2 Related Work

This is optional. If wanted, save til last.

3 Methods/Approach

The following subsections present details and explanations of the methods and functions implemented as part of this project.

3.1 Data and Technical Analysis Indicators

The Quandl platform was used to fetch 11 years of market data from Dec 31, 2006 to Dec 31, 2017 on various identified US tickers across different sectors. Tickers used in the project can be found in Table 1 categorized by sector. As the amount of fetched and calculated data is very large, pickle files

were used to save all data to be quickly reimported instead of refetching data through the Quandl platform. TA-Lib: Technical Analysis Library was used to calculate features on the market data for each ticker obtained using Quandl. <insert stuff about TA-lib having various types of features>. The features incorporated into this project are found with in Table 2.

Financials	Utilities	Energy	Healthcare	Technology	Real Estate
JPM	T	XOM	JNJ	AAPL	ECL
BAC	VZ	CVX	UNH	GOOGL	FMC
WFC	NEE	BP	PFE	MSFT	IP
C	TMUS	GE	MRK	FB	VMC
MS		SLB	ABBV	INTC	BMS
			MMM	CSCO	
			AMGN	ORCL	
			MDT	IBM	
				NVDA	

Table 1: Tickers

Overlap	Momentum	Volume	Cycle	Price	Volatility	Statistical	Math Transformation	Pattern Recognition
SMA EMA BBAND HT_TRENDLINE SAR	RSI MOM ROC ROCP STOCH_SLOWK STOCH_SLOWD MACD MACDSIGNAL MACDHIST	OBV AD ADOSC	HT_DCPERIOD HT_TRENDMODE	AVGPRICE TYPPRICE	ATR	BETA LINEARREG VAR	EXP LN	CDLENGULFING CDLDOJI CDLHAMMER CDLHANGINGMAN

Table 2: Technical Analysis Indicators

3.2 Normalization

As each technical analysis indicator produces values applicable based on the way the indicator was calculated, normalization of the indicators makes correlations between them during feature selection more accurate and applicable. Each value in a specified time period is normalized using Equation (1). If there is extra data not consisting of a full time period, the extras are thrown out at the beginning of the data as data near the end may be more relevant and thus more desirable to keep. The start indices are computed for each ticker and return to ease future handling.

$$x_n = \frac{x - \min}{\max - \min} \quad (1)$$

3.3 Feature Correlation

Removing highly correlated features allows for the optimization of the classification algorithms by reducing the feature space. Features that are highly correlated most likely offer no additional data and they are an extra expense in computation time. The pairwise correlation of columns was computed, and columns that were correlated above a certain threshold were marked to be removed from future classification algorithms. The following Figure 1

3.4 Maximal Information Coefficient (MIC)

The MIC is "a measure of two-variable dependence designed specifically for rapid exploration of many-dimensional data sets" - <http://www.exploredata.net/>. A benefit to MIC correlations between two variables is that it can be described regardless of linear or non-linear relationships. The MIC yields a single value $0 \leq MIC \leq 1$ with a value closer to 1 representing that the variables are more closely correlated, and a value near 0 indicates statistically independent variables that have neither linear nor nonlinear relationships. The *minepy* library was used in python to rank the features according to their MIC with the target variable. The MIC was calculated for each feature in each ticker, and then a final MIC value for each feature was calculated by taking the mean of the values.

ADD TABLE OF MIC RESULTS!

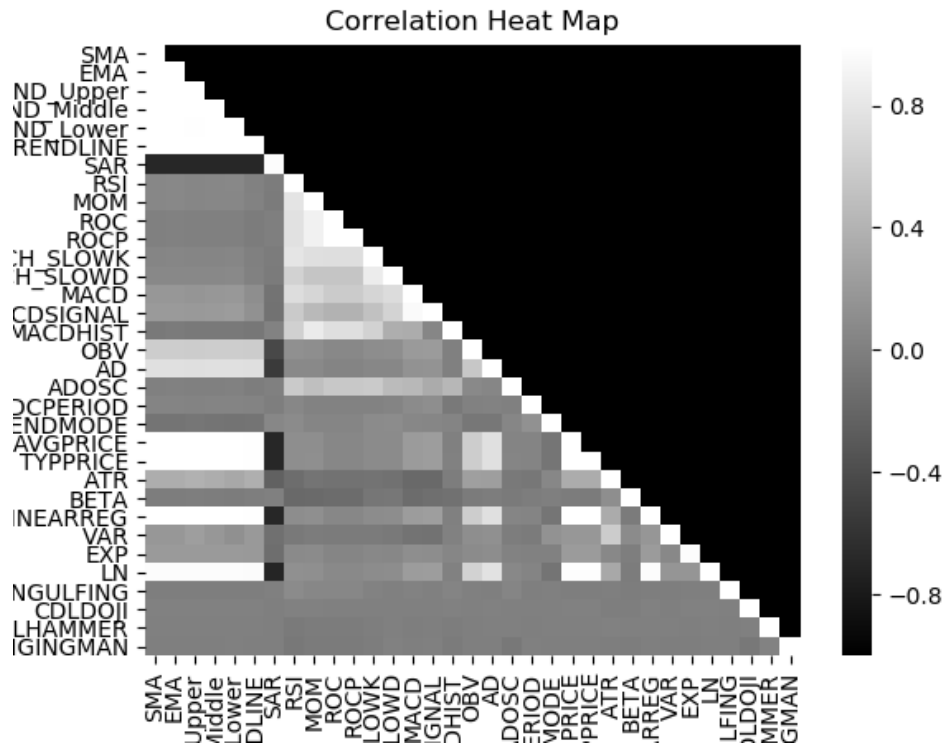


Figure 1: Heat Map of Correlated Variables

3.5 Recursive Feature Elimination (RFE)

Given an external estimator

3.6 Random Forest Classifier (RFC)

The Random Forest Classifier is an ensemble technique implementing decision trees and

3.7 Principle Component Analysis (PCA)

PCA orthogonally transforms a set of features into a set of linearly uncorrelated principal components. PCA is a method for reducing the dimensionality of the feature set size while retaining principal component variance, and the features informational relevance. To analyze the performance of this method, PCA was implemented on the original 33 features and then the resulting principal components were used in the RFC classifier with 5 fold cross validation to see the resultant accuracy. The number of principal components were iterated from 1 to 33 and then the principal component were implemented in the RFC to find the optimal number of principal components to yield the best accuracy and in the least amount of time.

4 Results

5 Conclusion

References

Principal Components	Score	Time
1	0.49	37.58
2	0.5	14.46
3	0.51	10.86
4	0.51	16.85
5	0.54	12.96
6	0.55	13.1
7	0.55	14.1
8	0.55	12.86
9	0.56	18.28
10	0.57	15.71
11	0.57	15.6
12	0.59	17.16
13	0.59	17.79
14	0.59	18.22
15	0.59	17.86
16	0.63	24.15
17	0.63	21.67
18	0.67	21.5
19	0.7	21.54
20	0.7	21.89
21	0.7	22.01
22	0.71	20.63
23	0.7	20.89
24	0.7	20.48
25	0.71	25.09
26	0.7	25.25
27	0.7	23.96
28	0.7	24.81
29	0.7	24.7
30	0.7	24.05
31	0.71	25.62
32	0.71	23.61

Table 3: Principal Component Analysis