

MACHINE LEARNING AND INFORMATION-THEORETIC APPROACHES FOR  
FINANCIAL APPLICATIONS

BY

KELECHI MOISE IKEGWU

DISSERTATION

Submitted in partial fulfillment of the requirements  
for the degree of Doctor of Philosophy in Informatics  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2021

Urbana, Illinois

Doctoral Committee:

Professor Robert J. Brunner, Chair, Director of Research  
Professor Theodore Sougiannis  
Assistant Professor Jeff McMullin  
Assistant Professor Joseph Yun

# Abstract

Firm-specific information contributes to resource allocation in capital markets. The prediction and inference of future firm-specific information are often critical factors for a prosperous economy. Machine Learning and Information Theory offer alternative methods to aid in the prediction and inference of future firm-specific information. Thus we explore approaches from machine learning and information theory to improve existing research areas in finance.

In Chapter 2 we introduce an agnostic framework called SML that integrates a query-like language to simplify the development of machine learning pipelines. We show that SML significantly reduces the code one has to write to develop parts of a machine learning pipeline. In Chapter 3 we utilized random forest trees to predict the annual direction of profitability for firms with a minimal amount of information. We demonstrate that we can out-perform benchmark models typically used in financial accounting research settings.

Major public firms announce their annual earnings during the first quarter of the year. The embedding information in these announcements effects the announcing firm and is transferred or incorporated in other firms. Existing literature in finance and accounting focuses on measuring information transfers as effects stemming from firm-specific information releases. In Chapter 4 we discuss a solution to estimate information transfer via transfer entropy between random processes. We show that our solution to estimate information transfer scales better than and is up to 1,072 times faster than all open source existing solutions for large datasets. In Chapter 5 we present an alternative approach to estimate information transfer between firms centered around earnings announcements. We show that information transfer between firms is stronger for firms on days with unexpected earnings

news. We also show that communities of firms are not adequately captured by characteristics that are the focus of existing literature.

*To my Father, Mother, Joshua, Sandra, Ezinne, Natalia, Khamille, and Khalil.*

# Acknowledgments

First and foremost, I would like to thank my advisor, Robert Brunner, for his guidance, patience, and support during my graduate studies at the University of Illinois at Urbana-Champaign. I also wish to thank my committee members, Theodore Sougiannis, Jeff McMullin, and Joseph Yun, for their help and support. I am also grateful to Junhyung Kim, Vincent Reverdy, William Biscarri, Samantha Thrush, Matias Carrasco Kind, Jacob Trauger, Alice Perng, and Howard Hardiman.

— Kelechi Moise Ikegwu

This work was partially funded by the Graduate College Fellowship program at the University of Illinois. The work was partially funded by the National GEM Consortium. This work utilizes resources provided by the Innovative Systems Laboratory at the National Center for Supercomputing Applications at the University of Illinois at Urbana-Champaign.

# Table of Contents

<b>Chapter 1 Introduction</b> . . . . .	<b>1</b>
1.1 Machine Learning . . . . .	1
1.2 Network Science . . . . .	4
1.3 Information Transfer in Financial Markets . . . . .	8
1.4 Figures and Tables . . . . .	12
<b>Chapter 2 Standard Machine Learning Language</b> . . . . .	<b>20</b>
2.1 Introduction . . . . .	20
2.2 Prior Works . . . . .	21
2.3 Grammar . . . . .	22
2.4 SML's Architecture . . . . .	26
2.5 Interface . . . . .	26
2.6 Use Cases . . . . .	27
2.7 Future Work . . . . .	29
2.8 Conclusion . . . . .	29
2.9 Figures and Listings . . . . .	31
<b>Chapter 3 Directional Profitability Predictions with Random Forests</b> . .	<b>39</b>
3.1 Introduction . . . . .	39
3.2 Data . . . . .	40
3.3 Methods . . . . .	42
3.4 Random Forest Implementation to Predict Directional Profitability . .	45
3.5 Random Walk Benchmark Models . . . . .	45
3.6 Comparative Analyses . . . . .	48
3.7 Results . . . . .	50
3.8 Conclusion and Future Work . . . . .	53
3.9 Figures and Tables . . . . .	55
<b>Chapter 4 Information Transfer Estimation on Large Data</b> . . . . .	<b>73</b>
4.1 Introduction . . . . .	73
4.2 Estimating Transfer Entropy . . . . .	74
4.3 PyIF . . . . .	77
4.4 Comparative Analysis . . . . .	78
4.5 Results . . . . .	81
4.6 Conclusion . . . . .	82

4.7 Figures and Tables . . . . .	83
<b>Chapter 5 Cross-Firm Information Transfers During Earnings Season: A Network Approach . . . . .</b>	<b>89</b>
5.1 Introduction . . . . .	89
5.2 Data . . . . .	90
5.3 Estimating Dynamic Information Transfer Between Firms . . . . .	91
5.4 Exploratory Data Analysis . . . . .	95
5.5 Earnings Surprise as a Determinant of Dynamic Information Transfers . . . . .	97
5.6 Community Detection . . . . .	99
5.7 Conclusion . . . . .	101
5.8 Figures and Tables . . . . .	103
<b>Chapter 6 Conclusion . . . . .</b>	<b>118</b>
6.1 Summary and Conclusion . . . . .	118
<b>Appendix A Standard Machine Learning Language Supplemental Code</b>	<b>127</b>
A.1 Iris Python Code . . . . .	127
A.2 Auto-MPG Python Code . . . . .	130

# Chapter 1

## Introduction

The economic goal of a society is to maximize wealth by allocating resources efficiently. Capital markets (such as the stock market) provide a method to help allocate resources efficiently. The prediction and inference of future financial events are often critical to achieving the economic goal. Firm-specific information contributes to resource allocation (i.e., labor, capital, and natural resources) in capital markets. The prediction of a future state for a firm (such as its profitability, financial condition, or revenue growth) is an essential factor contributing to resource allocation. There are many methods used to predict or inform about the future state of a firm.

### 1.1 Machine Learning

Machine Learning encompasses a vast set of statistical tools for understanding data and making predictions and inferences. These tools fall into supervised or unsupervised categories. For data-driven problems, typically, we have a response variable (or labels) referred to as  $Y$ , and  $p$  independent variables (or features) referred to as  $X$ . Given the assumption that there is some relationship between  $Y$  and  $X$  the relationship can be expressed as:

$$Y = f(X) + \epsilon \tag{1.1}$$

In eq. 1.1  $f$  is an unknown function that models the relationship between  $X$  and  $Y$  with  $\epsilon$  noise. In most cases it is not possible to perfectly model  $f$  with real-world data and  $f$  is estimated ( $\hat{f}$ ). In eq. 1.1 if  $Y$  is a quantitative variable then regression techniques are used

estimate  $f$ , however if  $Y$  is a qualitative variable then classification techniques are used.  $\hat{f}$  is sought for inference where one wants to estimate  $f$  to understand how  $\hat{Y}$  (predicted labels) changes as one adjusts features.  $\hat{f}$  is also sought for prediction where the only concern is the accuracy of predictions for  $Y$ . Machine learning provides a variety of approaches to estimate  $f$  for prediction and inference (see James et al. (2013)).

Machine learning problems typically fall into either supervised learning, semi-supervised learning, or unsupervised learning. With supervised learning, one wants to fit a model on labels using a set of features. Semi-supervised learning consists of an environment with features and some corresponding labels. Typically, in this setting, one wishes to use the given features and partial labels to determine the remaining missing labels. Lastly, with unsupervised learning, one has features and no associated labels and seeks to find structure between the features. Currently, the research in this dissertation will only utilize supervised learning.

### 1.1.1 Supervised Learning

In a supervised learning paradigm, it is not enough to only have  $\hat{f}$  estimate  $f$  well on a dataset.  $\hat{f}$  is often sought to learn a dataset well enough to generalize well to new data that  $\hat{f}$  has not seen before. To find a suitable  $\hat{f}$ , a standard methodology is to partition the data into a training set and testing set (see figure 1.5). The training set consists of data to train the model. The testing set consists of data to assess  $\hat{f}$ 's ability to generalize well to new data. The researcher/practitioner typically determines the ratio of data used to train and assess the model. A specific <sup>1</sup> way to estimate  $\hat{f}$ 's performance is to use an error metric such as mean squared error (MSE) given by:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad (1.2)$$

---

<sup>1</sup>One can use many error metrics; however, for the sake of brevity MSE is the only error metric discussed.

In eq. 1.2  $y_i$  and  $x_i$  are the  $i^{th}$  observations of the labels and features respectively.  $\hat{f}(x_i)$  is the prediction of  $\hat{f}$  for the  $i^{th}$  observation and  $y_i$  is the label for the  $i^{th}$  observations. Training MSE indicates the mean error of  $\hat{f}$ 's ability to learn the training data. Subsequently, there is a testing MSE associated with the mean error of  $\hat{f}$  on withheld testing data. A model with a small testing MSE generalizes well to the withheld data. Often, the training MSE and testing MSE differ where the training MSE underestimates the MSE for the testing MSE (see chapter 5 of James et al. (2013)). Many approaches can produce a more accurate estimate of how well  $\hat{f}$  will generalize to new data.

### 1.1.2 Validation Approaches

An alternative approach to training and assessing  $\hat{f}$  on training and testing sets, respectively, is to create a validation dataset by randomly sampling the training dataset. The researcher/practitioner determines the sampling method.  $\hat{f}$  trains on the training data, and the trained model assesses its ability to generalize with validation data. The validation MSE is more indicative of the test MSE. However, there are issues with this approach. For example, randomly sampled data is in the validation dataset. Given that the validation set is composed of random data, the validation MSE can be highly variable. Cross-validation can decrease the variability of the validation MSE.

A common cross-validation technique is  $k$ -fold cross-validation where the training data is split into  $k$  separate folds (see figure 1.6). The first fold is withheld and used as the validation dataset to assess the model's performance. The remaining  $k - 1$  folds train the model. For the next iteration, the same procedure repeats where the 2nd fold serves as the validation dataset to assess the model (see figure 1.7) and the other  $k - 1$  folds train the model. This process repeats  $k$  times and produces  $k$  MSE scores ( $MSE_1, MSE_2, \dots, MSE_k$ ). Averaging the MSE fold values yields the  $k$ -fold CV estimate:

$$CV = \frac{1}{k} \sum_{i=1}^k MSE_i \quad (1.3)$$

where  $k$  is the number of folds.

Selecting a large value  $k$  can be computationally expensive as the model must be retrained on  $k$  different folds. In addition to this, the selection of  $k$  can affect the model's ability to generalize to new data well. Typically with a smaller value of  $k$  you have more bias, and a larger value of  $k$  has more variance. For example consider when  $k = 3$  and  $k = 10$ , the training folds will have  $n_{train}$  observations where  $n_{train} = \frac{(k-1)n}{k}$  and the validation fold will have  $n_{validation}$  observations where  $n_{validation} = \frac{k-(k-1)n}{k}$ . When  $k = 10$ , 90% of the data is used to train the model, and the remaining 10% assesses the model. The model trains on very similar datasets, and the predictions from the model on the cross-validation datasets are highly correlated. The high correlated datasets can result in a model with high variance. If  $k = 3$  66.7% of the data trains the model and the remaining 33.3% of the data assesses the model, CV averages fewer predictions that are less correlated with each other since the overlap between the training sets are smaller. When  $k = 3$  this results in a CV score with lower variance and higher bias.

## 1.2 Network Science

Network science is the study of networks, and a network (or graph) consists of a collection of nodes joined by connections (commonly referred to as edges). Networks capture interactions between nodes, and over the years, scientists have discovered that specific patterns of interactions have a significant effect on the network's behavior (see Newman (2010)). This field has evolved, and a vast set of tools have developed to analyze, model, and understand networks (again, see Newman (2010)). Researchers have used networks in financial markets to determine which sectors are related around specific events. For example, Billio et al. (2012) and Diebold and Yilmaz (2011) found that banks play an essential role in transmit-

ting shocks relative to the remaining set of financial firms, consistent with banks playing a more central role in systemic risk for the economy around the 2007-2009 financial crisis. Network science can help us understand more about the probable future states of firms.

Networks can be formed from data to offer an alternative representation of the data, allowing many analysis techniques. Network Theory involves the study of networks as a representation of relationships between nodes. Throughout this section,  $G$  denotes a network,  $m$  denotes the number of edges in a network, and  $n$  denotes the number of nodes. A graph with  $n$  nodes and  $m$  edges is denoted as  $G(n, m)$ .

Undirected graphs do not have to specify the direction of an edge. For example, if Node  $A$  and Node  $B$  are connected, an undirected graph will not distinguish whether Node "A" connects to Node "B" or vice-versa; an undirected graph will only retain information about nodes  $A$  and  $B$  having a connection. Figure 1.1 shows an example of a undirected graph  $G(7, 16)$ . The amount of edges a particular node has is the degree of the node. Figure 1.1, the node "FB" has a degree of 6 because it connects to 6 other nodes. Degree is a standard metric used in network theory.

It is also possible for a graph to retain additional information about edges. Another source of information is the direction of the edges; graphs that store information about edges' directions are called directed graphs. For example, consider a directed graph with two nodes (Node "A" and Node "B"). If Node "A" connects to Node "B" a directed graph would contain information that Node "A" has an edge to Node "B" but Node "B" does not have an edge to Node "A". Figure 1.2 shows an example of a directed graph; in figure 1.2 there is an edge from "JNJ" to "FB" however there is no edge from "FB" to "JNJ".

Weighted networks (or weighted graphs) consider edge weights. Instead of having a binary connection to determine whether an edge is present or not, it can measure the strength between nodes' connections. For example, if one has an attribute to determine the strength between nodes, such as a statistic from the data, it can be used as the weight value of an edge to determine the strength of the connection between nodes. Edge weights are applicable

to both directed and undirected graphs. Ergo it is possible to have an undirected weighted graph or directed weighted graph. Figure 1.3 contains an example of a directed weighted network.

Under certain circumstances, it is more convenient to express a graph as an adjacency matrix. The graph in Figure 1.1, is converted to a table with column and row headers table 1.1. However it is mathematically more convenient to express table 1.1 as an adjacency matrix (see Eq 1.4).

$$A = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 \end{bmatrix} \quad (1.4)$$

Similar to Table 1.1 each matrix element ( $A_{ij}$ ) in  $A$  is assigned a 1 or 0 based on if there is a connection between a row index ( $i$ ) and column index  $j$ . You'll notice in Eq 1.4 that the matrix is symmetric because the direction of edges are not taken into account; ergo if node  $i$  and node  $j$  are connected  $A_{ij}$  and  $A_{ji}$  have a value of 1. It is also possible to represent directed graphs by assigning a 1 to a matrix element  $A_{ij}$  where node  $i$  is connected to node  $j$ . Lastly you can express the weight of edges in adjacency matrices by assigning non-binary values to matrix elements.

### 1.2.1 Analysis of Networks

After forming a network from data, common metrics are often sought to help describe the structure of the network. The first metric are the degree of nodes where the degree can be defined in eq 1.5 as the sum of edges for a particular node. In Figure 1.4 the node sizes are scaled based on the degree.

$$m_i = \sum_{ij} A_{ij} \quad (1.5)$$

where  $A$  is an adjacency matrix. Subsequently, there are degree measures for directed networks, such as in-degree, which counts the number of edges connected to node  $i$  (see Eq 1.6). Whereas out-degree counts the amount of edges that node  $i$  is connected to (see Eq 1.7).

$$m_i^{in} = \sum_{j=1}^n A_{ij} \quad (1.6)$$

$$m_i^{out} = \sum_{i=1}^n A_{ij} \quad (1.7)$$

The network's density provides a numerical value for the proportion of actual connections out of all potential connections in a network. If all potential connections in a network are present, then the density will have a numerical value of 1, and if there are no connections in an undirected network, the density will be 0. The maximum number of potential connections in a network is  $n(n - 1)/2$ . The degree summed for all nodes is  $m$  (for directed networks, the summed in-degree and out-degree for all nodes is  $m$ ). The density for a graph is:

$$D = \frac{2m}{n(n - 1)} \quad (1.8)$$

Centrality is another metric often sought to measure the importance of nodes in networks. Measuring the degree of nodes in relation to other nodes is called degree centrality (see Eq

1.9).

$$C_i^d = m_i = \sum_{ij} A_{ij} \quad (1.9)$$

## 1.3 Information Transfer in Financial Markets

The price discovery process is also vital to achieve the economic goal of a society. In particular, how information transfers from one price to another is useful in the price discovery process. Information is embedded in price, and an efficient market changes rapidly (if not instantaneously) when new information arrives. For example, Apple Inc. recently announced that their new line of computers (Macs) would use central processing units (CPUs) developed in-house and transition away from Intel CPUs. This information will reflect in Apple's stock price and Intel's stock price in ideal settings.

In the Apple and Intel example, it is probable that additional information reflects in Intel's price change at a given moment. However, it is difficult to determine price changes for specific events. The entire set of information that reflects in the price change at a given time is unknown. Prior studies have examined the price changes over long windows such as days or weeks (see Foster (1981) and Baginski (1987)). In particular, Foster (1981) measured information transfers from firms that announce earnings and non-announcement firms in the same industry based on short-window price reactions. Subsequent studies use a similar approach to examine short-window price reactions to firm-level earnings announcements (see Clinch and Sinclair (1987), Pownall and Waymire (1989), Han and Wild (1990), WANG (2014), and Hann et al. (2019)).

### 1.3.1 Information Theory

Information theory provides a framework for studying the quantification, storage, and communication of information (see Stone (2015)). Information theory has been useful in the

price discovery process (see Billio et al. (2012) or Diebold and Yilmaz (2011)). It offers an alternative approach for understanding information in price. In particular, it offers alternative methods to understand information transfer between random processes (such as equity prices).

In this dissertation, we measure information transfer in financial markets with the relatively new measure Transfer Entropy (TE) (see Schreiber (2000)). TE allows us to measure information transfers in a broader sense and avoids assumptions of prior works. We avoid the assumption of links between firms due to a shared characteristic (such as shared industry) or a known event (such as earnings announcements). TE quantifies the reduction in uncertainty in one random process from knowing past realizations of another random process.

For example, it is unlikely to predict Intel's closing price tomorrow with no information accurately. However, knowing Intel's price today reduces uncertainty in the prediction for its price tomorrow. In addition to this, knowing Apple's price today may help further reduce the uncertainty in Intel's future price predictions. If Apple's price today reduces the uncertainty further, TE will be positive, indicating an information transfer from Apple to Intel. The TE value will vary at any given time, given the information incorporated into Apple's and Intel's prices.

Schreiber (2000) discovered TE and coined the name "transfer entropy," although Paluš et al. (2001) also independently discovered the concept as well. To define TE, one must begin with the building blocks of information theory <sup>2</sup>. Shannon Information is:

$$h(x) = \log_2 \frac{1}{p(x)} \quad (1.10)$$

where  $x$  is an event and  $p(x)$  is probability of  $x$  occurring. Shannon's information definition implies that values of  $x$  with small probabilities contain more information. Consequently, values of  $x$  that are more probable (or common) contain less information. Entropy is the

---

<sup>2</sup>For simplicity, we use discrete events to introduce these concepts.

amount of uncertainty or disorder in a random process.

$$H(X) = \frac{1}{n} \sum_{i=1}^n \log_2 p(x_i) \quad (1.11)$$

In eq. 1.11 we define entropy as the weighted average of Shannon Information. Conditional entropy is the entropy after conditioning on another process ( $Y$ ):

$$H(X|Y) = \sum_{y \in \Omega_y} p(y) H(X|y) \quad (1.12)$$

In eq. 1.12  $y \in \Omega_y$  represents the occurrence of  $y$  in a set of possible events for  $y$  and  $H(X|y)$  represents conditional entropy on a single event  $y$  (or  $H(X|y) = \sum_{x \in \Omega_x} \frac{p(x|y)}{\log_2 p(x|y)}$ ). Mutual Information quantifies the amount of information shared across random variables:

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (1.13)$$

In eq. 1.13 consider  $H(X) - H(X|Y)$ ,  $H(X)$  computes the entropy of  $X$  and  $H(X|Y)$  computes the entropy of  $X$  conditioned on  $Y$ . The difference between  $H(X)$  and  $H(X|Y)$  captures the shared information between  $X$  and  $Y$ . Mutual Information is a symmetric measure ergo  $I(X, Y) = I(Y, X)$ .

Lagged mutual information  $I(X_t : Y_{t-k})$  can be used as a time-asymmetric measure of information transfer from  $Y$  to  $X$  where  $X$  and  $Y$  are both random processes,  $k$  is a lag period, and  $t$  is the current time period. However, lagged mutual information is unsatisfactory as it does not account for a shared history between the processes  $X$  and  $Y$  (see Kaiser and Schreiber (2002)). While similar to lagged mutual information, TE also considers the dynamics of information and how these dynamics evolve in time (see Schreiber (2000)). TE is also an asymmetric measure of information transfer. Ergo, TE computed from process  $A$  to process  $B$  may yield a different result than TE computed from  $B$  to  $A$ . The information-theoretic framework and these measures have led to various applications in

different research areas (see Bossomaier et al. (2016)).

TE considers the shared history between two processes via conditional mutual information. Specifically, TE conditions on the past of  $X_t$  to remove any redundant or shared information between  $X_t$  and its past. This also removes any information in the process  $Y$  about  $X$  at time  $t$  that is in the past of  $X$  (see Williams and Beer (2011)). Transfer entropy  $T$  (where the transfer of information occurs from  $Y$  to  $X$ ) is:

$$T_{Y \rightarrow X}(t) \equiv I(X_t : Y_{t-k} | X_{t-k}) \quad (1.14)$$

Kraskov et al. (2004) shows that transfer entropy can be expressed as the difference between two conditional mutual information computations:

$$T_{Y \rightarrow X}(t) = I(X_t | X_{t-k}, Y_{t-k}) - I(X_t | X_{t-k}) \quad (1.15)$$

The intuition of this definition is that TE measures the amount of information in  $Y_{t-k}$  about  $X_t$  after considering the information in  $X_{t-k}$  about  $X_t$ . Put differently, TE quantifies the reduction in uncertainty about  $X_t$  from knowing  $Y_{t-k}$  after considering the reduction in uncertainty about  $X_t$  from knowing  $X_{t-k}$ .

## 1.4 Figures and Tables

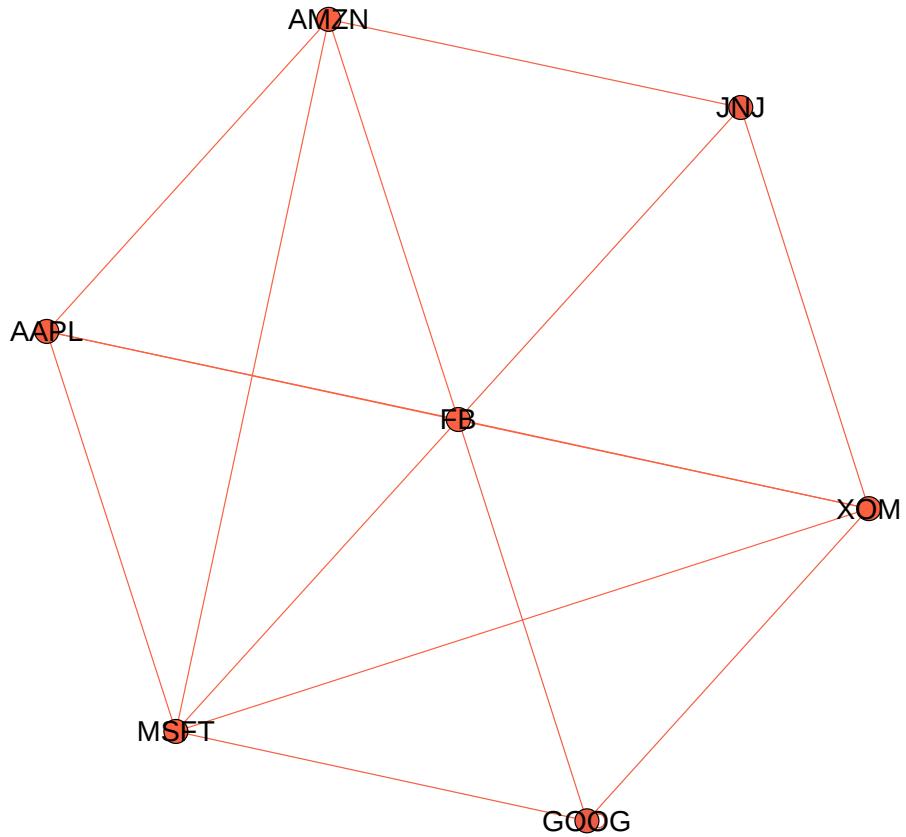


Figure 1.1: This figure contains an example of an undirected graph with 7 nodes and 16 edges. This graph is formed from randomly selecting stock symbol tickers from a set of tickers and forming connections randomly. The purpose of this data is solely to demonstrate the basics of network theory. The nodes are defined as circles and have ticker symbols overlaid on them. Edges are formed between nodes with lines between them. Since this is an undirected graph there is no directional information ergo no distinction as to whether one node is connected to another or vice-versa. The only information represented from an undirected graph is that they are connected.

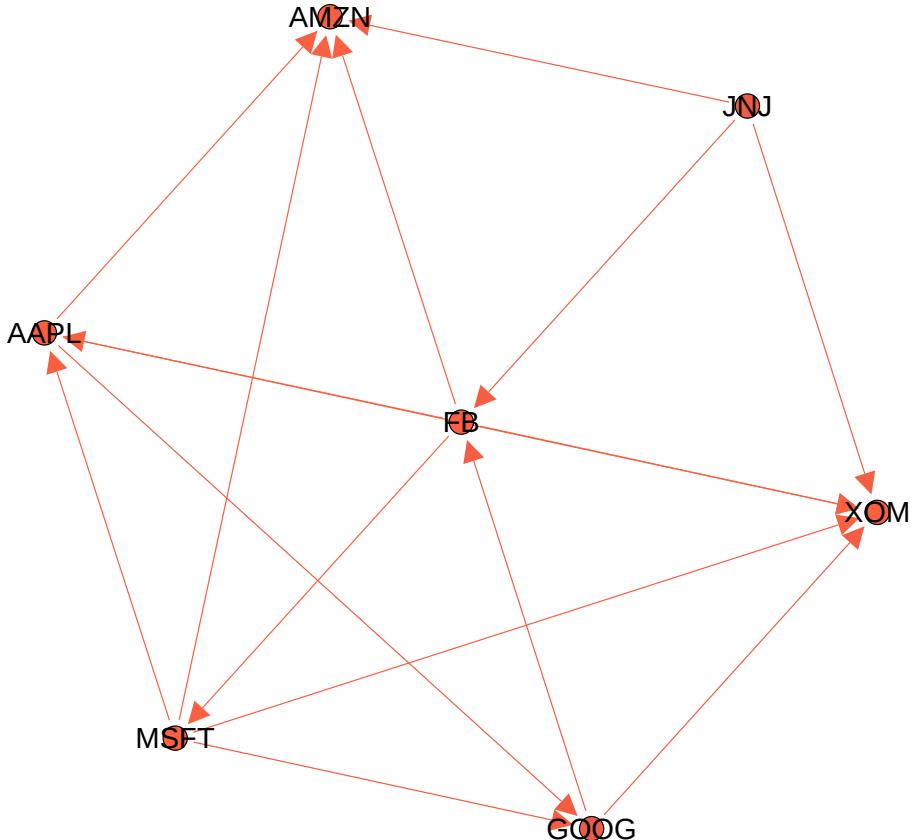


Figure 1.2: This graph is nearly identical to the graph in Figure 1.1 except this graph is a directed network. This means that it contains directional information between the nodes. This directional information is communicated based on the position of the arrow in a line. For example, the node *FB* has an edge to *MSFT* however *MSFT* does not have an edge to *FB*. The directional information allows other metrics to be used such as in-degree (which represents the amount of edges directed toward a node) and out-degree (which represents the amount of edges directed away from a node). In this example *FB* has an in-degree of 2, an out-degree of 4 and a degree of 6.

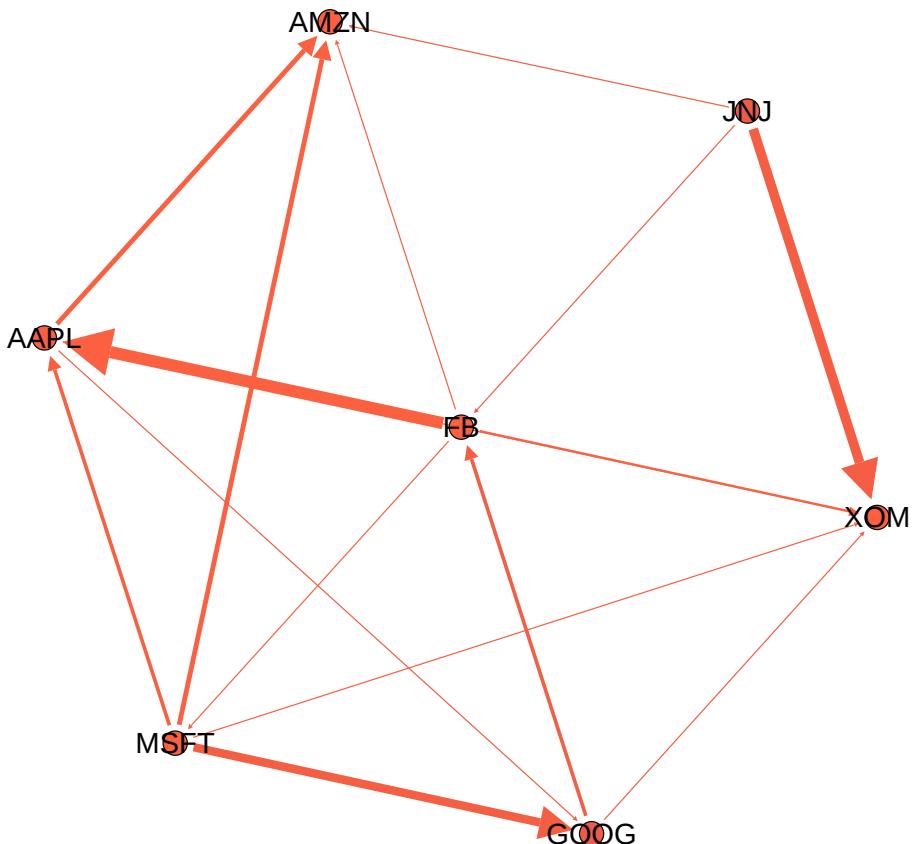


Figure 1.3: A nearly identical graph to the graph in Figure 1.2. This directed graph contains weighted edges where a thick edge represents a high weight value (or strong connection). The smaller the edge weight value the thinner the edge will be. Here *FB* has a strong connection to *AAPL* whereas *FB* has a weak connection to *AMZN*.

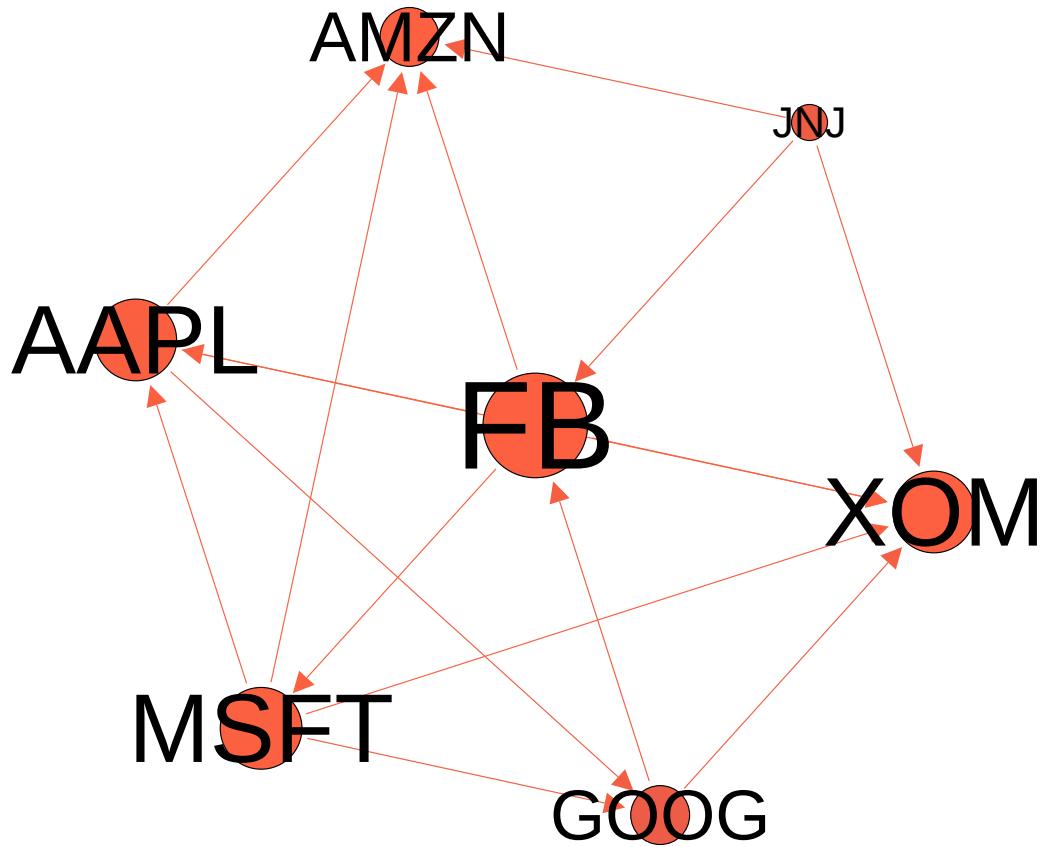


Figure 1.4: A nearly identical graph to the graph in Figure 1.2. Here the node sizes are scaled based on the degree. The higher the degree the larger the node size and consequently the smaller the degree the smaller the node size.



Figure 1.5: This figure shows how a data set can be partitioned for a supervised machine learning paradigm. The “data-set” is partitioned into 2 unequal sets with the “Train Set” having more data assigned to it than the “Test Set”. The “Train Set” and “Test Set” ratios can vary and is typically determined by the practitioner. There are many ways to partition the data for the train and test sets. For example the dataset can be partitioned sequentially from the data, or form partitions based on the data being randomly sampled.

## Train Set



Figure 1.6: The "Train Set" here represents the "Train Set" in Figure 1.5. In this figure the data is split into  $k$  folds. How the  $k$  folds are created can vary. A common technique is to randomly divide the "Train Set" observations into  $k$  equal folds.

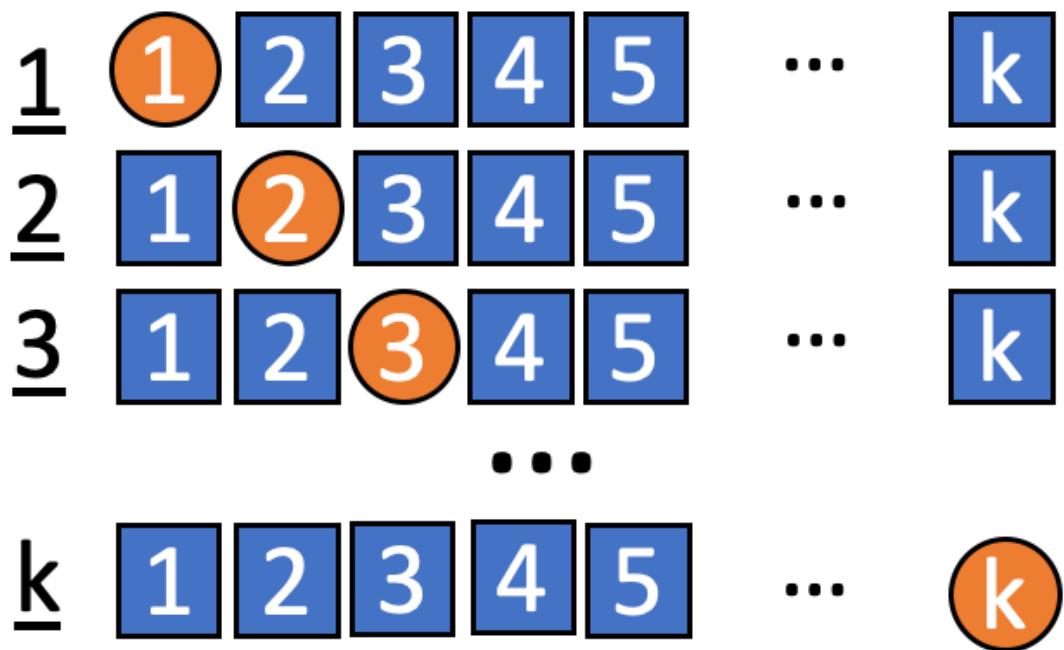


Figure 1.7: This figure breaks down how  $k$ -Fold Cross Validation works. These folds are created from “Train Set” in Figure 1.6. For the first iteration (where you see 1 underlined), the 1st fold (circle labeled 1) is used to assess the models ability and the remaining folds (squares) are used to train the model. The 2nd iteration follows the same procedure as the 1st except only the 2nd fold is withheld and the other folds are used. This repeats  $k$  times.

	MSFT	AAPL	AMZN	FB	JNJ	GOOG	XOM
MSFT	0	1	1	1	0	1	1
AAPL	1	0	1	1	0	1	1
AMZN	1	1	0	1	1	0	0
FB	1	1	1	0	1	1	1
JNJ	0	0	1	1	0	0	0
GOOG	1	1	0	1	0	0	1
XOM	1	1	0	1	0	1	0

Table 1.1: This table contains an example of simulated data. This table was formed from randomly selecting stock symbol tickers from a set of tickers. The purpose of this data is solely to demonstrate the basics of network theory. Here the columns and row indicies belong to the selected tickers. The values contain either a 1 to represent if there is a connection between the ticker at a particular row index  $i$  and the ticker of a particular column  $j$  or a 0 if there is no connection.

# Chapter 2

## Standard Machine Learning Language

### 2.1 Introduction

Machine Learning has simplified the process of solving problems in a variety of fields (see Bakshi and Bakshi (2018) or Monahan (2018) for examples). However, Domingos (2012) noted several nuisances to consider when developing machine learning pipelines. If one does not consider these nuisances into consideration, one may not receive satisfactory results. To combat these issues, we introduce Standard Machine Learning Language (SML), targeted at domain experts that want to utilize machine learning to solve research questions.

The overall objective of SML is to provide a level of abstraction which simplifies the development process of machine learning pipelines. Consequently, this will enable students, researchers, and industry professionals without a background in developing machine learning pipelines to solve problems in different domains with machine learning (see Listing 2.1 for an example). In the subsequent sections, related works are discussed, followed by defining the grammar used to create SML queries. The architecture of SML is described. Lastly, SML is applied to use-cases to demonstrate how it reduces the complexity of solving problems that utilize machine learning.

---

This chapter contains material from the following working paper:

- Kelechi Ikegwu, Micheal Hao, Nerraj Asthana, and Robert Brunner. Standard machine learning language: A language agnostic framework for streamlining the development of machine learning pipelines. 2017. URL [https://github.com/lcdm-uiuc/Publications/blob/master/2017\\_Kelechi\\_Mike\\_Brunner/main.pdf](https://github.com/lcdm-uiuc/Publications/blob/master/2017_Kelechi_Mike_Brunner/main.pdf)

## 2.2 Prior Works

There are prior related works that attempt to provide a level of abstraction for developing machine learning models. Rizzolo and Roth (2010) created a tool called LBJava based on a programming paradigm called Learning Based Programming (see Roth (2005)). Learning Based Programming is an extension of conventional programming that creates functions using data-driven approaches. LBJava utilizes machine learning to create these functions and abstract the details from the user. What differentiates SML from LBJava is that it offers a higher level of abstraction by providing a query-like language, allowing people who aren't experienced programmers to use SML.

TPOT (see Olson et al. (2016)) is a tool implemented in Python that creates and optimizes machine learning pipelines using genetic programming. Given cleaned data, TPOT preprocesses the data, performs feature selection, and the construction of machine learning models. Given the task (classification, regression, or clustering) TPOT uses genetic programming to tune model parameters and select features to determine the most optimal model to use. Similar to TPOT Kotthoff et al. developed Auto-WEKA, which automates the selection of learning algorithms and tuning hyper-parameters for implemented models in WEKA see (Frank et al. (2005)).

Subsequently Komer et al. created Hyperopt-Sklearn that provides automated algorithm selection from models in the Scikit-learn machine learning library<sup>1</sup> similar to Auto-Weka. Feurer et al. introduced improvements upon Hyperopt-Sklearn by taking into account past performance on similar datasets and constructing ensembles from optimized models. What differentiates SML from these prior works is that it provides an agnostic language to reduce the amount of programming required to write and offers a visualization framework to assess the models' performance visually.

---

<sup>1</sup>See Pedregosa et al. (2011) for an introduction to Scikit-learn.

## 2.3 Grammar

The SML language is a domain-specific language with grammar implemented in Bakus-Naur form (BNF) (see Backus (1959)). Each expression has a rule and can be expanded into additional terms. Listing 2.1 is an example of how one would perform classification on a dataset using SML. The query in listing 2.1 reads from a dataset, performs an 80/20 split of training and testing data, respectively, and performs classification on the 5th column of the hypothetical dataset using columns 1,2,3, and 4 as predictors. In the subsequent subsections, SML’s grammar in BNF form is defined in addition to the keywords.

### 2.3.1 Grammar Structure

This subsection is dedicated to defining the grammar of SML in terms of BNF. A *Query* can be defined by a delimited list of actions where the delimiter is an *AND* statement; with BNF syntax this is defined as:

$$< \text{Query} > ::= < \text{Action} > | < \text{Action} > \text{AND} < \text{Query} > \quad (2.1)$$

An *Action* in (2.1) follows one of the following structures defined in (2.2) where a *Keyword* is required followed by an *Argument* and/or *OptionList*.

$$\begin{aligned} &< \text{Action} > ::= < \text{Keyword} > < \text{Argument} > \\ &| < \text{Keyword} > < \text{Argument} > (< \text{OptionList} >) \\ &| < \text{Keyword} > (< \text{OptionList} >) \end{aligned} \quad (2.2)$$

A *Keyword* is a predefined term associating an *Action* with a particular string. An *Argument* generally is a single string surrounded by quotes that specifies a path to a file. Lastly, an *Argument* can have a multitude of options (2.3) where an *Option* consist of

an *OptionName* with either an *OptionValue* or *OptionValueList*. An *OptionName* and *OptionValue* consist of a single string. An *OptionList* (2.4) consist of a comma delimited list of options and an *OptionValueList* (2.5) consist of a comma delimited list of *OptionValues*.

$$\begin{aligned} < \text{Option} > ::= & < \text{OptionName} > = < \text{OptionValue} > \\ | & < \text{OptionName} > = [< \text{OptionValueList} >] \end{aligned} \tag{2.3}$$

$$< \text{OptionList} > ::= < \text{Option} > | < \text{Option} >, < \text{OptionList} > \tag{2.4}$$

$$\begin{aligned} < \text{OptionValueList} > ::= & < \text{OptionValue} > \\ | & < \text{OptionValue} >, < \text{OptionValueList} > \end{aligned} \tag{2.5}$$

To put the grammar into perspective the example *Query* in Listing 2.1 has been transcribed into BNF format and can be found in Listing 2.2. In Listing 2.2 the first *Keyword* is *READ* followed by an *Arugment* that specifies the path to the dataset. Next, an *OptionValueList* contains information about the delimiter of the dataset and the header. We include the *AND* delimiter to specify an additional *Keyword SPLIT* with an *OptionValueList* that tells us the size of the training and testing partitions for the dataset specified with the *READ Keyword*. Lastly, we use the *AND* delimiter to specify another *Keyword CLASSIFY* which performs classification using the training and testing data from the result of the *SPLIT Keyword* followed by an *OptionValueList* which provides information to SML about the features to use (columns 1-4), the label we want to predict (column 5), and the algorithm to use for classification. The next subsection describes the functionality for all *Keywords* of SML.

### 2.3.2 Keywords

Currently, there are eight *Keywords* in SML<sup>2</sup>. These *Keywords* can be chained together to perform a variety of actions. In the subsequent subsections, we describe the functionality of each *Keyword*.

#### Reading Datasets

When reading data from SML one must use the *READ Keyword* followed by an *Argument* containing a path to the dataset. *READ* also accepts a variety of *Options*. The first *Query* in listing 2.3 consist of only a *Keyword* and *Argument*. This *Query* reads in data from "/path/to/dataset". The second *Query* includes an *OptionValueList* in addition to reading data from the specified path; the *OptionValueList* specifies that the dataset is delimited with semicolons and does not include a header row.

#### Cleaning Data

When NaNs, NAs, and other missing values are present in the dataset, we clean these values in SML by using the *REPLACE Keyword*. Listing 2.4 shows an example of the *REPLACE Keyword* being used. In this *Query* we use the *REPLACE Keyword* in conjunction with the *READ Keyword*. SML reads from a comma-delimited dataset with no header from the path "/path/to/dataset". Then we replace any instance of "NaN" with the mode of that column in the dataset.

#### Partitioning Datasets

It is often useful to split a dataset into training and testing datasets for most tasks involving machine learning. Splitting a dataset can be achieved in SML by using the *SPLIT Keyword*.

---

<sup>2</sup>Detailed documentation providing examples and describing all of the keywords of SML are publicly available on GitHub: <https://github.com/lcdm-uiuc/sml/tree/master/dataflows>

Listing 2.5 shows an example of a SML *Query* performing an 80/20 split for training and testing data respectively by utilizing the *SPLIT Keyword* after reading in data.

## Creating Models

In SML, one can create a model to either perform classification, regression, or clustering. To use a classification model in SML one would use the *CLASSIFY Keyword*. SML has the following classification models implemented: Support Vector Machines, Naive Bayes, Random Forest, Logistic Regression, and K-Nearest Neighbors. Listing 2.6 demonstrates how to use the *CLASSIFY Keyword* in a *Query*. Clustering models can be utilized by using the *CLUSTER Keyword*. SML currently has K-Means clustering implemented. Listing 2.7 demonstrates how to use the *CLUSTER Keyword* in a *Query*. Regression models use the *REGRESS Keyword*. SML currently has the following regression algorithms: Simple Linear Regression, Ridge Regression, Lasso Regression, and Elastic Net Regression. Listing 2.8 demonstrates how to use the *REGRESS Keyword* in a *Query*.

## Saving/Loading Models

It is possible to save models and reuse them later. To save a model in SML one would use the *SAVE Keyword* in a *Query*. To load an existing model from SML one would use the *LOAD Keyword* in a *Query*. Listing 2.9 shows the syntax required to save and load a model using SML. With any of the existing queries using *REGRESS*, *CLUSTER*, or *CLASSIFY Keywords* attaching *SAVE* to the *Query* will save the model.

## Visualizing Datasets and Metrics of Algorithms

When using SML it is possible to visualize datasets or performance of models (such as learning curves or ROC curves). To do this, the *PLOT Keyword* must be specified in a *Query*. Listing 2.10 shows an example of how to use the *PLOT Keyword* in a *Query*. We apply the same operations to perform clustering in Listing 2.7 however, we utilize the *PLOT*

*Keyword.*

## 2.4 SML's Architecture

With SML's grammar defined, we have presented enough information to explain SML's architecture. When SML receives a *Query* in the form of a string, it is passed to the parser. The grammar parses through the string to determine the actions to perform. The actions are stored in a dictionary and given to one of the following SML phases: Model Phase, Apply Phase, or Metrics Phase. Figure 2.1 shows a block diagram of this process.

The model phase is for constructing a model. The *Keywords* that generally invoke the model phase are: *READ*, *REPLACE*, *CLASSIFY*, *REGRESS*, *CLUSTER*, and *SAVE*. The apply phase is for applying a preexisting model to new data. The *Keyword* that invokes the apply phase is *LOAD*. It's often useful to visualize the data that one works with and beneficial to see a machine learning model's performance metrics. By default, if one specifies the *PLOT Keyword* in a *Query*, SML will execute the metrics phase.

The last significant component of SML's architecture is the connector. The connector connects drivers from different libraries and languages to achieve an action a user wants during a particular phase (see Figure 2.2). If one considers applying linear regression on a dataset, SML calls the connector to retrieve the linear regression library during the model phase. In this case, SML uses sci-kit learn's implementation. However, if we wanted to use an algorithm not available in sci-kit learn, such as a Hidden Markov Model (HMM), SML will use the connector to call another library that supports HMM.

## 2.5 Interface

There are multiple interfaces available for working with SML. We have developed an alpha version of a web tool that allows users to write queries and get results back from SML through

a web interface (see Figure 2.3). There is also a REPL environment available that allows the user to write queries and display results from the appropriate phases of SML interactively. Lastly, users can import SML into an existing pipeline to simplify the development process of applying machine learning to problems.

## 2.6 Use Cases

We tested SML’s framework against ten popular machine learning problems with publicly available data sets. We applied SML to the following datasets: Iris Dataset <sup>3</sup>, Auto-MPG Dataset <sup>4</sup>, Seeds Dataset <sup>5</sup>, Computer Hardware Dataset <sup>6</sup>, Boston Housing Dataset <sup>7</sup>, Wine Recognition Dataset <sup>8</sup>, US Census Dataset <sup>9</sup>, Chronic Kidney Disease <sup>10</sup>, Spam Detection <sup>11</sup> which were taken from UCI’s Machine Learning Repository (see Lichman (2013)). We also applied SML to the Titanic Dataset <sup>12</sup>.

As mentioned in footnote 2 there are detailed examples and explanations for all 10 data sets. In this chapter, we discuss the process of applying SML to the Iris Dataset and the Auto-MPG dataset. We compare the process for using machine learning to solve the problems presented by the datasets with SML against traditionally writing code. We forgo comparing SML to prior works in section 2.2. We used the same libraries and programming language in SML to solve these use cases for both of these datasets.

---

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets/Iris>

<sup>4</sup><https://archive.ics.uci.edu/ml/datasets/Auto+MPG>

<sup>5</sup><https://archive.ics.uci.edu/ml/datasets/seeds>

<sup>6</sup><https://archive.ics.uci.edu/ml/datasets/Computer+Hardware>

<sup>7</sup><https://archive.ics.uci.edu/ml/datasets/Housing>

<sup>8</sup><https://archive.ics.uci.edu/ml/datasets/Wine>

<sup>9</sup>[https://archive.ics.uci.edu/ml/datasets/US+Census+Data+\(1990\)](https://archive.ics.uci.edu/ml/datasets/US+Census+Data+(1990))

<sup>10</sup>[https://archive.ics.uci.edu/ml/datasets/Chronic\\_Kidney\\_Disease](https://archive.ics.uci.edu/ml/datasets/Chronic_Kidney_Disease)

<sup>11</sup><https://archive.ics.uci.edu/ml/datasets/Spambase>

<sup>12</sup><https://www.kaggle.com/c/titanic>

## Iris Dataset

Listing 2.11 shows all of the code required to perform classification on the Iris dataset using SML in Python. We read in data in listing 2.11 from a specified path named "iris.csv" from a subdirectory called "data" and perform an 80/20 split, uses the first four columns to predict the fifth column, uses support vector machines as the algorithm to perform classification and finally plot distributions of features from the dataset and metrics of the classification model's performance. Appendix A.1 illustrates what is required to perform the same operations using Python and sci-kit learn. The *Query* in listing 2.11 and the code in Appendix A.1 use the same 3rd party libraries implicitly or explicitly. It is also worth noting that the code in Appendix A.1 is publicly available and well documented<sup>13</sup> and it is out of the scope of this paper. Instead the complexities required to produce such results with and without SML are outlined. The result for both snippets of code are the same and is in Figure 2.4.

## Auto-Mpg Dataset

Listing 2.12 shows the SML *Query* required to perform regression on the Auto-MPG dataset in Python. In listing 2.12 we read data from a specified path, fixed-width spaces separate the dataset, and we choose not to provide a header for the dataset. Next, we perform an 80/20 split, replace all occurrences of "?" with the column's mode. We then perform linear regression using columns 2-8 to predict the 1st label. Lastly, we visualize distributions of features from the dataset and metrics of our algorithm. Appendix A.2. demonstrates what's required to perform the same operations using sci-kit learn<sup>14</sup>. The outcome of both processes are the same and can be seen in Figure 2.5.

---

<sup>13</sup>For detailed documentation describing this code visit: [https://github.com/lcdm-uiuc/sml/blob/master/dataflows/plot/iris\\_svm-READ-SPLIT-CLASSIFY-PLOT.ipynb](https://github.com/lcdm-uiuc/sml/blob/master/dataflows/plot/iris_svm-READ-SPLIT-CLASSIFY-PLOT.ipynb)

<sup>14</sup>For a detailed documentation describing this code visit: [https://github.com/lcdm-uiuc/sml/blob/master/dataflows/plot/autompq\\_linear\\_regression-READ-SPLIT-REGRESS-PLOT.ipynb](https://github.com/lcdm-uiuc/sml/blob/master/dataflows/plot/autompq_linear_regression-READ-SPLIT-REGRESS-PLOT.ipynb)

### 2.6.1 Discussion

For the Iris and Auto-MPG use cases, the same libraries and programming language were used to perform regression and classification. The amount of work required to perform a task and produce the following results in Figure 2.5 and Figure 2.4 significantly decreases when SML is utilized. Constructing each SML query used less than 10 lines of code; however, implementing the same procedures without SML using the same programming language and libraries needed 70+ lines of code. This provides evidence that SML simplifies the development process of solving problems with machine learning.

## 2.7 Future Work

While we have formally introduced an agnostic framework, much work remains to improve SML. In the future, we plan to extend the connector to support more machine learning libraries and additional languages. We can extend SML’s web application to include additional functionality to improve the ease of use. Feature selection, model selection, and parameter optimization are additional areas to add to SML. In addition to improving SML it can also be tested by researchers in a comparative analysis against the works outlined in section 2.2 to determine how beneficial SML is against alternative frameworks.

## 2.8 Conclusion

To summarize, we introduced an agnostic framework that integrates a query-like language to simplify the development of machine learning pipelines. We provided a high-level overview of its architecture and grammar. We then applied SML to machine learning problems and demonstrated how the code one has to write significantly decreases when SML is used. The source code and detailed documentation for SML is open-sourced and publicly available on

github<sup>15</sup>.

SML provides a realm of possibility to rapidly develop machine learning pipelines with an agnostic language to solve problems. This attractive aspect can boost the productivity of researchers who utilize machine learning. Abstracting machine learning complexities with a tool like SML can foster new research and solve problems in different disciplines faster.

---

<sup>15</sup><https://github.com/lcdm-uiuc/sml>

## 2.9 Figures and Listings

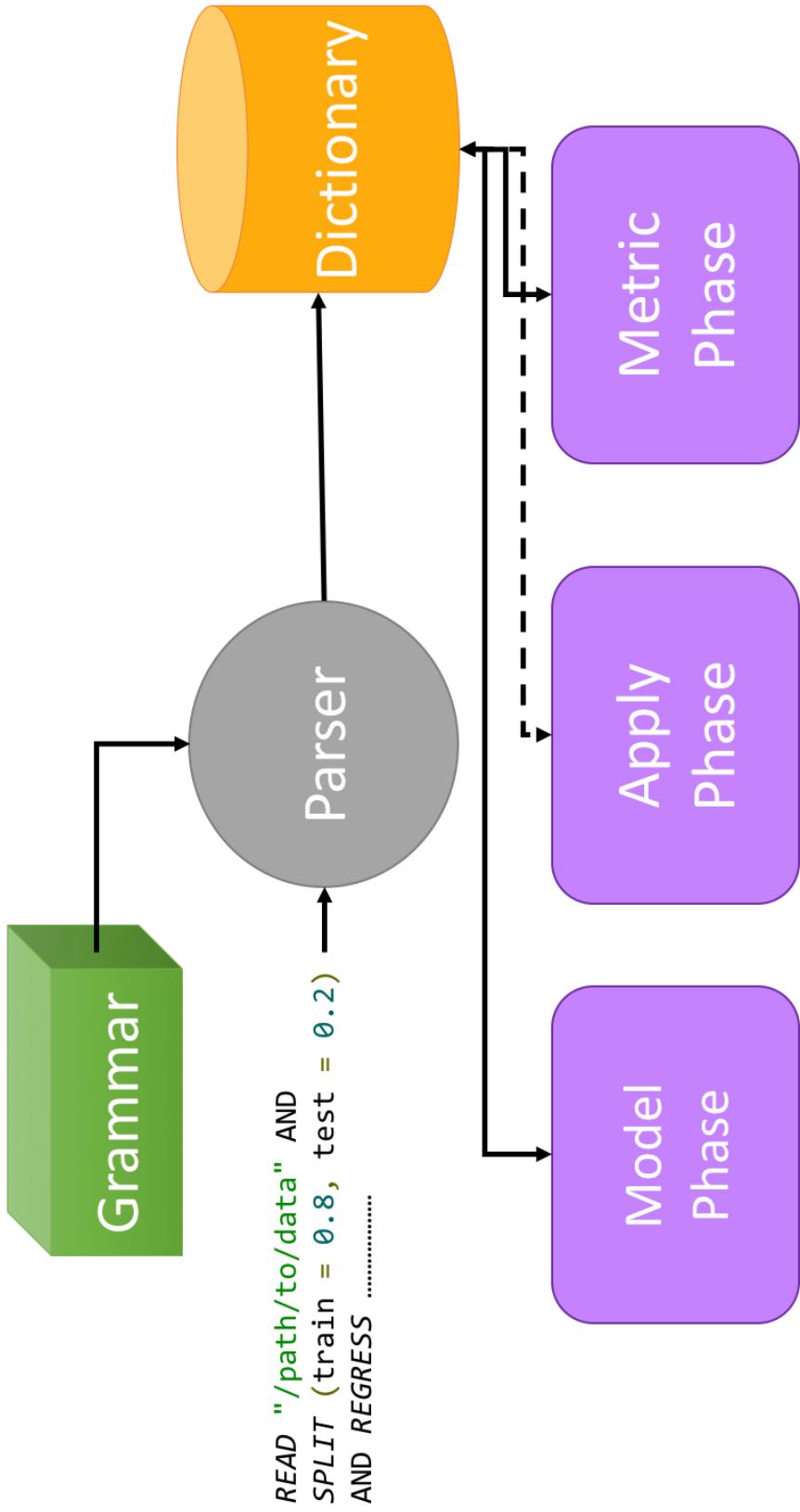


Figure 2.1: This figure shows a Block Diagram of SML's Connector.

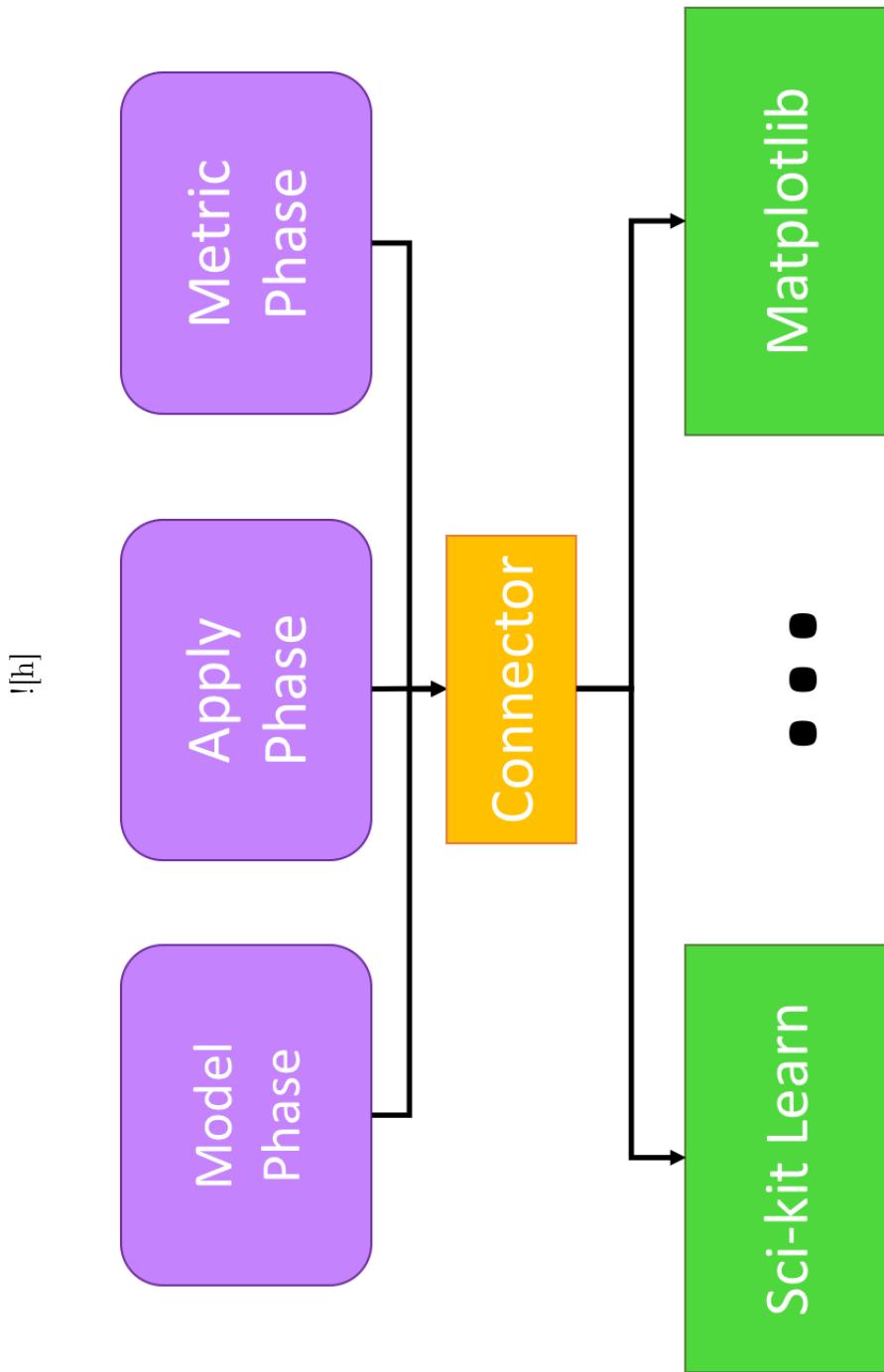


Figure 2.2: This figure shows a Block Diagram of SML's Connector.

Fork me on GitHub

# SML - Standard Machine Learning Language

SML aims to provide a universal language agnostic framework to simplify the development of machine learning pipelines.

**Instructions for using SML**

You can try using SML on the web before downloading it. The middle section contains a valid SML query. If you press the green "Execute" button on the bottom of that section the right section will update. The right section updated shows the results of the SML query. This SML Query read in data, split it into training and testing data, and performed classification on it and produces metrics and visualizations of the data. Here's a list of publicly available datasets that are readily available that you can try with SML on the web.

You can also upload data by clicking on the blue button that says upload. A path will be generated for you to insert in the SML query.

Editor

```

1.1 READ "data/Arts.csv" AND
1.2 SPILT <train, ~0.8, test = 0.2> AND
1.3 CLASSIFIER [predictor = [1,2,3,4],
1.4 label = 5, algorithm = sm]

```

Results

Figure 2.3: This figure shows the interface of SML’s webapp. Currently, users can read instructions, and examples of how to use SML are on the left pane. In the middle pane, users can type an SML Query and then hit the execute button. The results after executing a SML Query through SML are in the right pane.

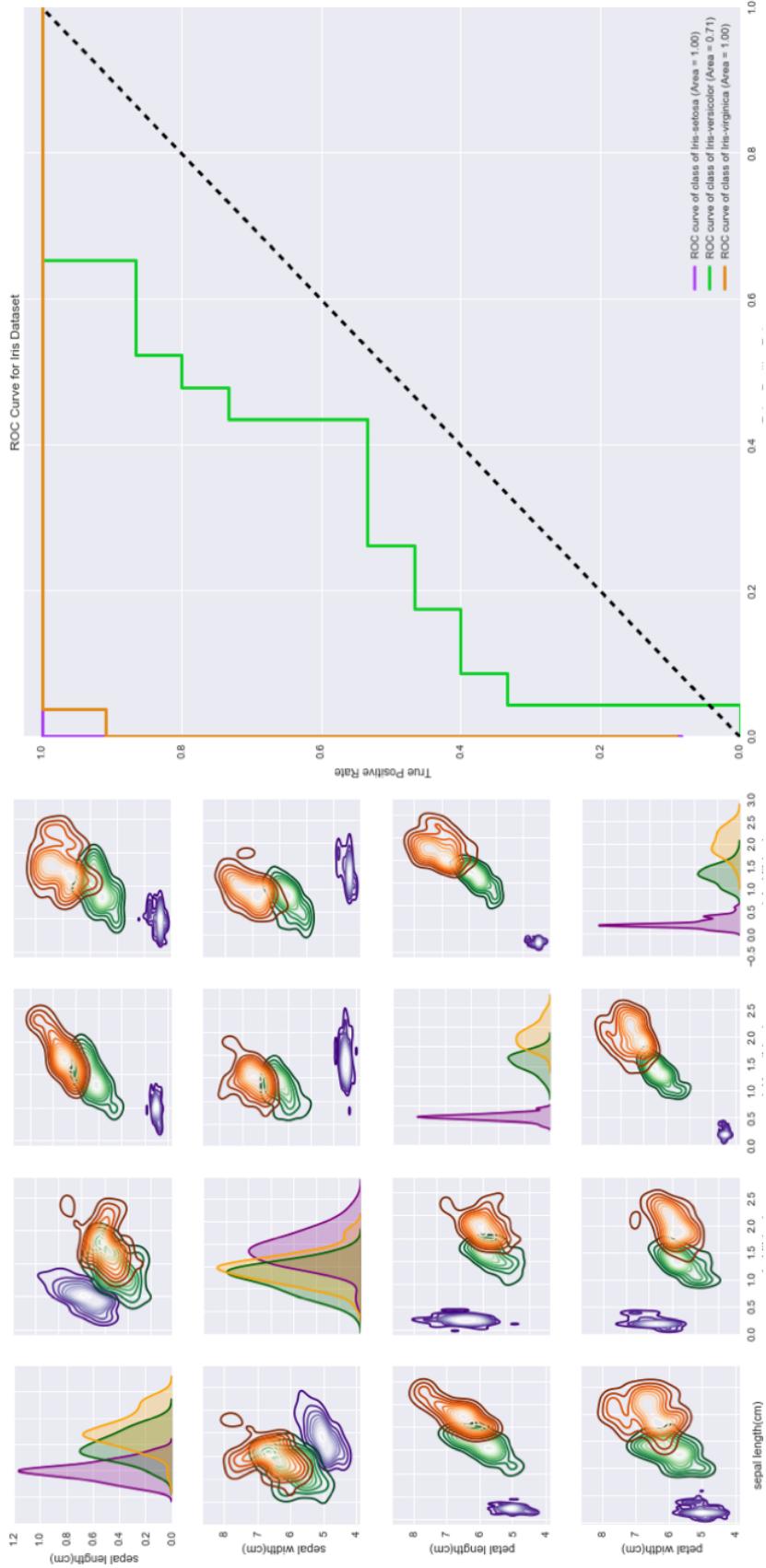


Figure 2.4: The SML *Query* in figure ?? and the code in figure ?? produce these results. The subgraph on the left is a lattice plot showing the density estimates of each feature used. The graph on the right shows the ROC curves for each class of the iris dataset.

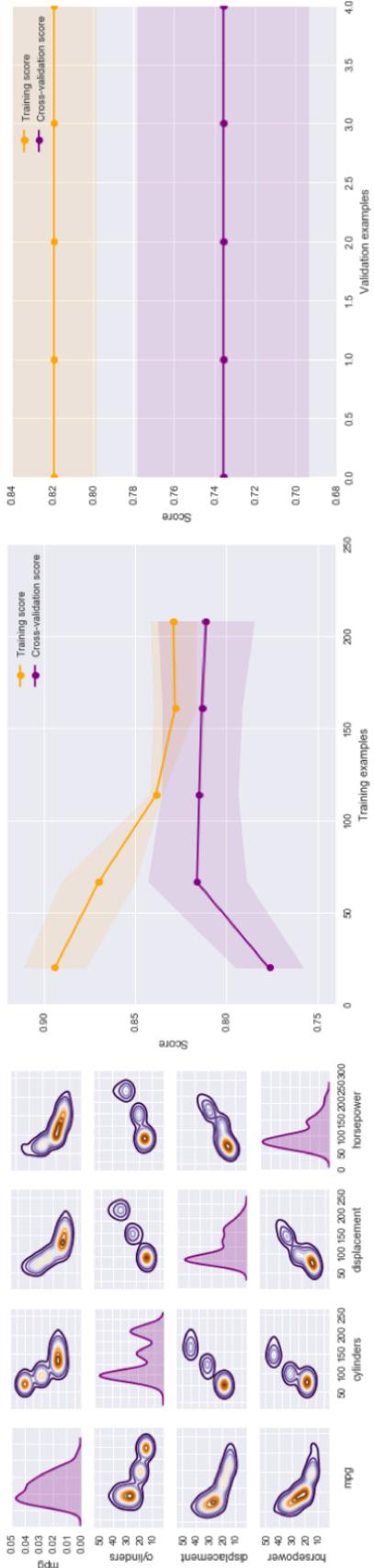


Figure 2.5: The SML *Query* in figure ?? and the code in Appendix A.2. produce these results. The subgraph on the left is a lattice plot showing the density estimates of each feature used. The top right graph shows the model's learning curve and the graph on the lower right shows the validation curve.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test=0.2) AND CLASSIFY
3 (predictors =[1,2,3,4] , label = 5 , algorithm = svm)

```

Listing 2.1: Example of a SML Query Performing Classification.

```

1 <Keyword> <Argument> (<OptionList>)
2 AND <Keyword> (<OptionList>) AND <Keyword>
3 (<OptionList>)

```

Listing 2.2: Here the example *Query* in listing 2.1 is defined in BNF format.

```

1 READ "/path/to/data"
2 READ "/path/to/data" (separator = ";" , header = None)

```

Listing 2.3: Examples using the *READ Keyword* in SML.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND REPLACE (missing = "NaN" , strategy = "mode")

```

Listing 2.4: An example utilizing the *REPLACE Keyword* in SML.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test = 0.2)

```

Listing 2.5: Example using the *SPLIT Keyword* in SML.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test = 0.2) AND CLASSIFY
3 (predictors = [1,2,3,4] , label=5, algorithm=svm)

```

Listing 2.6: Example using the *CLASSIFY Keyword* in SML. Here we read in data and create training and testing datasets using the *READ* and *SPLIT Keywords* respectively. We then use *CLASSIFY Keyword* with the first 4 columns as features and the 5th column to perform classification using a support vector machine.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test = 0.2) AND CLUSTER
3 (predictors = [1 ,2 ,3 ,4] , algorithm=kmeans)

```

Listing 2.7: Example using the CLUSTER Keyword in SML. Here we read in data and create training and testing datasets using the READ and SPLIT Keywords respectively. We then use CLUSTER Keyword with the first 4 columns as features and perform unsupervised clustering with the K-Means algorithm.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test = 0.2) AND CLUSTER
3 (predictors = [1 ,2 ,3 ,4] , label=5, algorithm=ridge )

```

Listing 2.8: Example using the *REGRESS Keyword* in SML. Here we read in data and create training and testing datasets using the *READ* and *SPLIT Keywords* respectively. We then use *REGRESS Keyword* with the first 4 columns as features and the 5th column to perform regression on using ridge regression.

```

1 SAVE "/path/to/save/model"
2 LOAD "/path/to/save/model"

```

Listing 2.9: Example using the *LOAD* and *SAVE Keywords* in SML.

```

1 READ "/path/to/data" (separator = ";" , header = None)
2 AND SPLIT (train = 0.8 , test = 0.2) AND CLUSTER
3 (predictors = [1 ,2 ,3 ,4] , algorithm=kmeans)
4 AND PLOT

```

Listing 2.10: Example using the *PLOT Keyword* in SML.

```

1 from sml import execute
2 query = 'READ "../data/iris.csv" AND \
3 SPLIT (train = 0.8 , test = 0.2) AND \
4 CLASSIFY (predictors = [1 ,2 ,3 ,4] , label = 5, algorithm=svm) AND \
5 PLOT'
6
7 execute(query , verbose=True)

```

Listing 2.11: SML *Query* that performs classification on the iris dataset using support vector machines. It's important to note that detailed documentation is publicly available in <sup>13</sup> and the purpose of this figure is to highlight the level of the level of complexity relative to an SML query.

```

1 from sml import execute
2 query = 'READ "../data/auto-mpg.csv" AND \
3 REPLACE (missing = "?", strategy = "mode") AND \
4 SPLIT (train = 0.8 , test = 0.2) AND \
5 REGRESS (predictors = [2 ,3 ,4 ,5 ,6 ,7 ,8] , label = 1, algorithm=simple) AND \
6 PLOT'
7
8 execute(query , verbose=True)

```

Listing 2.12: SML *Query* that performs regression on the Auto-MPG dataset using Linear Regression.

# Chapter 3

## Directional Profitability Predictions with Random Forests

### 3.1 Introduction

An economic goal of a society is to maximize wealth. In order to maximize wealth, resources such as labor, capital, and natural resources must be allocated efficiently to firms. Capital markets (such as the stock market) allow capital to be allocated to said firms efficiently. Firms' profitability can affect society, whether through job growth/loss, services offered, provision of goods, or personal gain/loss (see Spiceland et al. (2018)). Predicting a firm's future profitability is an essential endeavor as it is central to the valuation of the firm and the securities they issue (see Monahan (2018)).

A study by Bradshaw et al. (2012) provides evidence that analysts' predictions are not superior to random walk predictions. Monahan (2018) conducts a review and discussion on accounting research related to out-of-sample profitability predictions. He provides evidence that models in prior works have lower out-of-sample accuracy than random walk models. This provides motivation to seek an alternative method, given the shortcomings of traditional methods. Sarsam et al. (2020), Hang and Banks (2019), Choudhury et al. (2021), and Yang (2020) are relatively recent examples which exhibit the effectiveness of machine learning for

---

This chapter contains material from the following working paper:

- Vic Anand, Robert Brunner, Kelechi Ikegwu, and Theodore Sougiannis. Predicting profitability using machine learning. *SSRN*, Oct 2019. URL [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3466478](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3466478)

prediction and inference in a variety of applications. We explore machine learning methods to predict the future profitability of firms.

An approach outlined in Monahan (2018) is to select: what to predict (labels), the predictors (or features), a regression model where the researcher determines the functional form, and the training and testing datasets. It is easier with a parametric approach to assume the functional form of  $f$  (see 1.1) and estimate a set parameters; however, it is likely that the functional form will not match the true unknown form of  $f$  which would provide poor results. Given the objective to find the best  $f$  in eq. 1.1 that predicts  $Y$  reasonably well, non-parametric methods will be used. In addition to this cross-validation described in section 1.1.2, is used to avoid overfitting the model. Lastly, we use machine learning to predict directional changes rather than the magnitude of changes. We make this choice to account for working with a smaller dataset. Ball and Brown (1968) and Ou and Penman (1989) have shown that the future directional change of profitability is a useful indicator.

The research objective is to determine if a machine learning model can outperform a random walk model. This chapter will describe the datasets used in this study and the methods used to predict firms' directional change of profitability. Then comparative analyses will be conducted with the proposed model and random walk model, followed by a discussion of results and potential future directions of this research.

## 3.2 Data

We follow the data selection process described in Hou et al. (2012) in this research. We use all observations from the fiscal years 1963 - 2016 from Compustat fundamentals annual file. We merge this data set with the CRSP monthly returns file, which includes all securities listed on the NYSE, Amex, and Nasdaq. The securities selected have share codes 10 or 11, indicating that they are ordinary common shares of US companies and exclude ADR's and REITs. We remove observations if total assets, common equity, dividends, income before

extraordinary items, or accruals are missing, resulting in 194,172 observations. Given the objective to predict the future direction of profitability for firms, each firm must have at least three observations; the rationale behind this choice is to compute two price changes for a firm. Excluding firms with fewer than four years of observations, the final dataset consists of 166,925 observations.

The following variables were retrieved and computed for this research: Return on common equity (ROE), Return on assets (ROA), return on net operating assets (RNOA), Cash flow from operations (CFO), and free cash flow (FCF). Table 3.1 contains more details about the variables used from Compustat and table 3.2 shows how the profitability measures are computed.

Table 3.3 provides descriptive statistics on variables retrieved from Compustat. Outliers are present in most variables, and all of them have positively skewed distributions given that the mean is greater than the median. Table 3.4 provides the descriptive statistics on profitability measures. Since the raw variables construct the profitability measures, their values are also extreme. However, unlike the raw variables, the distributions of the targets are negatively skewed (the mean is smaller than the median). Table 3.5 shows the percentage of times the target variables increase in the holdout period (2012-2015). With the exception of CFO and FCF during 2013, there are no relatively significant class imbalances for directional changes in profitability measures.

Tables 3.6, 3.7, 3.8, 3.9, and 3.10 show the Pearson correlations between variables that will serve as features in the random forests. Since correlations of raw variables are extreme due to large outliers, we truncate each variable at 2.5% on each side of its distribution. Many of the features in these tables are correlated. Nevertheless, the proposed machine learning method is insensitive to outliers. Table 3.11 shows the autocorrelations of truncated profitability features. With the exception of CFO, the second-order autocorrelations are marginally smaller than the first-order autocorrelations for all profitability measures.

## 3.3 Methods

### 3.3.1 Cross Validation for Time Series

The supervised learning paradigm described in section 1.1.1 will be employed. We use 90% of the data for training and the remaining 10% for testing. The first 90% of data contains data from the fiscal years 1963-2011, and the last 10% of data is from 2011-2016. Given the objective to predict the year-to-year change, we cannot predict using features for firms in the fiscal year 2016 because the fiscal year 2017 is not in the sample. Given the time-series nature of the data, blocked K-fold cross-validation as described in Cerqueira et al. (2019) is used. The difference between blocked  $k$ -fold cross-validation and  $k$ -fold cross-validation is that the observations are not randomly shuffled and then divided into  $k$  folds. Instead, the observations are split into  $K$  blocks of contiguous observations.

### 3.3.2 Tree based Machine Learning Methods

Many machine learning models can provide high accuracy scores. However, there is a trade-off between interoperability and performance (see James et al. (2013), Chapter 2). Given that the research objective is more aligned with prediction than inference, we emphasize model performance. Tree-based models are flexible, simple to implement, and have high performance (see James et al. (2013), Chapter 8). In addition to this, tree-based methods can handle multicollinearity and non-linearity better than linear regression methods traditionally used in the financial accounting setting (see Monahan (2018)).

### 3.3.3 Decision Trees

Decision Trees are the building blocks of tree-based machine learning methods. Decision Trees are directed graphs with no cycles meaning there are no nodes with edges to the same node, there is no more than one edge connecting two nodes, and there is only one path

connecting two nodes. We refer to edges in trees as branches. Consider three nodes node  $A$ ,  $B$ , and  $C$  if node  $A$  has a branch to nodes  $B$  and  $C$  then,  $B$  and  $C$  are the children of node  $A$ ; node  $A$  in this case is the parent of nodes  $B$  and  $C$ . The root node starts the decision tree and does not have a parent. Lastly, a node with no children signifies a stopping point in the decision tree and is referred to as a leaf node (see Figure ?? for an example).

In a supervised machine learning paradigm, decision trees are given data to train on. Each node in the decision tree will represent the training features (or feature space, i.e.  $X_1, \dots, X_p$ ) partitioned into particular regions. Using the CART algorithm (see Breiman et al. (1984)) a decision tree will create  $m$  distinct non-overlapping regions  $R$ . The root node will contain a rule to partition the feature space into two non-overlapping regions ( $R_1$  and  $R_2$ ). The process will repeat recursively i.e.  $R_1$  will be split by a splitting criteria to partition the data into 2 additional regions ( $R_3$  and  $R_4$ ) and sub-sequentially create 2 additional regions from  $R_2$  ( $R_5$  and  $R_6$ ).

In a classification setting, a decision tree at a given region will determine the most optimal feature to select from the set of features to partition the data. The selection is typically based on the feature's ability to minimize error. A common metric is gini-impurity, where gini-impurity measures the variance across  $C$  classes. We can define gini-impurity as:

$$G_{index} = \sum_{i=1}^C p(i) \cdot (i - p(i)) \quad (3.1)$$

where  $C$  is the number of classes in a label and  $p(i)$  is the probability of an observation labeled class  $i$ . A small value of gini-impurity indicates that a particular region has many observations from a single class.

Decision trees have stopping criteria. For example, we stop creating additional regions when a child has less than a defined number of observations at a particular region or when the difference of gini-impurity between the parent and its children is less than a specified amount. Lastly, to predict values of  $Y$  from another set of data (containing the same features), new

observations can be passed into the trained tree. The trained tree will select the appropriate region for an observation based on the splitting criteria rules.

Decision Trees offer an intuitive explanation of the results it produces. However, there are some caveats to this method. A few additional observations or any modification to observations in a feature space can drastically change a decision tree because it is using a greedy search to build a tree (see Koning and Smith (2017)). Meaning that at each stage of the tree-building process, a node is split using a local optimal choice instead of a choice that will find the optimal choice on a global level.

### 3.3.4 Random Forests

Random forests can overcome some of the shortcomings that decision trees face. Random forests are an ensemble of decision trees used to make predictions. However, a random forest is less interpretable than a single decision tree. A random forest will create  $B$  decision trees. For each decision tree, the random forest will assign a random dataset to each decision tree in the forest by sampling the data with replacement (or bootstrapping). In addition to this, the random forest will assign a subset of random features (with replacement) for each decision tree to use.

The rationale behind selecting only a subset of features is to decorrelate trees meaning that if a feature is important, it will not always be used as a top-level split producing unique decision trees. Each tree can then be built and make predictions using the random data and features in the method outline in the Decision Trees subsection at 3.3.3. To make predictions in a classification setting  $B$  decision trees will make a prediction using the features from the bootstrapped data. The most common prediction class is used as the prediction for the forest of trees (which is referred to as the majority vote).

## 3.4 Random Forest Implementation to Predict Directional Profitability

For this research, annual profitability changes (ROE, ROA, RNOA, CFO, and FCF) are coded as +1 for an increase and -1 for a decrease for individual firms in the dataset described in section 3.2. The features used are described in section 3.6. With data sorted by fiscal year, for a change in a specific profitability measure that feature set (see tables 3.6, 3.7, 3.8, 3.9, and 3.10) the data will be partitioned into training and testing sets. The first 90% of observations (fiscal years between 1963-2011) are used to train and build  $B$  trees. The remaining 10% of observations (2012-2015) are used to make predictions out of sample.

While building  $B$  decision trees, cross-validation, as described in 3.3.1 will be employed. The training data will have  $k$  folds (where  $k = 5$ ) grouped by fiscal year. Each particular fold will consist of observations for particular firm-years. Each decision tree in the forest will bootstrap from each fold and select random features. The decision tree will then train and build a tree based on the methodology outlined in 3.3.3. Then the random forest of trees will make predictions on the validation sets with the process outlined in section 3.3.4. Finally, the trained random forest will use make predictions out-of-sample on observations from the test set.

## 3.5 Random Walk Benchmark Models

The random walk is a model of a stochastic process. Campbell et al. (1997) (hereafter CLM) described it by the following equation:

$$P_t = \mu + P_{t-1} + \epsilon_t \quad (3.2)$$

In eq. 3.2,  $P_t$  is the value of some state variable at time  $t$ ,  $\mu$  is a drift term that we will assume is 0, and  $\epsilon_t$  is a disturbance, or increment. CLM note that a common assumption

is that the disturbance term is independent and identically distributed, with mean zero and variance  $\sigma^2$ , i.e.  $\epsilon_t \sim IID(0, \sigma^2)$ .

The random walk is a mathematical description of a process, not a forecasting tool. However, some papers (i.e. Bradshaw et al. (2012)) that wish to use the random walk as a benchmark for their results treat it as a forecasting tool by taking the expected value of both sides of eq. 3.2. That yields the following:

$$E[P_{t+1}] = P_t \Leftrightarrow E[P_{t+1} - P_t] = 0 \quad (3.3)$$

These papers then reason that, since the expected value of  $P_{t+1}$  is  $P_t$ , the best forecast of  $P_{t+1}$  is also  $P_t$ . However, the random walk does not actually predict  $P_{t+1}$ . It simply says that the expected magnitude of change in  $P$  at any time is zero. In other words, if you observe changes in  $P$  over a long time series, the average change will be zero.

To be consistent with prior literature, we treat the random walk model as a forecasting tool. However, our work introduces a complication. Prior literature forecasts the magnitude of some variable. By contrast, we forecast the direction of change in our variables. Applying the random walk to our setting, therefore, requires some additional work. CLM notes that the expected value of  $\epsilon_t$  is zero, but without additional assumptions, we cannot know the proportion of increases and decreases in  $\epsilon_t$ .

### 3.5.1 Random Walk: Benchmark 1

Rearranging equation 3.2 and setting the drift term to zero yields  $P_t - P_{t-1} = \epsilon_t$ . Thus, the change in profitability is exactly equal to the disturbance term. That implies that the sign of the change in profitability is equal to the sign of the disturbance term. By assumption,  $\epsilon_t \sim IID(0, \sigma^2)$ , and CLM (p.32) state that “perhaps the most common distributional assumption for the innovations or increments  $\epsilon_t$  is normality.” Normality implies a 50-50 chance of an up or down movement since the normal distribution is symmetric. Therefore, our first random

walk prediction should be that profitability is equally likely to increase or decrease.

Under the assumption that the error term is equally likely to increase or decrease, we find that the random walk model has an expected accuracy of 50%. To see this, assume the testing subsample contains  $N$  observations, and  $0 \leq p \leq 1$  is the fraction of increases. Then  $p$  of those observations increases, and  $(1 - p)N$  decreases. For any observation, the random walk will predict an increase 50% of the time and a decrease 50% of the time. Thus, of the  $p$  increases, on average  $0.5pN$  will be classified accurately. Similarly, of the  $(1 - p)N$  decreases,  $0.5(1 - p)N$  will be classified accurately. For the overall sample, the fraction that will be accurately classified is:

$$Accuracy = \frac{0.5pN + 0.5(1 - p)N}{pN + (1 - p)N} = \frac{1}{2}$$

Note that the classification accuracy of this interpretation of the random walk is invariant to the fraction of increases in the data. For any value of  $p$ , if we assume increases and decreases are equally likely, this random walk model will accurately predict on average the outcome 50% of the time.

### 3.5.2 Random Walk: Non-normal Error Term

The random walk model requires that the error term has zero mean and finite variance,  $\epsilon_t \sim I(0, \sigma^2)$ . This restriction does not imply normality of the error term:  $\epsilon_t$  could have any arbitrary probability of increase as long as its expected value is zero. For example, if  $\frac{2}{3}$  of the time  $\epsilon_t$  increases by 1, and  $\frac{1}{3}$  of the time it decreases by 2, it has mean zero. A distributional assumption other than normality may be appropriate if the actual fraction of increases in the data is not 50%.

Anand et al. (2019) showed in Appendix B that even if we assume  $\epsilon_t$  will increase 40%-60% of the time, the classification accuracy ranges from 48%-52% which is very close to the 50% by assuming a normally distributed error term. Table 3.5 shows that, for most measures

in most years, the actual fraction of increases in our sample are very close to 50%. Even in the most extreme case (in 2013, CFO and FCF increase 43.3% and 42.5% of the time).

### 3.5.3 Random Walk: Benchmark 2

As an alternative interpretation of applying the random walk model to forecasting, consider the following equations:

$$P_t = P_{t-1} + \epsilon \quad (3.4)$$

$$P_{t-1} = P_{t-2} + \epsilon_{t-1} \quad (3.5)$$

Subtracting eq. 3.5 from eq. 3.4 yields:

$$\begin{aligned} (P_t - P_{t-1}) &= (P_{t-1} - P_{t-2}) + (\epsilon_t - \epsilon_{t-1}) \\ \Delta P_t &= \Delta P_{t-1} + (\epsilon_t - \epsilon_{t-1}) \end{aligned} \quad (3.6)$$

If we take the expected value of both sides, equation 3.6 reduces to  $\Delta P_t = \Delta P_{t-1}$  since, by assumption,  $E[\epsilon_t] = 0, \forall t$ . Under this interpretation, in expectation, an increase in a profitability measure in period  $t - 1$  will persist into period  $t$ . Thus, our second random walk prediction is that the sign of the change in a measure in one period will persist into the subsequent period.

## 3.6 Comparative Analyses

To answer the research objective, we conduct comparative analyses between the benchmark random walk models (described in 3.5) and the random forest models (described in 3.3.4). We

employ multiple random forests to predict specific profitability measures and add additional features to each random forest for each subsequent analysis. These analyses provide insight as to whether specific features are informative and yield more accurate estimates.

We perform the analyses in two phases. In the first phase, we test whether Random Forests performs better with winsorized or un-winsorized data. In the second phase, we employ multiple random forests to learn the incremental value of additional variables used in fundamental analysis.

### **3.6.1 Phase 1: Effect of Winsorization**

The purpose of phase 1 is to test the effect of winsorization on predictive accuracy. Our prior is that winsorization will not affect the predictive accuracy of random forests. In phase 1, we use a minimal information set in which each firm-year observation consists of the firm identifier (PERMNO), fiscal year, the contemporaneous value of a measure (i.e., ROE), and the label or change that occurs in the subsequent year (i.e., increase or decrease). We create three random forests for this phase. We train the first random forest on un-winsorized data and refer to this as analysis 1-1. For analysis 1-2, we train the second random forest on data where the measure is winsorized at 1% on each side. For analysis 1-3 we repeat the same methodology as analysis 1-2 where the measure is winsorized at 2.5% on each side.

### **3.6.2 Phase 2: Effect of Adding Additional Fundamental Information**

In phase 2, we progressively add information about each firm's industry. We begin with analysis 2-1, and add SIC codes as a feature. Analysis 2-2 adds industry information to the feature set. For each firm-year observation in analysis 2-2 we add, the industry mean, median, and standard deviation for the measure. For example, consider predicting the change in annual RNOA for IBM during the fiscal year 2000. The industry mean, median,

and standard deviation of RNOA for IBM's Industry are used as features during the year 2000. The firm's 4-digit SIC code determines the industry.

Analysis 2-3 repeats analysis 2-2 while including the lag of the profitability measure. For example, if the measure is ROE, then for IBM in the year 2000 its ROE from 1999 is included as a feature in addition to other features from analysis 2-2. For the final analysis 2-4, we repeat analysis 2-3 while including the 2nd lag of the profitability measure. Following the IBM example for analysis 2-3, its ROE from 1998 is included as a feature.

## 3.7 Results

### 3.7.1 Phase 1 Results and Discussion

Phase 1 tests the predictive power of random forests with minimal information. It also tests the effect of winsorizing the input data, a common practice in accounting research. In phase 1, we create three random forests per profitability measure and use a minimal set of features. For each random forest, the only predictors (features) in each observation are a firm identifier (PERMNO), fiscal year, and the contemporaneous value of the profitability measure. The label is a categorical variable that indicates whether profitability increased or decreased in the subsequent year. The input data's non-categorical measures are winsorized at three levels, 0%, 1%, and 2.5% on each side. Analysis 1-1 uses a random forest to make predictions on the non-winsorized data, 1-2 uses the 1% winsorized data, and 1-3 uses the 2.5% winsorized data.

We begin with the default values for model parameters in sci-kit learn version 0.22 (see Pedregosa et al. (2011)). We varied the number of trees used in a forest for each measure to predict the change in profitability. We found for all measures that there was no significant improvement in performance after using 30 decision trees in a random forest, see Figure 3.1. Figure 3.1 all subplots share the same x-axis where the x-axis represents the number of trees used in a forest. The y-axis for each subplot represents the accuracy score for a particular

measure.

In each subplot, the dashed line with the circle markers represents the mean training scores during cross-validation. The shaded region surrounding the dashed line with circle markers represents the variance of the training accuracy scores during cross-validation. This shaded region is not visible because the variance of the  $k - 1$  training scores is very small. The solid line with star markers represents the mean cross-validation accuracy score, and the shaded-in region represents the variance of the cross-validation scores for each fold.

Given the large gap between the mean training and cross-validation accuracy scores, it is probable that the random forests are over-fitting the data. With the current model parameters (while varying the number of trees used in the random forest), the model is learning information from the training data that is noise. The current parameters impede the model's ability to generalize to the cross-validation data well. To decrease the model's variance (or flexibility), we vary the depth of each decision tree in the forest.

Figure 3.2 shows the mean training and cross-validation scores for a random forest with 30 trees as the maximum depth is varied. As the depth increases, the model becomes more flexible. Across all profitability measures, one can see that the model can learn from the training data fairly well with high training accuracy scores. However, the highly flexible models cannot generalize well to the cross-validation data and perform worst than other models with a smaller maximum depth. Ergo, we select the maximum depth, which yields a model with relatively low variance and bias for each profitability measure. We forgo displaying similar plots for analyses 1-2 and 1-3 because the effect and accuracy scores are similar.

Table 3.12 shows the model maximum depth selected for Analyses 1-1, 1-2, and 1-3. Tables 3.13, 3.14, 3.15, 3.16, and 3.17 show the mean training, mean cross-validation, and testing scores for the ROE, ROA, RNOA, CFO, and FCF measures respectively. In addition to showing the cross-validation scores we show the random walk models performance on the data, and the feature importance of each variable.

The results show that the random forest can outperform the random walk models for each measure in predicting the change of profitability with a minimum set of features. In addition to this, excluding the ROE measure, we find no significant difference in mean testing accuracies between winsorized and non-winsorized data. Analysis 1-1 (no winsorization) and 1-3 (2.5% per side winsorization) had a mean test accuracy of 55.2%; however, analysis 1-2 (1% per side winsorization) had a mean test accuracy of 56.8%. The remaining analyses for the other profitability measures have mean testing accuracies within 1% of each other. These results provide evidence that winsorization has a negligible effect on performance out of sample.

To summarize, traditional approaches have struggled in the past to out-perform random walks. With a minimum amount of features (company identifier, profitability measure, and fiscal year), we find that random forest are able to out-perform random walks in predicting annual profitability changes. In addition to this we find that winsorizing the data (which is a traditional method used to handle outliers in financial accounting) has a negligible effect on the performance of a random forest model.

### 3.7.2 Phase 2 Results and Discussion

Given that there is no significant effect of winsorization concerning random forests performance, we forgo using the winsorized datasets. We create additional random forests where each successive forest incrementally adds features. In analysis 2-1 the firm's industry is included as a feature. Analysis 2-2 includes descriptive statistics about the firm's industry at a given year. The specific statistics used are mean, median, and standard deviation. Analysis 2-3 adds the previous lag value of the profitability measure (i.e., the profitability measure at  $t - 1$ ) and analysis 2-4 includes the two previous lag values of the profitability measure (i.e.,  $t - 1, t - 2$ ). Tables 3.18, 3.19, 3.20, 3.21, and 3.22 report the results from the 2nd set of analyses.

Similar to the first set analyses, all random forest models outperform the random walk

models. Comparing the second set of analyses to analysis 1-1 (non-winsorized data with a minimum set of features), we find mixed results for each profitability measure. For CFO and FCF, we find improvements across all sets of analyses. However, for ROE, ROA, and RNOA there is little to no improvement from analysis 1-1. We find that CFO and FCF random forest models utilize the feature  $CFO_t$  and  $FCF_t$  the most throughout all analyses. This finding provides evidence that the best feature to predict the future directional change of a measure is the current value of the measure for CFO and FCF.

RNOA, ROA, and ROE set of analyses still utilize  $RNOA_t$ ,  $ROA_t$ , and  $ROE_t$  respectively more than any other feature; however, they utilize it significantly less than the random forest models that predict CFO and FCF. For example, in analysis 2-3 for RNOA, ROA, ROE, CFO, and FCF the previous lagged measures are the  $2^{nd}$  most used by the model. The feature importance of CFO and FCF are 69.2% and 78.5% more than the lagged value ( $CFO_{t-1}$  and  $FCF_{t-1}$ ). However, the ROE, ROA, and RNOA measures use their respective current values at  $t$  45.7%, 50.3%, and 45.8% more than the lagged values. In addition to this, we find that ROE's performance improves more than ROA and RNOA. This finding further provides evidence that the most information about future profitability changes is in the profitability measure's current value.

### 3.8 Conclusion and Future Work

We explored the utility of random forest trees in a supervised learning paradigm to predict the annual direction of profitability for firms with minimal information. This study's motivation stems from the fact that accounting literature shows that traditional regression methods cannot produce models superior to random walk models at predicting out of sample. We also explore the advantages of following a supervised learning paradigm to predict directional change. For example, random forests are insensitive to multicollinearity, and properly utilizing the supervised paradigm discovers a functional form that generalizes well

to new data.

We implement and train random forest trees in a classification setting using a large sample of US firms over the period 1963-2011. We generate out-of-sample predictions of directional changes (increases or decreases) in five profitability measures: return on equity (ROE), return on assets (ROA), return on net operating assets (RNOA), cash flow from operations (CFO), and free cash flow (FCF) from 2011-2016. We found that the classification accuracy for each measure outperformed the random walk models. The accuracy does not tend to decline with a four-year forecast horizon significantly for the out of sample predictions. In addition to these findings, models with a strong reliance on a profitability measure's current value provided more accurate estimates of the future directional change in profitability.

We can improve these results further by selecting a model that yields better performance but offers less interpretability. Given that we focus on a minimal set of features in this research, we can incorporate additional features for additional performance gains. Utilizing more data, we can extend the methodology outlined in this paper to a regression setting and make predictions about the magnitudes of profitability.

The results shown in the paper are promising since that we can outperform a random walk model for all of the profitability measures used in this study. The methods used in this chapter are robust to econometric issues such as multicollinearity and outliers. When used correctly, the methodology outlined in this paper can produce models that generalize well to out of sample data. Machine learning has simplified and solved problems in many aspects of our lives. The results show the benefit that it can have on predicting the profitability of firms in a financial accounting research setting.

### 3.9 Figures and Tables

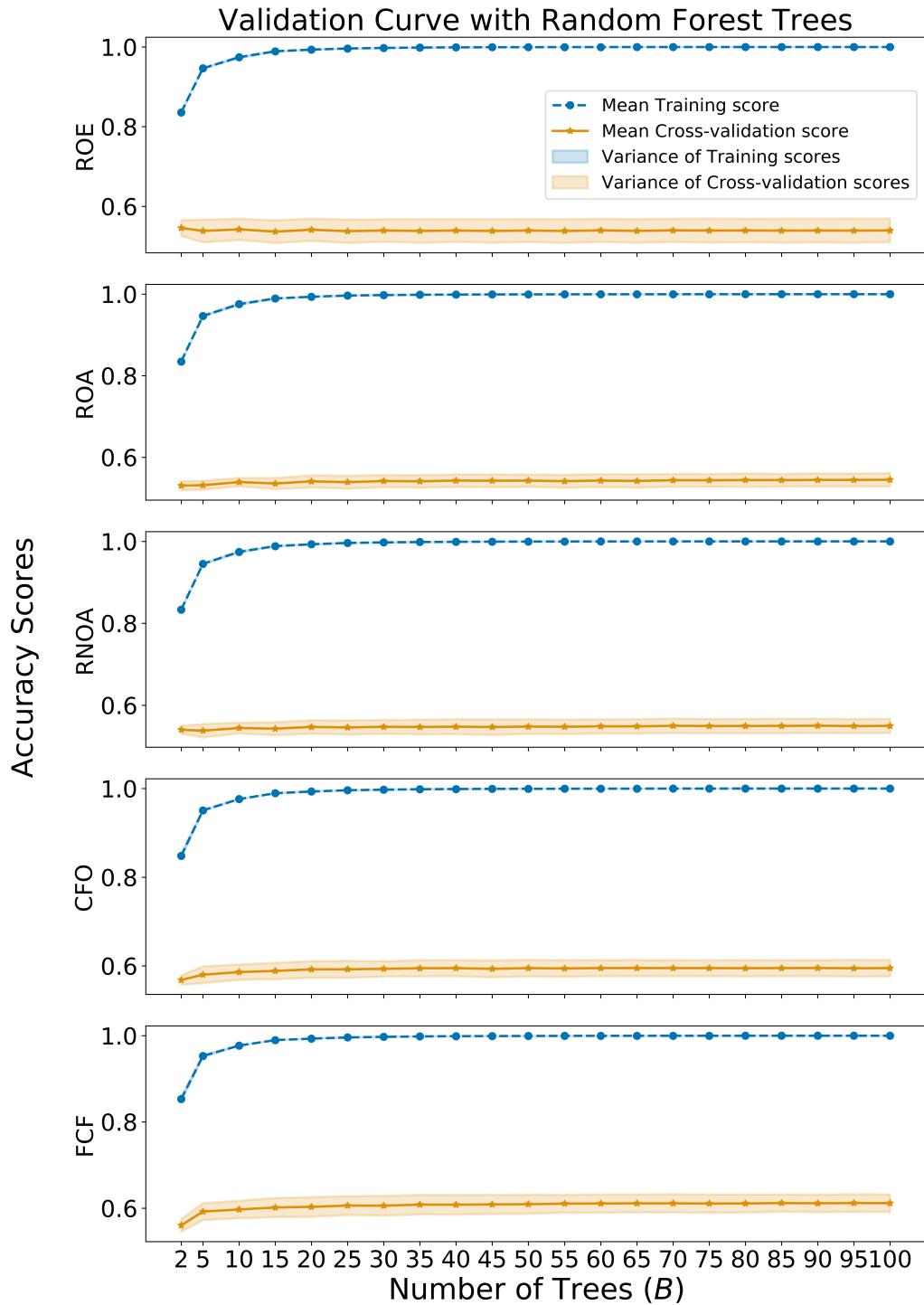


Figure 3.1: This figure shows validation curves for each measure in Analysis 1-1, varying the number of trees.

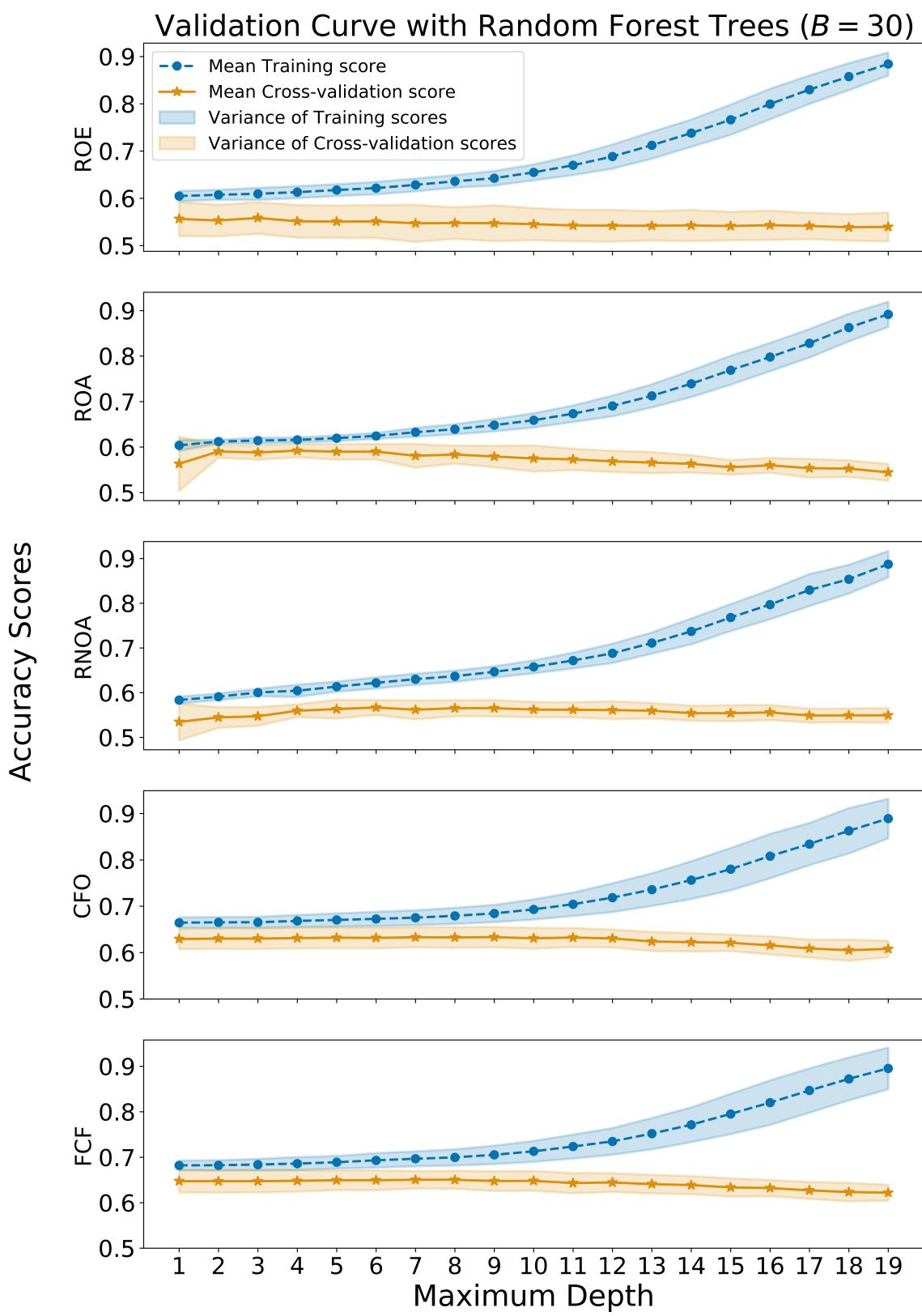


Figure 3.2: This figure shows the validation curves for each measure in Analysis 1-1 varying the max depth using 30 trees in a forest.

Variable	Code
Current Assets – Total	ACT
Assets – Total	AT
Capital Expenditures	CAPX
Common Equity – Total	CEQ
Cash and Short-Term Investments	CHE
Debt in Current Liabilities (short-term debt)	DLC
Long-Term Debt – Total	DLTT
Depreciation and Amortization	DP
Data year – fiscal	FYEAR
Earnings (income before extraordinary items)	IB
Earnings (income before extraordinary items, from statement of cash flows)	IBC
Interest and Related Income – Total	IDIT
Short-term investments – Total	IVST
Current Liabilities – Total	LCT
Liabilities – Total	LT
Operating Activities – Net Cash Flow	OANCF
Operating Income before Depreciation	OIBDP
CRSP permanent number	PERMNO
Property, Plant, and Equipment – Total (Net)	PPENT
Standard Industry Classification	SIC
Income Taxes Payable	TXP
Extraordinary Items and Discontinued Operations (Statement of Cash Flows)	XIDOC
Interest Expense	XINT

Table 3.1: Table showing Computstat variables used in this study.

Variable	Definition
$ROE_t$	$\frac{IB_t}{0.5(CEQ_t + CEQ_{t-1})}$
$ROA_t$	$\frac{IB_t + XINT_t}{0.5(A_t + A_{t-1})}$
$RNOA_t$	$\frac{OI_t}{0.5 \times (NOA_t + NOA_{t-1})}$
$CFO_t$	$\frac{OANCF_t}{0.5(A_t + A_{t-1})}$
$FCF_t$	$\frac{OANCF_t - CAPX_t}{0.5(A_t + A_{t-1})}$

Table 3.2: Table showing profitability measures definitions.

Variable	Count	Mean	Std	Min	1%	25%	50%	75%	99%	Max
AT	167,718.00	3,344.65	37,775.67	0	2.08	32.21	139.24	754.55	45,009.92	2,573,126.00
CAPX	167,213.00	103.64	665.28	-401.61	0	0.99	5.29	30.46	1,833.88	37,985.00
CEQ	167,718.00	733.70	4,746.54	-59,640.00	-46.55	13.94	59.19	272.44	11,881.49	255,550.00
DP	166,456.00	73.70	465.94	-7.91	0.02	0.95	4.20	23.01	1,250.90	22,016.00
IB	167,718.00	85.78	859.57	-99,289.00	-227.97	-0.32	3.38	25.86	1,844.83	53,394.00
OANCF	167,718.00	181.91	1,590.52	-110,560.00	-81.17	0.10	6.70	51.69	3,353.03	129,731.00
OIBDP	167,026.00	264.67	1,742.21	-76,735.00	-53.35	1.57	13.19	81.83	4,516.75	81,730.00
PPENT	166,826.00	620.94	3,822.74	0	0.05	4.86	24.81	154.47	12,305.75	252,668.00
XINT	167,718.00	49.38	639.00	-0.88	0	0.12	1.28	10.41	648.83	57,302.00

Table 3.3: This table shows the descriptive statistics for Compustat variables in the sample.

Variable	Count	Mean	Std	Min	1%	25%	50%	75%	99%	max
<b>ROE</b>	160,653.00	0.00	0.50	-4.47	-2.52	-0.01	0.10	0.17	1.10	2.97
<b>ROA</b>	160,653.00	0.02	0.18	-1.33	-0.85	0.01	0.06	0.10	0.26	0.35
<b>RNOA</b>	160,653.00	0.07	1.07	-12.03	-4.00	0.02	0.10	0.17	3.64	10.66
<b>CFO</b>	160,653.00	0.04	0.16	-0.91	-0.68	0.01	0.07	0.12	0.34	0.41
<b>FCF</b>	160,653.00	-0.03	0.17	-1.07	-0.76	-0.07	0.01	0.06	0.27	0.34

Table 3.4: This table shows the descriptive statistics for profitability measures.

Fisical Year	ROE	ROA	RNOA	CFO	FCF
2012	47.6%	48.4%	47.9%	49.4%	51.9%
2013	48.9%	49.7%	47.6%	43.3%	42.5%
2014	46.6%	45.5%	44.7%	49.1%	51.5%
2015	48.4%	51.0%	50.2%	50.8%	53.6%

Table 3.5: This table shows the percentage increases of annual profitability in our sample.

	<b>ROE</b>	Ind ROE Mean	Ind ROE Median	Ind ROE Std	ROE_L1	ROE_L2
<b>ROE</b>	1.00	0.20	0.39	-0.11	0.53	0.39
Ind ROE Mean		1.00	0.35	-0.63	0.14	0.12
Ind ROE Median			1.00	-0.19	0.33	0.29
Ind ROE Std				1.00	-0.12	-0.10
ROE_L1						0.54
ROE_L2						1.00

Table 3.6: This table shows the Pearson correlations for ROE features in our sample.

	<b>ROA</b>	Ind ROA Mean	Ind ROA Median	Ind ROA Std	ROA_L1	ROA_L2
<b>ROA</b>	1.00	0.53	0.49	-0.42	0.74	0.63
Ind ROA Mean		1.00	0.89	-0.84	0.48	0.43
Ind ROA Median			1.00	-0.60	0.45	0.42
Ind ROA Std				1.00	-0.38	-0.36
ROA_L1					1.00	0.74
ROA_L2						1.00

Table 3.7: This table shows the Pearson correlations for ROA features in our sample.

	<b>RNOA</b>	Ind RNOA Mean	Ind RNOA Median	Ind RNOA Std	RNOA_L1	RNOA_L2
<b>RNOA</b>	1.00	0.10	0.24	0.01	0.36	0.22
Ind RNOA Mean		1.00	0.25	0.26	0.06	0.03
Ind RNOA Median			1.00	0.09	0.18	0.14
Ind RNOA Std				1.00	0.01	0.01
RNOA_L1					1.00	0.36
RNOA_L2						1.00

Table 3.8: This table shows the Pearson correlations for RNOA features in our sample.

	<b>FCF</b>	<b>Ind FCF Mean</b>	<b>Ind FCF Median</b>	<b>Ind FCF Std</b>	<b>FCF_L1</b>	<b>FCF_L2</b>
<b>FCF</b>	1.00	0.49	0.45	-0.31	0.62	0.53
<b>Ind FCF Mean</b>		1.00	0.89	-0.71	0.39	0.36
<b>Ind FCF Median</b>			1.00	-0.46	0.37	0.33
<b>Ind FCF Std</b>				1.00	-0.28	-0.26
<b>FCF_L1</b>					1.00	0.61
<b>FCF_L2</b>						1.00

Table 3.9: This table shows the Pearson correlations for FCF features in our sample.

	<b>CFO</b>	<b>Ind CFO Mean</b>	<b>Ind CFO Median</b>	<b>Ind CFO Std</b>	<b>CFO_L1</b>	<b>CFO_L2</b>
<b>CFO</b>	1.00	0.50	0.47	-0.32	0.67	0.60
<b>Ind CFO Mean</b>		1.00	0.90	-0.70	0.42	0.40
<b>Ind CFO Median</b>			1.00	-0.47	0.40	0.38
<b>Ind CFO Std</b>				1.00	-0.29	-0.28
<b>CFO_L1</b>					1.00	0.66
<b>CFO_L2</b>						1.00

Table 3.10: This table shows the Pearson correlations for CFO features in our sample.

	1 Lag	2 Lags
ROE	0.16	0.15
ROA	0.28	0.27
RNOA	0.06	0.04
FCF	0.22	0.21
CFO	0.23	0.23

Table 3.11: This table shows the Autocorrelations in Profitability Measures.

<b>Analysis (Degree of Winsorization)</b>	<b>1-1 (None)</b>	<b>1-2 (1% per side)</b>	<b>1-3 (2.5% per side)</b>
ROE	3	2	3
ROA	2	2	2
RNOA	4	2	4
CFO	2	2	2
FCF	2	2	2

Table 3.12: This table shows the maximum depth selected for each analysis and measure combination with  $B = 30$ .

Return on Equity (ROE)			
	Analysis (Degree of Winsorization)		
	1-1 (None)	1-2 (1% per side)	1-3 (2.5% per side)
Mean Training Accuracy	61.1%	61.4%	61.1%
Mean Cross Validation Accuracy	55.5%	58.5%	55.5%
Mean Testing Accuracy	55.2%	56.8%	55.2%
Random Walk 1		50%	
Random Walk 2		48%	
Test Scores by Year			
2012	57.5%	60.4%	57.5%
2013	53.8%	57.1%	53.8%
2014	55.2%	54.2%	55.2%
2015	54.1%	55.4%	54.1%
Feature Importance			
PERMNO	2.7%	2.7%	2.7%
Fisical Year	19%	30%	19%
ROE	78.2%	67.2%	78.2%

Table 3.13: This table shows the results of analyses 1-1, 1-2, and 1-3 results for the profitability measure ROE.

Return on Equity (ROA)			
	Analysis (Degree of Winsorization)		
	1-1 (None)	1-2 (1% per side)	1-3 (2.5% per side)
Mean Training Accuracy	61.1%	61.2%	61.1%
Mean Cross Validation Accuracy	57.9%	58.1%	57.9%
Mean Testing Accuracy	58.5%	58.5%	58.5%
Random Walk 1		50%	
Random Walk 2		48.1%	
Test Scores by Year			
2012	60.6%	60.4%	60.6%
2013	57.6%	58.2%	57.6%
2014	58.9%	58%	58.9%
2015	56.8%	57.2%	56.8%
Feature Importance			
PERMNO	2.6%	2.5%	2.6%
Fisical Year	32.2%	32.2%	32.2%
ROA	65.2%	65.2%	65.2%

Table 3.14: This table shows the results of analyses 1-1, 1-2, and 1-3 results for the profitability measure ROA.

### Return on Net Operating Assets (RNOA)

		Analysis (Degree of Winsorization)		
		1-1 (None)	1-2 (1% per side)	1-3 (2.5% per side)
Mean Training Accuracy	61.2%	62.2%	61.2%	
Mean Cross Validation Accuracy	56.5%	59.8%	56.5%	
Mean Testing Accuracy	54.1%	53.3%	54.1%	
Random Walk 1			50%	
Random Walk 2			47.6%	
Test Scores by Year				
2012	59%	55.6%	59%	
2013	50%	53.3%	50%	
2014	52.7%	50.1%	53%	
2015	54.1%	54.3%	54.2%	
Feature Importance				
PERMNO	3.2%	1.8%	3.3%	
Fisical Year	16.2%	28.4%	16.2%	
RNOA	80.5%	69.8%	80.5%	

Table 3.15: This table shows the results of analyses 1-1, 1-2, and 1-3 results for the profitability measure RNOA.

Cash Flow from Operations (CFO)			
Analysis (Degree of Winsorization)			
	1-1 (None)	1-2 (1% per side)	1-3 (2.5% per side)
Mean Training Accuracy	66.5%	66.5%	66.5%
Mean Cross Validation Accuracy	63%	63.1%	0.63%
Mean Testing Accuracy	58.7%	58.5%	58.6%
Random Walk 1		50%	
Random Walk 2		38.2%	
Test Scores by Year			
2012	62.5%	62.2%	62.5%
2013	55.2%	54.9%	55.2%
2014	59%	58.9%	59%
2015	57.9%	0.58%	57.8%
Feature Importance			
PERMNO	7.6%	7.6%	7.6%
Fisical Year	19.6%	19.6%	19.5%
CFO	72.8%	72.8%	72.8%

Table 3.16: This table shows the results of analyses 1-1, 1-2, and 1-3 results for the profitability measure CFO.

Free Cash Flow (FCF)			
	Analysis (Degree of Winsorization)		
	1-1 (None)	1-2 (1% per side)	1-3 (2.5% per side)
Mean Training Accuracy	68.3%	68.4%	68.3%
Mean Cross Validation Accuracy	64.7%	64.8%	64.7%
Mean Testing Accuracy	59.1%	59.1%	59.1%
Random Walk 1		50%	
Random Walk 2		39%	
Test Scores by Year			
2012	63.2%	63.2%	63%
2013	55.6%	55.7%	55.6%
2014	60%	60%	59.8%
2015	57.9%	57.6%	57.9%
Feature Importance			
PERMNO	3.7%	3.7%	3.7%
Fisical Year	20.4%	20.4%	20.5%
FCF	75.8%	75.8%	75.8%

Table 3.17: This table shows the results of analyses 1-1, 1-2, and 1-3 results for the profitability measure FCF.

		Return on Equity (ROE)			
		Analysis			
		2-1	2-2	2-3	2-4
Mean Train Accuracy		60.5%	62.4%	62.4%	62.5%
Mean Cross Validation Accuracy		55.4%	56.0%	57.2%	56.4%
Mean Test Accuracy		55.0%	56.6%	56.8%	56.5%
Random Walk 1		50%			
Random Walk 2		48%			
Test Scores by Year					
2012		57.2%	58.2%	60.0%	60.0%
2013		53.2%	55.4%	54.8%	55.3%
2014		55.4%	57.4%	56.5%	54.2%
2015		53.9%	55.5%	55.9%	56.4%
Feature Importance					
PERMNO		10.8%	2.8%	2.8%	2.1%
FYEAR		32.8%	13.5%	8.8%	10.8%
<i>ROE</i>		50.1%	68.7%	61.1%	52.5%
SIC		6.3%	1.2%	1.0%	1.0%
Industry ROE Mean			5.2%	5.0%	4.3%
Industry ROE Median			5.3%	3.3%	3.0%
Industry ROE std.			3.3%	2.6%	2.5%
<i>ROE</i> <sub>t-1</sub>				15.4%	18.7%
<i>ROE</i> <sub>t-2</sub>					5.4%

Table 3.18: This table shows the results of analyses 2-1, 2-2, 2-3, and 2-4 for the profitability measure ROE.

		Return on Assets (ROA)			
		Analysis			
		2-1	2-2	2-3	2-4
Mean Training Accuracy		62.2%	61.9%	61.4%	61.4%
Mean Cross Validation Accuracy		59.0%	58.3%	58.7%	57.3%
Mean Testing Accuracy		58.6%	58.1%	58.1%	58.3%
Random Walk 1		50%			
Random Walk 2		48.1%			
Test Scores by Year					
2012		60.8%	60.7%	60.4%	60.5%
2013		58.2%	58.5%	57.8%	57.7%
2014		58.5%	57.6%	58.3%	58.4%
2015		56.8%	55.7%	55.9%	56.8%
Feature Importance					
PERMNO		4.0%	1.1%	0.8%	0.7%
FYEAR		15.9%	13.3%	5.0%	7.4%
ROA		78.3%	70.1%	66.7%	56.0%
SIC		1.8%	1.1%	0.8%	0.7%
Industry ROA Mean			5.0%	3.6%	3.5%
Industry ROA Median			7.4%	5.6%	4.0%
Industry ROA std.			2.1%	1.0%	0.8%
ROA <sub>t-1</sub>				16.5%	21.3%
ROA <sub>t-2</sub>					5.7%

Table 3.19: This table shows the results of analyses 2-1, 2-2, 2-3, and 2-4 for the profitability measure ROA.

Return on Net Operating Assets (RNOA)				
	Analysis			
	2-1	2-2	2-3	2-4
Mean Train Accuracy	60.8%	60.8%	61.9%	61.9%
Mean Cross Validation Accuracy	55.9%	56.3%	57.1%	56.6%
Mean Testing Accuracy	52.4%	52.7%	52.5%	52.8%
Random Walk 1	50%			
Random Walk 2	47.6%			
Test Scores by Year				
2012	59.0%	57.0%	57.4%	57.8%
2013	49.2%	49.8%	52.2%	52.1%
2014	47.5%	49.6%	46.9%	47.8%
2015	53.8%	54.4%	53.5%	53.4%
Feature Importance				
PERMNO	3.6%	1.0%	1.4%	0.9%
FYEAR	15.3%	12.2%	6.7%	7.8%
RNOA	79.6%	63.0%	61.7%	56.6%
SIC	1.5%	0.6%	0.9%	0.8%
Industry RNOA Mean		6.7%	4.5%	4.1%
Industry RNOA Median		14.3%	6.2%	6.8%
Industry RNOA std.		2.1%	2.7%	2.2%
RNOA <sub>t-1</sub>			15.9%	16.2%
RNOA <sub>t-2</sub>				4.5%

Table 3.20: This table shows the results of analyses 2-1, 2-2, 2-3, and 2-4 for the profitability measure RNOA.

Cash Flow from Operations (CFO)				
	Analysis			
	2-1	2-2	2-3	2-4
Mean Training Accuracy	67.5%	67.8%	68.8%	69.2%
Mean Cross Validation Accuracy	63.9%	64.0%	64.5%	64.6%
Mean Testing Accuracy	60.2%	59.8%	60.2%	60.1%
Random Walk 1	50%			
Random Walk 2	38.2%			
Test Scores by Year				
2012	62.6%	63.2%	64.1%	63.2%
2013	62.6%	58.5%	57.5%	58.1%
2014	59.0%	60.1%	60.4%	60.3%
2015	56.5%	57.3%	58.7%	58.7%
Feature Importance				
PERMNO	2.4%	1.6%	1.8%	1.4%
FYEAR	4.9%	3.5%	3.2%	2.8%
<i>CFO</i>	91.1%	82.1%	75.8%	75.4%
SIC	1.6%	1.7%	1.6%	1.4%
Industry CFO Mean		3.8%	4.5%	3.7%
Industry CFO Median		5.4%	4.4%	4.5%
Industry CFO std.		1.9%	2.1%	1.9%
<i>CFO</i> <sub>t-1</sub>			6.6%	5.1%
<i>CFO</i> <sub>t-2</sub>				3.8%

Table 3.21: This table shows the results of analyses 2-1, 2-2, 2-3, and 2-4 for the profitability measure CFO.

Free Cash Flow (FCF)				
	Analysis			
	2-1	2-2	2-3	2-4
Mean Training Accuracy	69.2%	69.5%	69.5%	69.6%
Mean Cross Validation Accuracy	65.2%	65.2%	65.5%	65.4%
Mean Testing Accuracy	60.8%	60.7%	61.3%	61.7%
Random Walk 1	50%			
Random Walk 2	39%			
Test Scores by Year				
2012	62.3%	64.1%	62.2%	63.9%
2013	63.4%	60.4%	63.9%	63.3%
2014	60.6%	60.6%	61.1%	61.5%
2015	56.9%	57.8%	57.9%	58.2%
Feature Importance				
PERMNO	2.1%	1.0%	1.0%	0.9%
FYEAR	5.0%	4.1%	2.5%	2.4%
<i>FCF</i>	91.5%	82.1%	82.2%	81.1%
SIC	1.3%	1.1%	1.0%	0.8%
Industry FCF Mean		3.9%	4.4%	2.8%
Industry FCF Median		6.0%	3.8%	5.0%
Industry FCF std.		1.9%	1.6%	1.6%
<i>FCF</i> <sub>t-1</sub>			3.6%	3.2%
<i>FCF</i> <sub>t-2</sub>				2.1%

Table 3.22: This table shows the results of analyses 2-1, 2-2, 2-3, and 2-4 results for the profitability measure FCF.

# Chapter 4

## Information Transfer Estimation on Large Data

### 4.1 Introduction

Transfer entropy (TE) is an information measure that quantifies information transfer between processes evolving in time (see Chapter 1.3.1). Transfer entropy has many potential applications in canonical systems, neuroscience, social media, and financial markets. Prior work concerning financial applications has typically used long windows in their research. For example, Marschinski and Kantz (2002) measured information transfer between two financial time series (the DAX stock index and the Dow Jones) to determine to what extent one index determined another's behavior. The data used by Marschinski and Kantz (2002) sampled data every minute between May 2000 and June 2001; after cleaning the data, 63,867 observations were used in Marschinski and Kantz (2002)'s study. More recently Sandoval (2014) used daily price data from 2003 to 2012 to detect causal relationships between 197 largest firms (globally) with Transfer Entropy.

Given the assumption in Chapter 1.3 that information is reflected in price and in an efficient market changes rapidly (if not instantaneously), then slower frequencies of observations (or long time windows) such as monthly, daily, hourly, or by the minute may be

---

This chapter contains material from the following publication:

- K. M. Ikegwu, J. Trauger, J. McMullin, and R. J. Brunner. Pyif: A fast and light weight implementation to estimate bivariate transfer entropy for big data. In *2020 SoutheastCon*, pages 1–6, 2020. doi: 10.1109/SoutheastCon44009.2020.9249650

insufficient in capturing the price change process. Given that we may not capture the entire price change process with long time windows between each observation, we examine price change with shorter windows. However, shorter time windows will have a finer resolution of data, requiring more observations in a dataset.

In figure 4.1 consider that we are looking at a univariate dataset for 30 days. If the frequency of observations (time windows) in 30 days is daily, there are 30 observations in the dataset. If the time windows are hourly, we can expect 720 observations in 30 days, and if the time window between observations occurs on a minute level, then we have 43,200 observations. If we continue these calculations down to the 1-second level, then there are 2,592,000 observations in a 30-day window. We find that existing open-source implementations are not suited to estimate information flow for datasets with large observations. Given the need to examine information transfers on finer resolutions of data, we propose an open-source software implementation to estimate TE on large datasets.

## 4.2 Estimating Transfer Entropy

Prior research has developed a few approaches for estimating TE. There is a requirement to specify the following parameter choices: the time between periods, length of the time series, number of past observations that inform the future observations, and the direction of information transfer in the approaches. The correct approach depends on the underlying data. There are many techniques for estimating mutual information. Khan et al. (2007) explored the utility of different methods for mutual information estimation, and many of the methods are applicable to estimate TE.

### 4.2.1 Kraskov Estimator

Kraskov et al. (2004) outlined a widely used method to estimate Transfer Entropy with k-nearest neighbors. Note to estimate entropy:

$$\hat{H}(X) = -\frac{1}{n} \sum_{i=1}^n \ln(\hat{p}(x_i)) \quad (4.1)$$

Kraskov et al. expanded this definition to estimate entropy to:

$$\begin{aligned} \hat{H}(X) = & -\frac{1}{n} \sum_{i=1}^n \psi(n_x(i)) - \frac{1}{k} + \psi(n) + \ln(c_{d_x}) + \\ & \frac{d_x}{n} \sum_{i=1}^n \ln(\epsilon(i)) \end{aligned} \quad (4.2)$$

where  $n$  are the number of data points,  $k$  is a parameter to select the amount of the nearest neighbors,  $d_x$  is the dimension of  $x$ , and  $c_{d_x}$  is the volume of the  $d_x$ -dimensional unit ball. For two random variables  $X$  and  $Y$ , let  $\frac{\epsilon(i)}{2}$  be the distance between  $(x_i, y_i)$  and it's  $k^{\text{th}}$  neighbor be denoted by  $(kx_i, ky_i)$ . Let  $\frac{\epsilon_x(i)}{2}$  and  $\frac{\epsilon_y(i)}{2}$  be defined as  $\|x_i - kx_i\|$  and  $\|y_i - ky_i\|$  respectively.  $n_x(i)$  is the number of points  $x_j$  such that  $\|x_i - x_j\| \leq \epsilon_x(i)/2$ ,  $\psi(x)$  is the digamma function where:

$$\psi(x) = \Gamma(x)^{-1} \frac{d\Gamma(x)}{dx} \quad (4.3)$$

and  $\Gamma(x)$  is the ordinary gamma function. Lastly  $\psi(1) = -C$  where  $C = 0.5772156649$  and is the Euler-Mascheroni constant.

To estimate Joint entropy between  $X$  and  $Y$ :

$$\begin{aligned} \hat{H}(X, Y) = & -\psi(k) - \frac{1}{k} + \psi(n) + \ln(c_{d_x} c_{d_y}) + \\ & \frac{d_x + d_y}{n} \sum_i^n \ln(\epsilon(i)) \end{aligned} \quad (4.4)$$

where  $d_y$  is the dimension of  $y$ , and  $c_{d_y}$  is the column of the  $d_y$ -dimensional unit ball. Using  $\hat{H}(X)$ ,  $\hat{H}(Y)$ , and  $H(\hat{X}, Y)$  mutual information can be estimated as:

$$\hat{I}(X, Y) = \psi(k) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n [\psi(n_x(i)) + \psi(n_y(i))] + \psi(n) \quad (4.5)$$

where  $n_y(i)$  is the number of points  $y_j$  such that  $\|y_i - y_j\| \leq \frac{\epsilon_y(i)}{2}$ .

The basic idea of mutual information estimation with the Kraskov's estimator is for each observation  $i$  to count the number  $n_x(i)$  (or  $n_y(i)$ ) of points less than the distance between  $x_i$  (or  $y_i$ ) and the  $k^{th}$  neighbor. The distance or  $\epsilon(i)$  fluctuates for each observation consequently  $n_x(i)$  and  $n_y(i)$  fluctuate. The counts are averaged over all observations which results in mutual information as defined in eq. 4.5.

To estimate TE with Kraskov's estimator recall from eq. 1.15 that TE is expressed as the difference between two mutual information computations. To estimate TE with eq. 1.15, the first mutual information is between  $Y$ 's future value,  $Y$ 's lagged values, and the  $X$ 's lagged values. The second mutual information is the mutual information between  $X$  and its lagged values. The difference between the first and second mutual information will compute TE. With an embedding of 1 TE can be estimated with:

$$TE_{Y \rightarrow X} = \left[ \psi(k) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n [\psi(n_{xy}^{(t-1)}(i)) + \psi(n_{xy}^{(t)}(i))] + \psi(n) \right] - \left[ \psi(k) - \frac{1}{k} - \frac{1}{n} \sum_{i=1}^n [\psi(n_x^{(t)}(i)) + \psi(n_x^{(t+1)}(i))] + \psi(n) \right] \quad (4.6)$$

Here  $n_{xy}$  is the number of points  $y_j$  such that  $\|y_i - y_j\| \leq \epsilon_y(i)/2$  and  $\|x_i - y_j\| \leq \epsilon_y(i)/2$ .

#### 4.2.2 Additional Estimators

Khan et al. (2007) also explored the utility of Kernel Density Estimation, Edgeworth approx-

imation of differential entropy to calculate Mutual Information, and adaptive partitioning of the XY plane to estimate the joint probability density, which can estimate mutual information. Ultimately Khan et al. (2007) found that a KDE estimator and Kraskov estimator outperform other methods for their ability to capture the dependence structure of random processes. Anderson and McMullin (2018) examined the properties of the Kraskov estimator between equities and their underlying options. They ultimately provided evidence that incorrect parameter choices lead to smaller TE estimates and that Kraskov’s method has a downward bias indicating that it underestimates information transfer. Nevertheless, it is relatively insensitive to parameter selection.

### 4.3 PyIF

PyIF<sup>1</sup> is our proposed software implementation to estimate bivariate TE on large data and is open-sourced. PyIF currently only supports using the Kraskov estimator to estimate TE. PyIF utilizes recent advancements in hardware to parallelize & optimize operations across CPUs and Cuda compatible GPUs (see NVIDIA et al. (2020)). Given the iterative nature of eq. 4.6 it is possible to speed up the estimation by performing computations in parallel. In particular we focus our efforts on the parallelization & optimization across operations to obtain  $n_x$  &  $n_{xy}$  in eq. 4.6 faster.

PyIF is a python implementation which utilizes 5 well-known and actively supported python libraries: SciPy (see Virtanen et al. (2020)), NumPy (see Harris et al. (2020)), scikit-learn (see Pedregosa et al. (2011)), nose (see Pellerin (2021)), and numba (see Lam et al. (2015)). SciPy is an open-source Python library used for a variety of STEM applications. NumPy is a part of SciPy’s ecosystem and is an open-source package that provides convenient ways to perform matrix manipulations and useful linear algebra capabilities. Scikit-learn is a popular open-source library for machine learning, and nose is another open-source library

---

<sup>1</sup>PyIF can freely be downloaded from: <https://github.com/lcdm-uiuc/PyIF>

that is useful for testing code to ensure that it will produce the correct outcome. Lastly, numba is a python compiler that can compile Python code to execute multicore CPUs and CUDA-capable GPUs.

PyIF’s interface requires one to supply  $X$  and  $Y$ , two NumPy arrays with  $N \times 1$  dimensions. PyIF supports optional arguments such as  $k$ , which controls the number of neighbors used in KD-tree nearest neighbor searches, *embedding* controls how many lagged periods are in the transfer entropy computation, and a boolean argument *GPU* allows one to use a CUDA compatible GPU for transfer entropy estimation. Lastly, *safetyCheck* is a boolean argument that can check for duplicates rows in your dataset. This boolean argument helps prevent a more subtle error when multiple data points in a bivariate dataset have identical coordinates. For duplicated observations, several points have an identical distance to a query point during  $k$  nearest neighbors search, which violates assumptions of the Kraskov estimator. A solution used in practice and that we recommend is to add a small amount of noise to the dataset to avoid this error.

## 4.4 Comparative Analysis

We compare PyIF’s ability to estimate Transfer Entropy against existing implementations with respect to computational performance. We present all of the data and code used to estimate TE for all implementations<sup>2</sup>. Each implementation in this comparative analysis estimates TE on four simulated bivariate datasets of different sizes. The estimated TE values are roughly the same for each implementation, and we forgo comparing the actual values since this is random simulated data. We assume that there is relatively little to no information transfer between the random processes. We run each of the implementations (excluding Transfer Entropy Toolbox) on nano, a cluster of eight SuperMicro servers with

---

<sup>2</sup>The data and code can freely be downloaded from:  
[https://github.com/lcdm-uiuc/Publications/tree/master/2020\\_Ikegwu\\_Traguer\\_McMullin\\_Brunner](https://github.com/lcdm-uiuc/Publications/tree/master/2020_Ikegwu_Traguer_McMullin_Brunner)

Intel Haswell/Broadwell CPUs and NVIDIA Tesla P100/V100 GPUs hosted by the National Center of Super Computing Applications at the University of Illinois at Urbana-Champaign. We used one node containing two E5-2620 v3 Intel Xeon CPUs and 2 NVIDIA P100 GPUs with 3584 cores. We refer to this as the first analysis.

We conduct the same analysis on different hardware to compare PyIF to Transfer Entropy toolbox because of MATLAB licensing issues with the National Center of Super Computing Applications. We use an Engineering Workstation with an Intel Xeon Processor E5-2680 v4 hosted by Engineering IT shared services at the University of Illinois at Urbana-Champaign. We use a single CPU core and up to 16GB of RAM to estimate TE with Transfer Entropy toolbox and PyIF. This workstation does not offer CUDA compatible GPUs for either PyIF or Transfer Entropy Toolbox, so we forgo comparing the GPU implementations. This workstation has a CPU time limit of 60 minutes, meaning that if any process uses 100% of a CPU core for more than 60 minutes, the process is terminated. We refer to this as the second analysis.

#### 4.4.1 IDTxl

The first implementation that we compare PyIF to is the Information Dynamics Toolkit xl (IDTxl). IDTxl is an open-source Python toolbox for network inference (see Wollstadt et al. (2019)). Currently, IDTxl relies on NumPy, SciPy, CFFI (which is another open source library that provides a C interface for Python code) (see Rigo and Fijalkowski (2018)), H5py which is a Python package that is used to interface with HDF5 binary data format (see Collette (2013)), JType (see Menard and Nell (2018)) which is a Python module that provides a Java interface for Python code, and Java jdk which is a developer kit to develop Java applications and applets. IDTxl has additional functionality besides estimating TE; however, we only use IDTxl's capability to estimate TE on a bi-variate dataset.

#### 4.4.2 TransEnt

TransEnt is a R package that estimates transfer entropy (see Mount et al. (2015)). Currently, TransEnt relies on Rcpp, which acts as an interface to C++ from R. TransEnt also relies on a C++ library called Approximate Nearest Neighbors (ANN) (see Arya and Mount (1998)), which performs exact and approximate nearest neighbor searches. Currently, the package is removed from CRAN. However, this software can be used and installed from Mount et al. (2015)'s github repo <sup>3</sup>.

#### 4.4.3 RTransferEntropy

RTransferEntropy is a R package that estimates transfer entropy between two time-series Simon et al. (2019). Currently, the RTransferEntropy package relies on Rcpp, and the future package supports performing computations in parallel to decrease the wall time. We include both the parallel implementation of RTransferEntropy and the default implementation for completeness in the results.

#### 4.4.4 Transfer Entropy Toolbox

Transfer Entropy Toolbox is an open-source MATLAB toolbox for transfer entropy estimation (see Lindner et al. (2011)). This code's dependencies include: the Statistics & Machine Learning toolbox, which provides functions to analyze and model data; the FieldTrip toolbox, which is used for EEG, iEEG, MEG, and NIRS analysis; the parallel computing toolbox that performs parallel computations of multicore CPUs and GPUs; the signal processing toolbox that provides functions to analyze, preprocess, and extract features from sampled signals; the TSTOOL toolbox which is a toolbox for nonlinear time series analysis. TSTOOL no longer exists and cannot be download from its official homepage <sup>4</sup>. Nevertheless, the Transfer Entropy toolbox developers include pre-compiled mex files of TSTOOL that will work with

---

<sup>3</sup>Mount et al. (2015)'s Github Repo: <https://github.com/Healthcast/TransEnt>

<sup>4</sup><http://www.dpi.physik.uni-goettingen.de/tstool/>

this implementation. At the time of writing, Transfer Entropy toolbox’s last update was in the year 2017.

#### 4.4.5 Data

We create four bivariate datasets for this comparative analysis. Each dataset contains two time series with randomly generated values between 0 and 1. The first dataset contains 1000 observations, the second dataset contains 10,000 observations, the third dataset contains 100,000 observations, and the fourth dataset contains 1,000,000 observations. We used the seed number 23 for the pseudo-random number generator for reproducibility. We will refer to the first dataset, second dataset, third dataset, and fourth dataset as the micro dataset, small dataset, medium dataset, and the large dataset respectively.

### 4.5 Results

We report the results for Analysis 1 in Tables 4.1, 4.2, 4.3, and 4.4. Each table contains the wall time to estimate TE using various implementations on the different data sets described in section 4.4. The higher the wall time, the longer it took for the specific implementation to estimate TE. The number in the relative performance column indicates how many times slower (or faster) PyIF (CPU) is to a particular implementation. After estimating TE using all of the implementations outlined in the comparative analysis section, we found that PyIF scales better on larger data.

Excluding the TransEnt implementation, the CPU implementation of PyIF (or PyIF (CPU)) takes less time to estimate TransferEntropy than all other implementations. The R package TransEnt has a better performance in terms of speed than PyIF (CPU) for the micro dataset and the small dataset. However, PyIF (CPU) can estimate transfer entropy in less time than all other implementations for the medium dataset and large dataset. PyIF (GPU) outperforms PyIF (CPU) for the small, medium, and large datasets. Figure 4.2 visualizes

this explanation. We suspect that the optimizations performed by Numba contribute to PyIF having a larger wall time than TransEnt on the micro and small datasets.

The results for Analysis 2 are in table 4.5. Although the Transfer Entropy Toolbox exceeds the CPU time limit for the large dataset, the results show that PyIF can scale better than Transfer Entropy Toolbox for the other three datasets. PyIF’s wall times are less than Transfer Entropy toolbox’s wall times, excluding the Micro Dataset. Figure 4.3 visualizes this explanation.

## 4.6 Conclusion

We address an important issue regarding large datasets with respect to estimating bi-variate TE. We introduce a fast solution to estimate TE with a small amount of dependencies. On large data, our implementation PyIF is up to 1072 times faster utilizing GPUs and up to 181 times faster utilizing CPUs than existing implementations that estimate bi-variate TE. PyIF is also open sourced and publicly available on github for anyone to use. For future work, we plan to improve the existing code base to increase the computational performance of PyIF even further. We plan to implement additional estimators outlined in section 4.2 to estimate bi-variate TE. This boost in computational performance will enable researchers to estimate bi-variate TE much faster for various research applications. In particular, PyIF will allow us to estimate information transfer between firms with small windows in a reasonable amount of time.

## 4.7 Figures and Tables

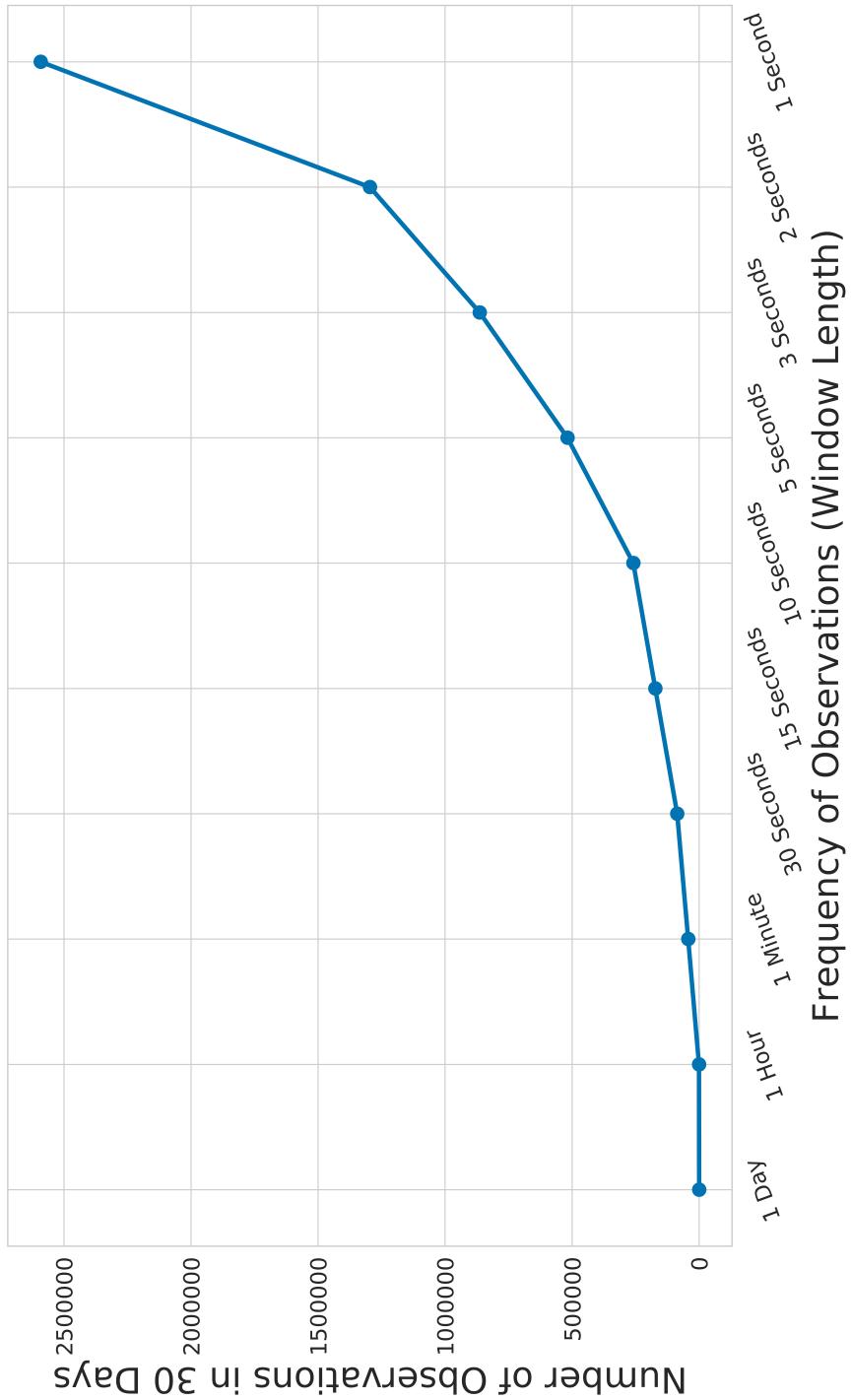


Figure 4.1: This figure shows the amount of observations for 30 days based on the frequency of observations.

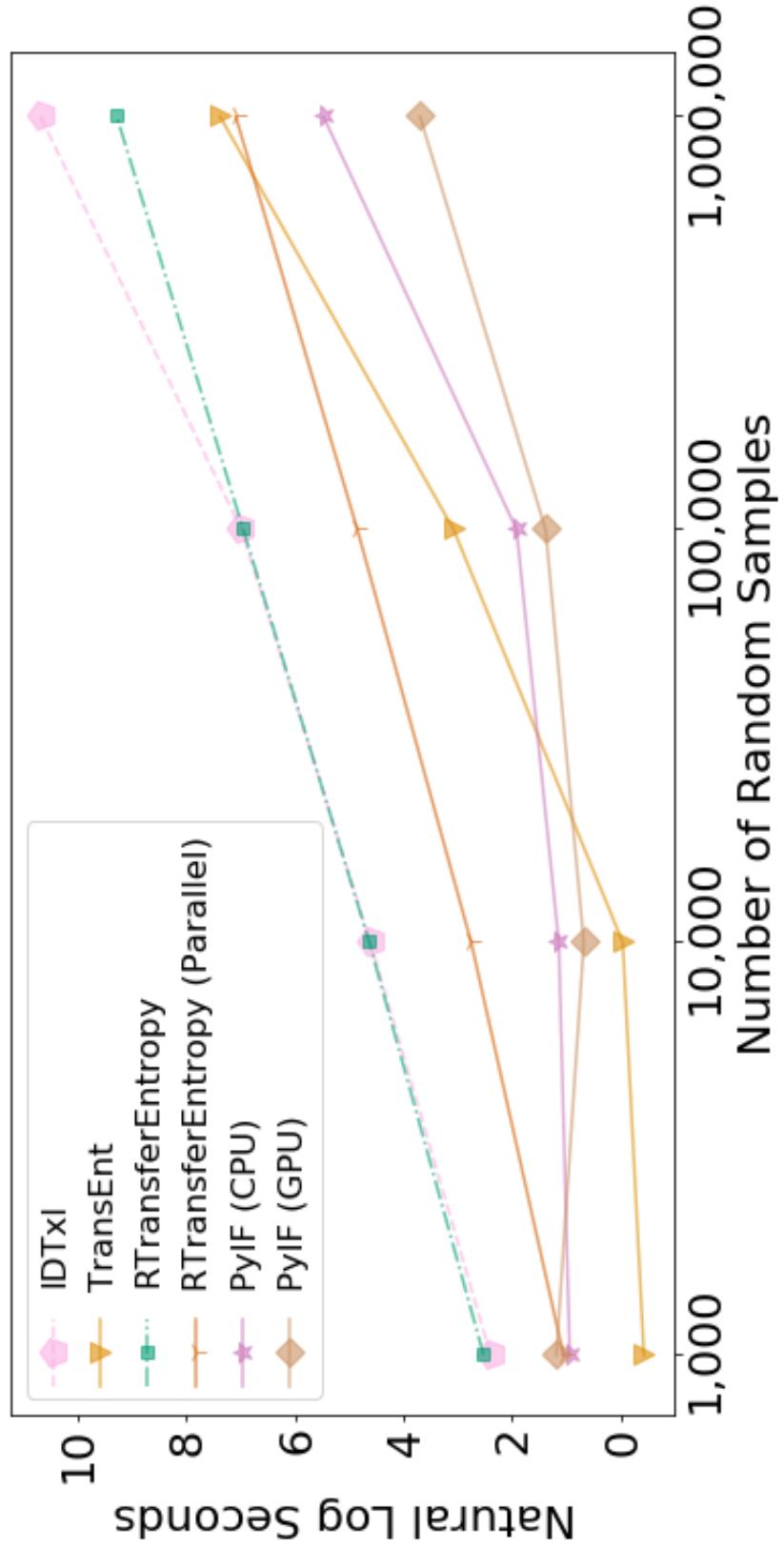


Figure 4.2: This figure shows the natural log time (in seconds) to estimate Transfer Entropy for each implementation (excluding Transfer Entropy Toolbox) for each dataset used in this study.

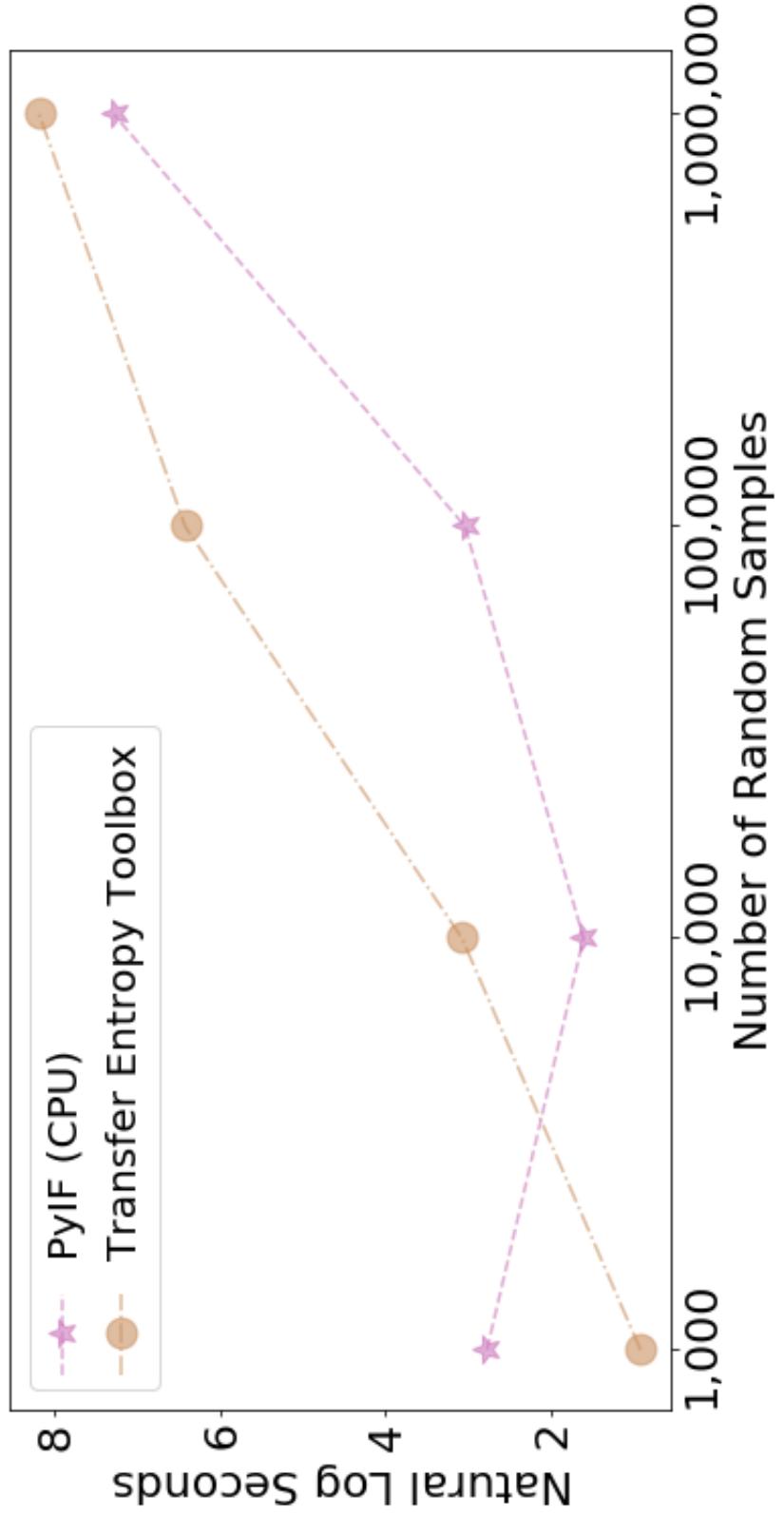


Figure 4.3: This figure shows the natural log time( in seconds) to estimate Transfer Entropy between PyIF and Transfer Entropy Toolbox on an Engineering Workstation described in the Comparative Analysis section. Transfer Entropy Toolbox exceeded the maximum allowable CPU runtime for the Large Dataset.

Implementation	Wall Time (in seconds)	Relative Performance to PyIF (CPU)
IDTxl	10.98	4.28
TransEnt	0.656	0.25
RTransferEntropy	12.492	4.87
RTransferEntropy (Parallel)	2.876	1.12
PyIF (CPU)	2.564	1
PyIF (GPU)	3.282	1.28

Table 4.1: Micro Data results for the first analysis.

Implementation	Wall Time (in seconds)	Relative Performance to PyIF (CPU)
IDTxl	100.23	31.94
TransEnt	0.968	0.308
RTransferEntropy	102.228	32.57
RTransferEntropy (Parallel)	15.703	5
PyIF (CPU)	3.138	1
PyIF (GPU)	1.98	0.63

Table 4.2: Small Data results for the first analysis.

<b>Implementation</b>	<b>Wall Time (in seconds)</b>	<b>Relative Performance to PyIF (CPU)</b>
IDTxl	1070.749	152.89
TransEnt	21.708	3.03
RTransferEntropy	1036.661	152
RTransferEntropy (Parallel)	127.281	18.66
PyIF (CPU)	6.82	1
PyIF (GPU)	3.996	0.58

Table 4.3: Medium Data results for the first analysis.

<b>Implementation</b>	<b>Wall Time (in seconds)</b>	<b>Relative Performance to PyIF (CPU)</b>
IDTxl	43150.129	181.97
TransEnt	1585.942	6.68
RTransferEntropy	10592.77	44.67
RTransferEntropy (Parallel)	1188.636	5.01
PyIF (CPU)	237.122	1
PyIF (GPU)	40.231	0.16

Table 4.4: Large Data results for the first analysis.

Implementation	Wall Time (in seconds)	Relative Performance to PyIF (CPU)
Micro Dataset Results (1000 Obs.)		
PyIF (CPU)	16.049	1.00
Transfer Entropy Toolbox	2.5012	0.15
Small Dataset Results (10,000 Obs.)		
PyIF (CPU)	4.989	1.00
Transfer Entropy Toolbox	21.6880	4.347
Medium Dataset Results (100,000 Obs.)		
PyIF (CPU)	20.915	1.00
Transfer Entropy Toolbox	616.8712	29.49
Large Dataset Results (1,00,000 Obs.)		
PyIF (CPU)	1455.725	1.00
Transfer Entropy Toolbox	> 3600	> 2.47

Table 4.5: Results for the second analysis.

# Chapter 5

## Cross-Firm Information Transfers During Earnings Season: A Network Approach

### 5.1 Introduction

Major public-firms announce their annual earnings during the first quarter of the year. The embedding information in these announcements effects the announcing firm and is incorporated in other firms. Existing literature in finance and accounting focuses on measuring information transfers as effects stemming from firm-specific information releases, such as earnings announcements (see Foster (1981)). This chapter's research objective is to understand how firm information releases influence other firms in the economy without prior assumptions. We avoid two assumptions when testing for informational links: an ex-ante source of fundamental linkage between firms and a particular information event that triggers an information transfer. On the first point, the bulk of studies assume that links between firms are persistent and readily observable by a characteristic such as a shared industry (see Foster (1981)) or a customer-supplier relationship (see Olsen and Dietrich (1985); R.Ahren and Harford (2014)).

Research by Billio et al. (2012) uses a network approach to examine cross-firm transfers implicit in monthly equity returns for large firms in the financial services sector. Consistent

---

This chapter contains material from the following working paper:

- Robert Brunner, Kelechi Ikegwu, Matthew Anderson, and Jeff McMullin. Cross-firm information transfers during earnings season: A network approach, April 2021

with a dynamic network of links among these financial firms, these authors find that the extent of cross-firm information transfers increases in recent decades and is associated with significant systemic risk in the financial sector, particularly around the recent 2007-2009 Great Recession. Evidence presented by Billio et al. (2012) raises the natural question of whether dynamic cross-firm information transfers are present in a broad-based sample of firms and what factors explain variation in this network.

We construct a novel network approach during Q1 2018 where the network relies on pairwise TE estimates between firms as a proxy for information transfer. Our approach contributes to extant accounting and finance literature in several ways. First, we provide novel evidence on the network of cross-firm links implied by the lead-lag structure implicit in high-frequency U.S. equity prices observed around releases of earnings information. By tracing the network of information transfers directly, we can document several features of cross-firm links that go beyond industry. Last, we introduce a novel way to identify central, or “bellwether”, firms in the equity market. We hope to encourage future work designed to broaden our understanding of the complex system of information flows inherent in modern equity markets.

For the remainder of the paper, we discuss the datasets used in this study. We discuss the methodology of estimating cross-firm information transfers and discuss the network science used to create dynamic cross-firm information transfer networks. We make observations about the structural properties of the networks. We then conduct an analysis to determine the effect of earnings surprise on dynamic cross-firm information transfers, and we discuss the results of the analysis.

## 5.2 Data

The data in the study comes from Wharton Research Data Services (WRDS). We obtain security prices from the Trade and Quote (TAQ) dataset, which contains all trades and

quotes that occurred at a sub-second level. The national best bid/offer (NBBO) price is determined by the lowest ask price and highest bid price available from multiple exchanges. The NBBO prices are identified from the TAQ dataset and are sampled at particular rates.

We obtain the NBBO price data for firms in the S&P 500 in Q1 2018. We construct ten datasets for each firm at different sampling rates. The sampling rates used in this study are: 1 second, 2 seconds, 3 seconds, 4 seconds, 5 seconds, 10 seconds, 15 seconds, 30 seconds, 60 seconds, and 120 seconds. For example, at 1 second sampling rate, there will be 120 NBBO price observations per minute for a particular firm. CRSP and Compustat are used to obtain firm-specific information. We also use the Institutional Brokers' Estimate System (IBES) to obtain earnings surprise for firms in the S&P 500 during Q1 2018. Throughout the text, we describe how variables are used from IBES, CRSP, and Compustat.

### 5.3 Estimating Dynamic Information Transfer Between Firms

TE estimates the amount of information transfer between every unique pair of firms in the S&P 500 during Q1 2018 (the first 61 days of 2018). To detect information transfer at the minute/sub-minute level between firms in the S&P 500, first, the NBBO prices will be obtained from the TAQ data for all firms in the S&P 500 in Q1 2018. We create ten datasets from the NBBO data for each date and firm for that particular date and firm, which samples the NBBO price data with the following sampling rates: 1, 2, 3, 5, 6, 10, 30, 60, and 120 seconds.

The idea behind sampling the NBBO price data is to help determine at what speed information transfer peaks. Next for each date, firm, and sampling rate bi-variate TE will be estimated with the NBBO sampled price data using three 130 minute windows throughout the trading day; the rationale behind this is to determine if there are any consistent information flow patterns during the beginning (9:30am-11:40am), middle (11:40am-1:50pm), or

end (1:50pm-4pm) of the trading day. To summarize TE will be estimated to a specific firm from all other firms in the S&P 500 for that particular date, sampling rate, and interval. Given that there are 9 sampling rates, 3 windows, 61 days, and 499 other firms, this amounts to 821,853 TE calculations per firm for Q1 2018.

Data from small sample rates will yield a high number of observations. Given that the open-source implementations were not suited to estimate bi-variate TE on big data (see Chapter 4) we utilize PyIF (see Ikegwu et al. (2020)). We compute TE estimates using the HAL cluster at the National Center for Supercomputing Applications using a single NVIDIA V100 GPU. The 1-second sampling rate has the most observations and is the most costly to estimate TE.

For a particular day, TE computations on HAL to 1 firm from the other 499 firms for the three intervals will take roughly 4 minutes. For all 500 firms, this took about 2000 minutes. For all days in Q1 2018, this took about 122,000 minutes (or about 85 days) at the 1-second sampling rate sequentially. Subsequently, for the 2-second sampling rate, the data is reduced by half. Ergo, the wall time was reduced by roughly half. Given that we can run up to five jobs in parallel on HAL, 1-second TE estimates took 17 days to compute.

### **5.3.1 Algorithmic Process for Estimating Information Transfer Between firms**

In this section, the algorithmic process is outlined to estimate information transfer between firms using the Q1 2018 data (see Algorithm 1). In Algorithm 1 we create a dictionary to map combinations of dates and sampling rates with matrices of computed information transfers. Given the dates of interest in line 2 and the sampling rates in line 3, we iterate through each date in line 4. For a particular date, we iterate through each sampling rate in line 5. For each unique combination of date and sampling rate, we select firms with observations in that date and sampling rate. We then create an empty matrix of size  $\text{Firms}_i^2$  by 3. Then

we create a counter variable to keep track of the observations in line 9.

Next, we iterate through each firm in the set of firms $_i$ . In algorithm 1 the trading day is split 3 into 2 hour and 10 minute windows which represents the beginning (9:30am-11:40am), middle (11:40am-1:50pm), and end of the trading day (1:50pm-4pm). In lines 10-12, we filter the NBBO prices for firm  $i$  to the beginning of the trading day, middle of the trading day, and end of the trading day for a single firm from the set of firm $_i$ . We repeat this process for all of the firms in the set firm $_j$ . We then estimate TE from firm $_j$  to firm $_i$ . The TE estimate is assigned to row  $\text{cnt}_{ij}$  and column 0 of the dateSR\_InfoFlow matrix for the particular key of the dictionary for the TE computation with the morning prices. Subsequently, the TE estimate is assigned to row  $\text{cnt}_{ij}$  and column 1 or 2 for the middle of the trading day or end of the trading day, respectively. We then increment  $\text{cnt}_{ij}$  by 1 and repeat this process for all pair of firms, sampling rates, and days.

---

**Algorithm 1:** Estimating Information Transfers Between Firms

---

```
1 dateSR_InfoFlow := { } ;
2 Dates := All Dates in Q1 2018 ;
3 SampleRates := [1, 2, 3, 4, 6, 10, 30, 60, 120] secs ;
4 for  $t \in Dates$  do
5   for  $SR \in SampleRates$  do
6     Firmsi := SelectFirms( $t, SR$ ) ;
7     Firmsj := SelectFirms( $t, SR$ ) ;
8     dateSR_InfoFlow[(t,SR)] := Matrix(length(Firmsi)2, 3) ;
9     cntij := 0 ;
10    for Firmi in Firmsi do
11      MoPricesi := Filter(Firmi, "9:30am-11:40am") ;
12      MidPricesi := Filter(Firmi, "11:40am-1:50pm") ;
13      AfterPricesi := Filter(Firmi, "1:50pm-4:00pm") ;
14      for Firmj in Firmsj do
15        MoPricesj := Filter(Firmj, "9:30am-11:40am") ;
16        MidPricesj := Filter(Firmj, "11:40am-1:50pm") ;
17        AfterPricesj := Filter(Firmj, "1:50pm-4:00pm") ;
18        dateSR_InfoFlow[(t,SR)][cntij,0] := ComputeTE(MoPricesi,
19          MoPricesj) ;
20        dateSR_InfoFlow[(t,SR)][cntij,1] := ComputeTE(MidPricesi,
21          MidPricesj) ;
22        dateSR_InfoFlow[(t,SR)][cntij,2] := ComputeTE(AfterPricesi,
23          AfterPricesj) ;
24        cntij += 1
25      end
26    end
27  end
28 end
29 end
```

---

### 5.3.2 Network Creation

After we compute all bi-variate TE estimates, we conduct exploratory data analysis to find patterns or exciting characteristics of the information transfer between firms to explore further. We employ network analysis techniques (see Chapter 1.2) to create an information transfer network for all firms at date  $t$ , sampling rate  $s$ , and window  $w$ . Figure 5.1 shows a single network for January 2nd, 2018, during the morning window (9:30am-11:40am). Each circle in the network is a node that represents a firm. The lines (or edges) between nodes have a thickness determined by the information transfer estimate via transfer entropy.

The information transfer network in figure 5.1 has been filtered with a Disparity Filter (see Serrano et al. (2009)) to reduce the amount of edges in a network. However, it is still challenging to gather insight from this network. In addition to this there's an information transfer network for each day  $t$ ,  $s$ , and  $w$ . Creating networks for all days, sampling rates, and windows will produce roughly 1,830 networks, which will make it more difficult to discover general insights from the data by viewing them. An alternative is to look at network measures to gain additional insights from the information transfer networks.

## 5.4 Exploratory Data Analysis

Figure 5.2 shows the distribution of information transfers for all days in Q1 2018 for all windows at each sampling rate. Each subplot is for a particular sample rate, the x-axis represents the bin value, and the y-axis represents the percentage of observations in a bin. Faster sampling rates have lower variance and higher means. Slower sampling rates have higher variance and smaller means. A possible explanation is that the frequency of observations decreases due to fewer non-overlapping observations during each 130-minute window. For example, there are 7,800 observations per firm at a 1-second sample rate and 130 observations per firm at a 1-minute sample rate. The 1-minute yields fewer data to compute TE and produces a noisier information transfer estimate.

Table 5.1 show additional summary statistics of the information transfers at each sampling rate. The 75th and 99th percentile values increase while the median percentile (and lower) values decay as the sample rate becomes slower. Table 5.2 contains the average TE values computed for each of the 130-minute trading windows during the first quarter of 2018 for pairs of firms in the S&P 500.

Across sampling rates, table 5.2 mean values exhibit similar patterns to table 5.1 with a decrease in means as the sampling rates become slower. Across the three 130-minute windows on each trading day, the 9:30 am-11:40 am (morning) window displays the highest average information transfer. This latter result is consistent with a surge in trading at the market open each day.

In table 5.3, we investigate the mean TE values at narrower time windows to determine when information transfers tend to peak during the trading day. Table 5.3 contains the average TE values computed for thirteen non-overlapping thirty-minute trading windows. Across all sampling rates, average TE values are highest during the first thirty minutes of the trading day. Mean TE values then decay throughout the trading day at all sampling rates. Sampling rates faster than 6 seconds have a slight burst of information transfer within the last hour of the trading day, consistent with information-based trading as part of the daily settlement (closing) trade.

Table 5.4 shows the average weighted out-degree of the information transfer network at various sampling rates across 61 morning trading periods (9:30am-11:40am). We scale each degree measure by the number of firms in the network at each measurement period. Since all firms in the information transfer networks are connected, the average weighted out-degree values across all firms in the first quarter of 2018 are equivalent to the mean weighted in-degree values. Stated differently, the average of all incoming information transfers from a set of firms to a particular firm is equivalent to the average of all outgoing connections from a set of firms to a firm. The average weighted incoming and outgoing information transfer spikes at a sampling rate of 5 – 6 seconds.

Tables 5.6 and 5.5 presents results of ordinary least squares models examining auto-correlation in daily measures of weighted network degree to determine whether firms' centrality in the network of information transfers displays persistence. Table 5.5 (table 5.6) presents auto-correlations for incoming (outgoing) information transfers based on the weighted in-degree (out-degree) computed from 9:30am-11:40am for each trading day in our 2018 sample period. The  $R^2$  values for both Panels exhibit similar behaviors across sampling rates. In particular, at faster sampling rates we see more explained variation from the morning information transfers and a gradual decay as the sample rates become slower. The  $R^2$  values for the auto-correlation models for outgoing information transfers during morning trading windows in table 5.6 are substantially higher than the incoming morning information transfers in table 5.5.

## 5.5 Earnings Surprise as a Determinant of Dynamic Information Transfers

Prior literature finds that more informative news events elicit more robust stock price responses among peer firms (see Foster (1981); Brochet et al. (2018)). Given this, we examine whether firms with larger absolute earnings surprises (measured relative to the consensus analyst forecast) display stronger information transfers in response to their earnings announcement. This analysis allows us to shed light on the dynamics of the information transfer network by focusing on variation in the timing of firm's earnings announcements during earnings season and on variation in the news contained in the announcement.

In particular, we estimate an ordinary least squares regression model of the form:

$$\begin{aligned}
Y_{it} = & \alpha + \beta_1 EA_{it} * Morning_{it} + \beta_2 EA_{it} * Morning_{it} * Abs\_Surprise_{it} + \\
& \beta_3 EA_{it} * Afternoon_{it} + \beta_4 EA_{it} * Afternoon_{it} * Abs\_Surprise_{it} + \\
& \beta_5 EA_{it} * Evening_{it} + \beta_6 Share\_Turnover_{it} + \beta_7 Abs\_RET_{it} + \\
& \beta_8 Morning_{it} + \beta_9 Evening_{it} + \epsilon_{it} \quad (5.1)
\end{aligned}$$

$Y_{it}$  is the weighted out-degree (or in-degree) for the  $i^{th}$  firm at the day  $t$ .  $EA_{it}$  is an indicator variable for the  $i^{th}$  firm at day  $t$  which is set equal to one for trading days with an earnings announcement for the  $i^{th}$  firm made after the prior market close or during the current trading day and to zero otherwise. Morning, Afternoon, and Evening are indicators variables set equal to one if the current observation is during the 9:30 am - 11:40 am, 11:40 am-1:50 pm, and 1:50 pm-4:00 pm trading windows, respectively.

$Abs\_Surprise$  is the absolute value of the quarterly earnings surprise measured as the difference between actual earnings in I/B/E/S and the last available consensus analyst earnings forecast from the I/B/E/S Summary file, scaled by the quarter-end stock price available from Compustat. The  $Abs\_Surprise$  variable is set to zero for all trading windows other than the first Morning, Afternoon, and Evening trading windows following the announcement.  $Share\_Turnover$  controls the number of shares traded during the trading day, scaled by the number of shares outstanding on CRSP.  $Abs\_RET$  and the absolute value of the close-to-close stock return from CRSP. Eq. 5.1 treats each 130-minute trading window for each firm as a distinct observation. This analysis excludes four outlying observations with absolute earnings surprises larger than 5% of price.

We present results of estimating eq. 5.1 in tables 5.7 and 5.8, where table 5.7 (table 5.8) presents results of OLS models where weighted out-degree (in-degree) is the dependent

variable. Results in table 5.7 show that a firm's weighted out-degree in the network of information transfers is significantly higher in the trading periods immediately following its earnings announcement when the firm announces a larger earnings surprise. In particular, the positive coefficient on the EA\*Morning\*Abs\_Surprise term across the sample rates in models (1) – (10) ranges from 1.672 in model (1) for the 1-second sample rate to 3.117 in model (4) for the 5-second sample rate. All estimates are statistically significant at a 5% two-tailed level. Further, the pattern in coefficients on the EA\*Morning\*Abs\_Surprise term across models (1) – (10) in Panel A suggests an information transfer that peaks at 5 seconds, consistent with our transfer entropy estimates presented in table 5.1. A significant relation between Wtd\_OutDegree and earnings news persists into the afternoon trading window at faster sample rates, with significant positive coefficients on the EA\*Afternoon\*Abs\_Surprise term for sample rates ranging from 1 second to 3 seconds in models (1) – (3) (two-tailed p-values < 0.05).

Turning to results in table 5.8 for weighted network in-degree shows that in contrast to results in table 5.7, Wtd\_InDegree displays limited variation with the magnitude of earnings news released. Except for the EA\*Morning\*Abs\_Surprise term in models (1) and (2) for the fastest sample rates, coefficients on the interaction terms with Abs\_Surprise are generally insignificant at conventional levels and display no clear pattern across sample rates. We interpret these results as suggesting that transfers out to other firms are significant in response to the magnitude of earnings news released, while transfers in are limited in response to the news released.

## 5.6 Community Detection

We ran community detection algorithms on the information transfer networks and found dynamic community formation across morning trading windows. In network science literature, a standard definition of a community is to have more nodes grouped where there is

a higher density of edges within a group than between groups. Given that the information transfer networks are relatively large, we utilize the Clauset-Newman-Moore greedy modularity maximization algorithm to find communities (see Clauset et al. (2004)). While the Clauset-Newman-Moore algorithm aims to find communities in vast networks, it cannot find communities in the information transfer networks.

The density in our information transfer networks are 1, whereas Clauset et al. (2004) benchmark network had a density of 0.0000125. Another way to think of this is that the ratio of edges to nodes in our network is 100 times greater than their benchmark network. This shortcoming creates the need for community detection algorithms that can scale to the information transfer networks density. With Clauset-Newman-Moore algorithm to combat the density issue, we reduce density (or the node to edges ratio). An option would be to filter out any edge in an information transfer network less than an arbitrary information transfer value. However, this would not preserve the network's structural properties and ignore small-scale interactions between firms. An alternative option and the approach used in this analysis is to apply a disparity filter (see Serrano et al. (2009)) to the information transfer networks. This method locally identifies statistically relevant weighted edges and can filter out relevant connections across all scales of interactions between firms.

We examine the two days with the most earnings announcements in our sample, February 1st, and January 31st. In addition to this, we include March 13th, a random trading day with no earnings announcements. We first take the average information transfer across 1 – 6 second sampling rates for the morning trading windows and reconstruct the information transfer networks. Then we apply the disparity filter to find the network backbone, which requires a value for  $\alpha$ . The selection of  $\alpha$  is more of an art than science. Table 5.9 shows the amount of nodes and edges remaining in a filtered network at a particular  $\alpha$  value. A smaller  $\alpha$  will yield an information transfer network that is too sparse to detect communities. A large  $\alpha$  will yield an information transfer network that is too dense to detect communities.

We select  $\alpha$  based on the Clauset-Newman-Moore algorithm's ability to detect communi-

ties for the most firms in the sample information transfer networks. We found that when  $\alpha = 0.375$ , we can detect communities for 100% of the firms on February 1st and 97% of the firms on January 31st and March 13th. We find fewer communities formed during announcement days (January 31st and February 1st) than the trading day with no announcements (March 13th). However, most of the community's sizes during the announcement days are larger than the communities in the March 13th information transfer network (see figure 5.3). For example, most of the firms in the March 13th information transfer network are in communities with 1, 2, 3, or 4 other firms. In contrast, during the announcement days it is less likely. In addition to this a firm could have up to 60 other firms in the same community during the announcement days and on the day with no announcements up to 39 other firms.

On the announcement days, the communities typically have only one firm that announced. There are three cases where a community has two firms that announced and one case where three firms announced. We see a diverse set of industries connected within and across the detected communities. For all the information transfer networks, there is a dynamic formation of hubs. Some hubs reoccur, and others form during a particular trading day. For example, in figures 5.4, 5.5, and 5.6 Amazon (AMZN), National Weather Service (NWS), and Mettler-Toledo International Inc (MTD) are hubs during the announcement days and the non-announcement day despite the different connections that form to each of these hubs. However, when Boeing (BA) announced on January 31st, it temporarily became a hub in 5.4 for that morning information transfer network. The networks' characteristics provide evidence that the networks are not random, which indicates that we are capturing valid cross-firm information transfers via our measurement of transfer entropies.

## 5.7 Conclusion

We introduce a new approach to examine information transfers during earnings season. While its known that earnings information produced by a single firm is incorporated into

the firm's equity price, we provide evidence that this information flows to the equity prices of other firms across industries. We study the network effect of earnings announcements by constructing daily networks of pairwise cross-firm information transfers. Our approach to construct this network relies on non-parametric estimates in equity prices with measures of transfer entropy drawn from information theory.

Our approach avoids standard assumptions in the existing literature on cross-firm information transfers that links between firms are persistent and form based on a readily observable characteristic, such as a shared industry. Our tests show that cross-firm links are substantially stronger for firms on days with releases of earnings information and firms with more unexpected earnings news, consistent with the network of cross-firm information transfers dynamically responding to shifts in the information landscape. Further, we find that communities form between firms with links that are not adequately captured by characteristics that are the focus of existing literature.

## 5.8 Figures and Tables

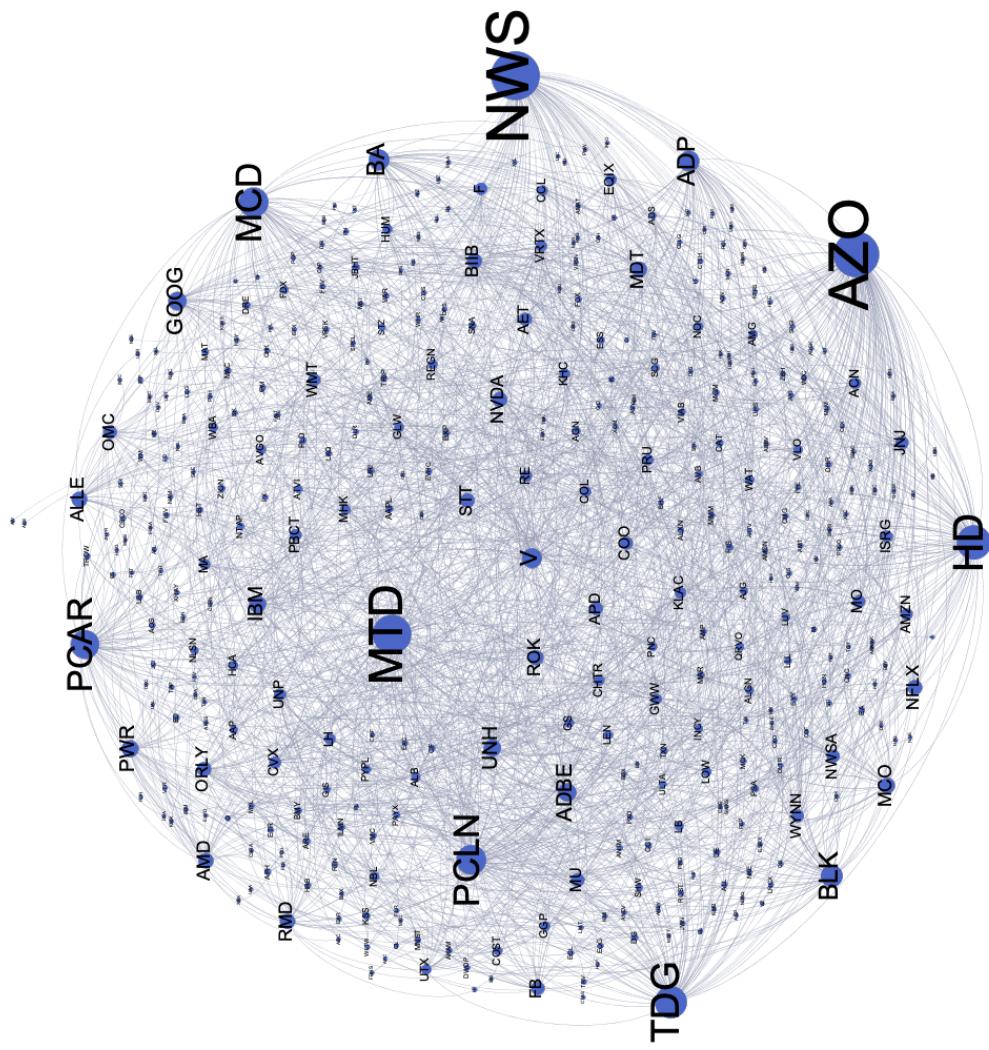


Figure 5.1: This figure shows an example network on January 2nd during the morning window (9:30 am-11:40 am) at the 1-second sampling rate.

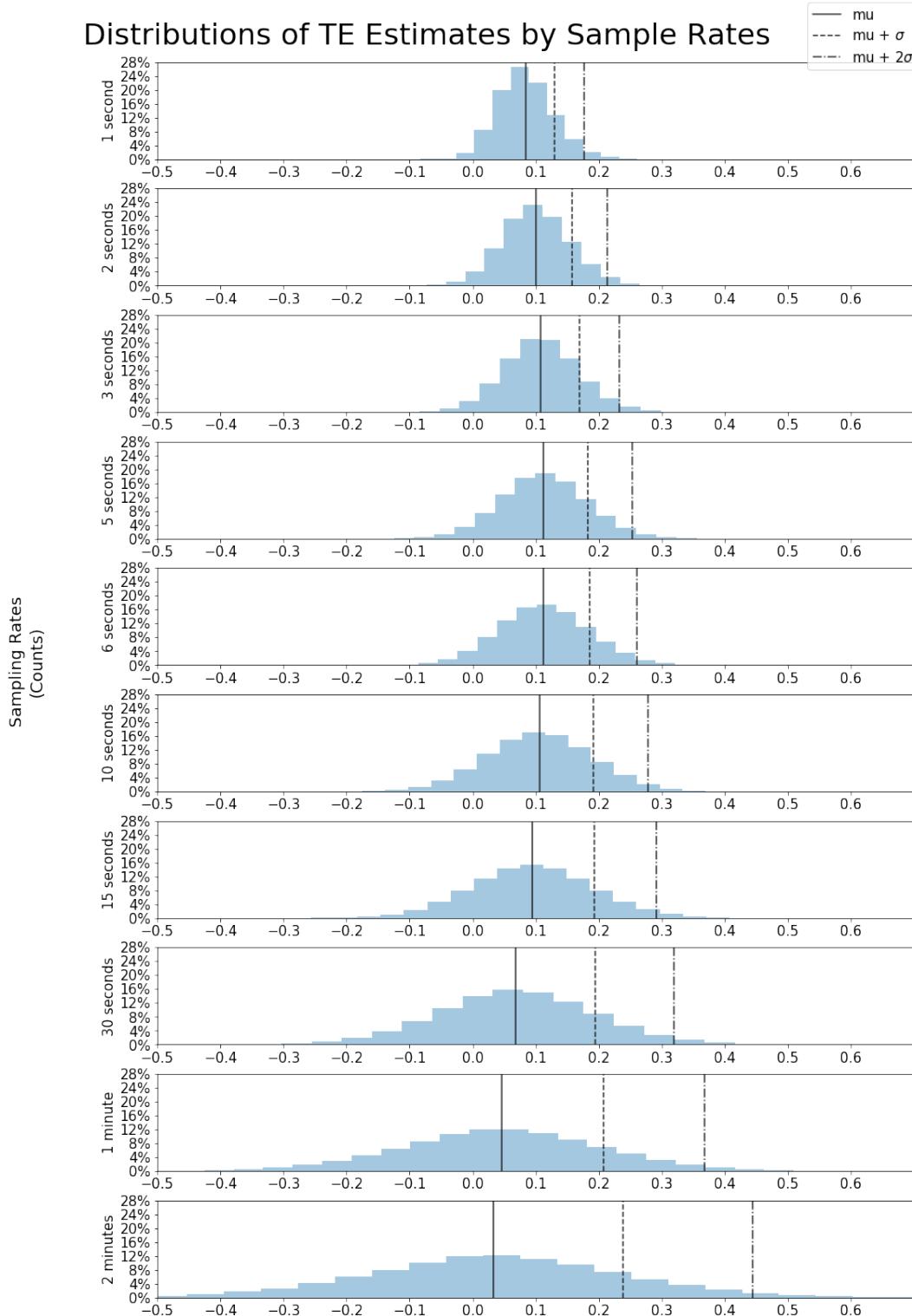


Figure 5.2: This figure shows the distributions of information transfers between firms during Q1 2018 for all sampling rates.

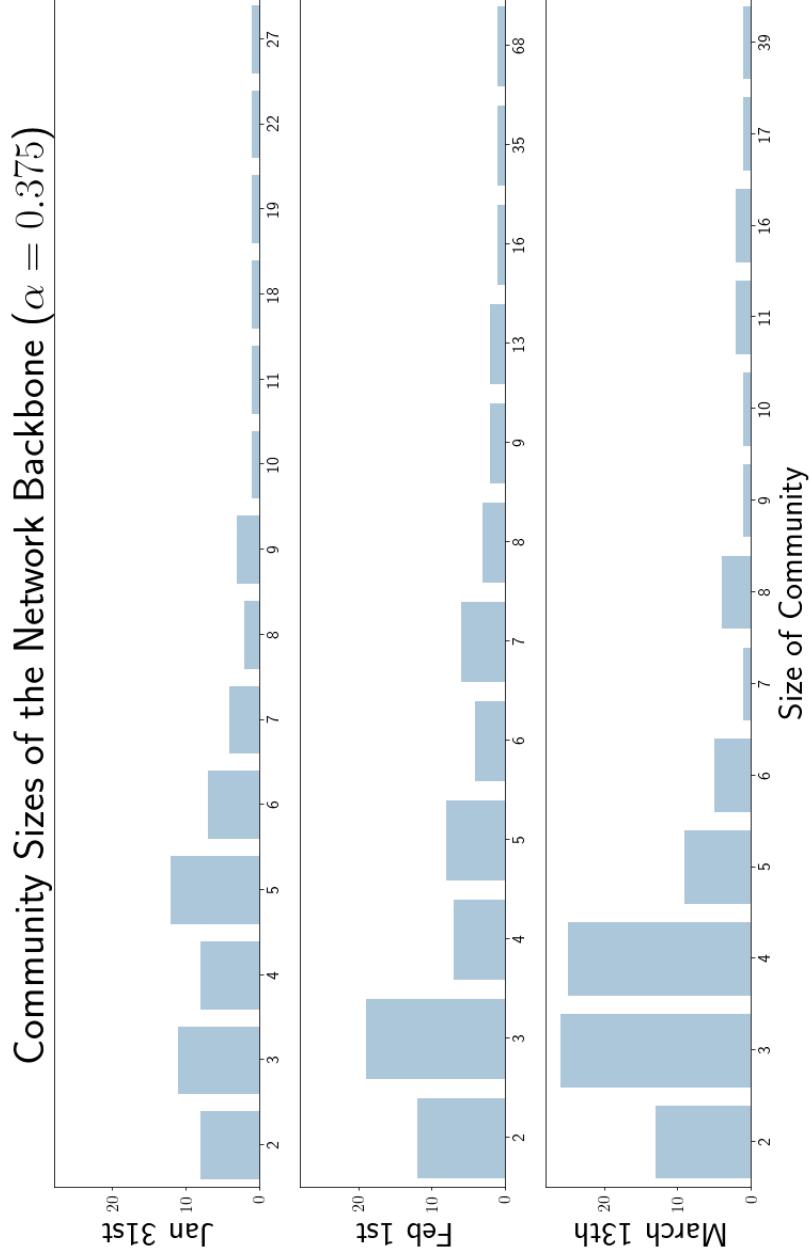


Figure 5.3: Community sizes of the network backbone when  $\alpha = 0.375$  for January 31st, February 1st, and March 13th

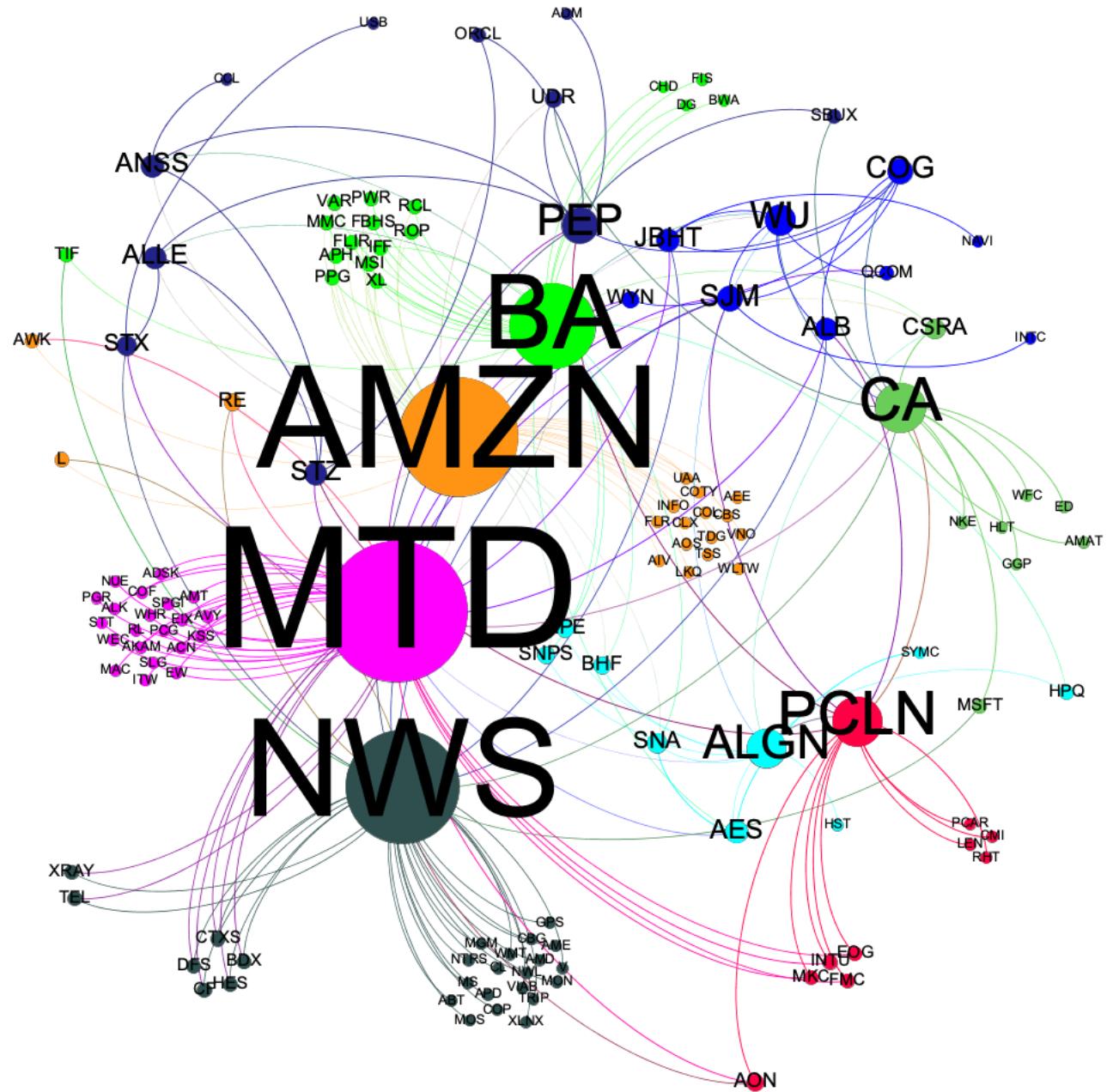


Figure 5.4: This figure shows the backbone network displaying communities for January 31st, where at least nine firms are in the community. Each color represents a specific community, and the node's size represents the degree of the node.

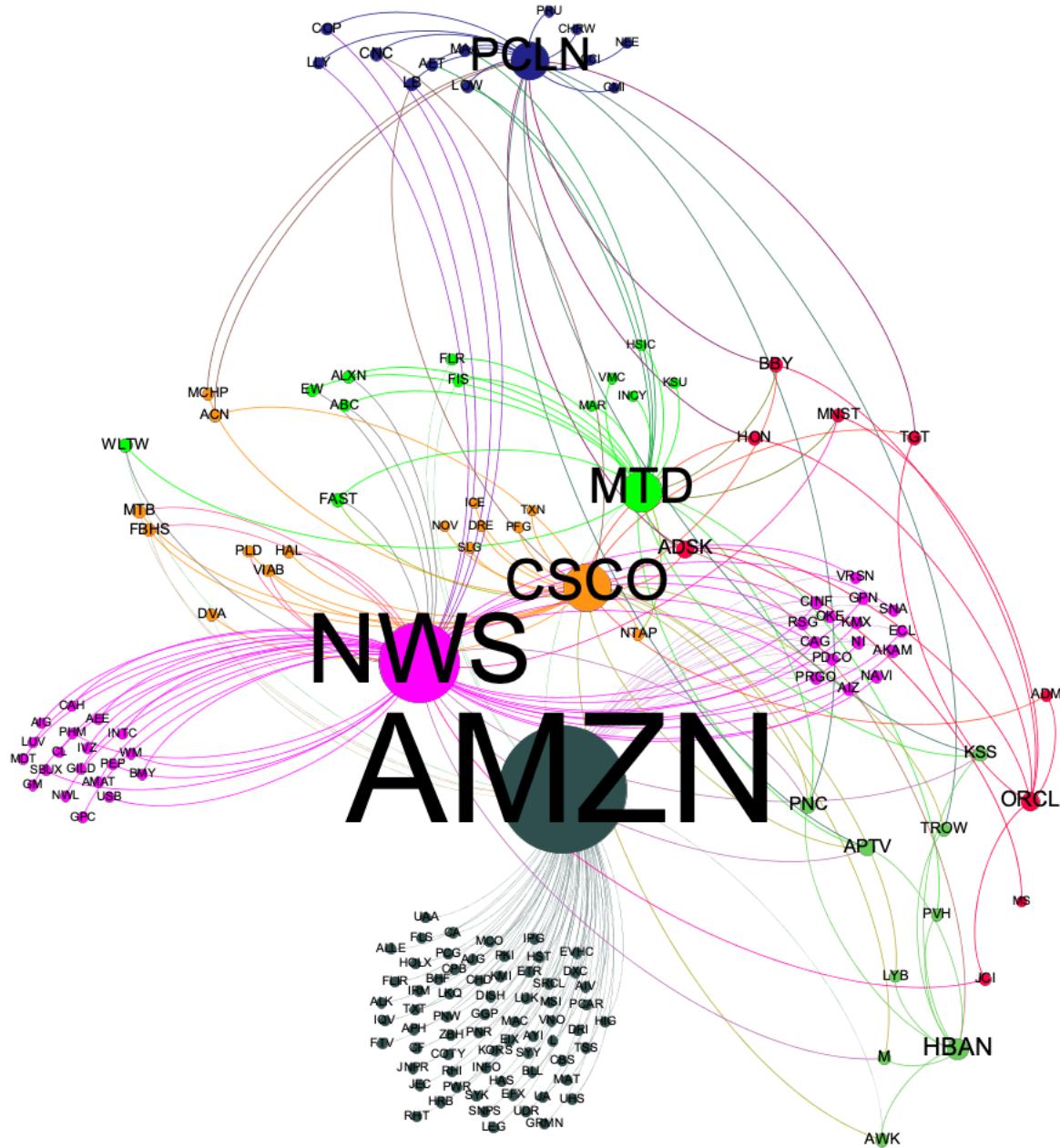


Figure 5.5: This figure shows the backbone network displaying communities for February 1st, where at least nine firms are in the community. Each color represents a specific community, and the node's size represents the degree of the node.

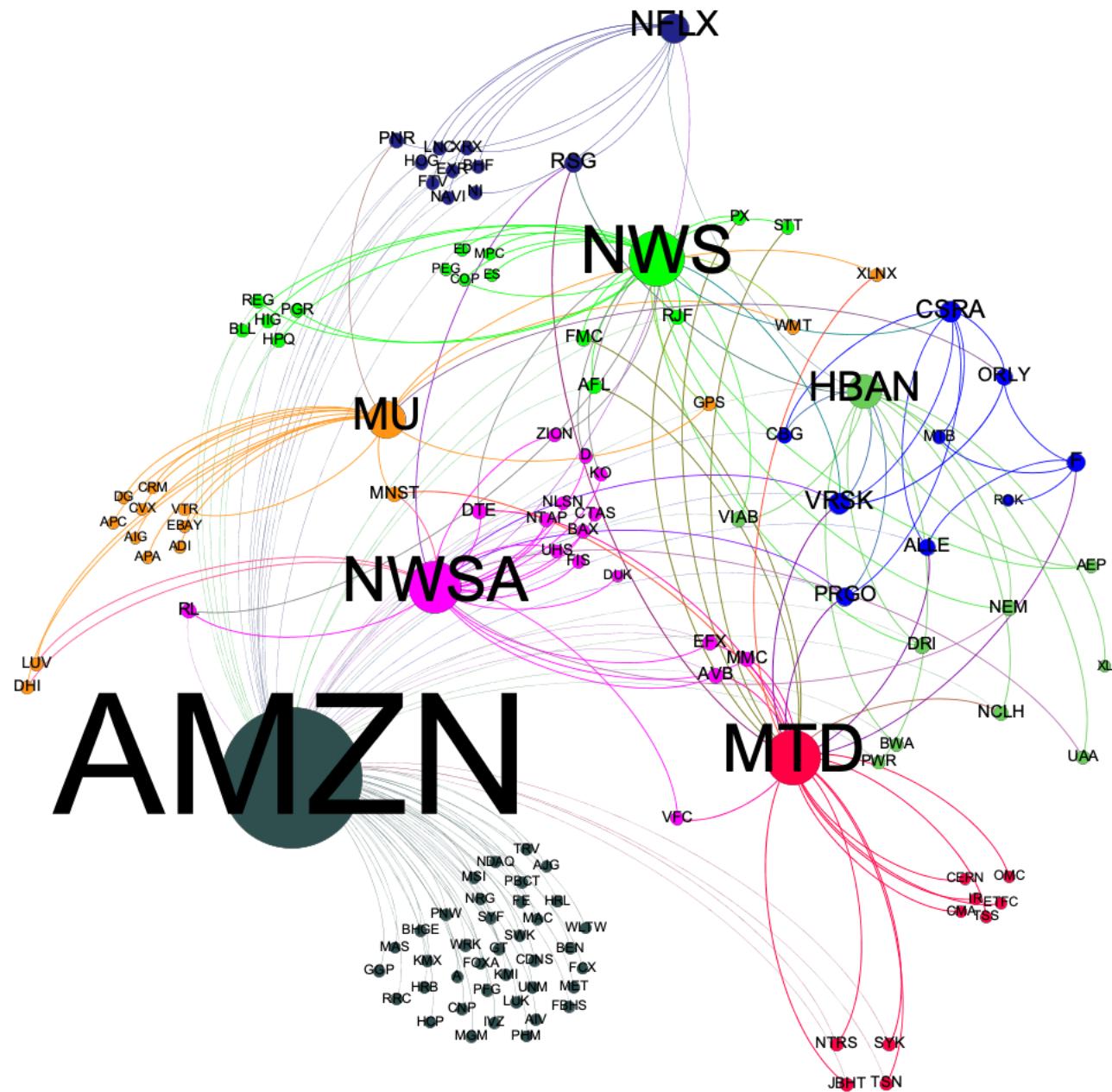


Figure 5.6: This figure shows the backbone network displaying communities for March 13th where at least nine firms are in the community. Each color represents a specific community, and the node's size represents the degree of the node.

Sample Rate	1 second	2 seconds	3 seconds	5 seconds	6 seconds	10 seconds	15 seconds	30 seconds	1 minute	2 minutes
<b>Mean</b>	0.084	0.101	0.108	0.112	0.112	0.106	0.094	0.069	0.046	0.033
<b>St. Deviation</b>	0.046	0.056	0.062	0.070	0.074	0.086	0.099	0.125	0.161	0.206
<b>1st percentile</b>	-0.014	-0.025	-0.033	-0.052	-0.061	-0.099	-0.141	-0.233	-0.344	-0.472
<b>Q1</b>	0.054	0.065	0.069	0.067	0.064	0.049	0.030	-0.013	-0.058	-0.098
<b>Median</b>	0.082	0.100	0.108	0.112	0.112	0.106	0.094	0.069	0.047	0.034
<b>Q3</b>	0.112	0.136	0.148	0.158	0.160	0.163	0.159	0.152	0.153	0.167
<b>99th percentile</b>	0.201	0.234	0.253	0.275	0.283	0.306	0.323	0.359	0.422	0.517

Table 5.1: This table shows the summary statistics for transfer entropy estimates.

Sample Rate	1 second	2 seconds	3 seconds	5 seconds	6 seconds	10 seconds	15 seconds	30 seconds	1 minute	2 minutes
9:30am-11:40am	0.097	0.117	0.125	0.130	0.130	0.123	0.109	0.082	0.058	0.044
11:40am-1:50pm	0.071	0.088	0.096	0.102	0.103	0.100	0.091	0.067	0.044	0.030
1:50pm-4pm	0.083	0.098	0.103	0.105	0.103	0.095	0.082	0.057	0.037	0.024

Table 5.2: Mean transfer entropy estimates for non-overlapping 130-minute daily trading windows at each sample rate.

Sample Rate	1 second	2 seconds	3 seconds	5 seconds	6 seconds	10 seconds	15 seconds	30 seconds
<b>09:30AM-10:00AM</b>	0.099	0.118	0.127	0.130	0.129	0.118	0.105	0.078
<b>10:00AM-10:30AM</b>	0.085	0.101	0.107	0.109	0.108	0.098	0.084	0.058
<b>10:30AM-11:00AM</b>	0.077	0.093	0.099	0.101	0.100	0.092	0.082	0.058
<b>11:00AM-11:30AM</b>	0.073	0.087	0.093	0.096	0.095	0.089	0.077	0.056
<b>11:30AM-12:00PM</b>	0.069	0.085	0.091	0.095	0.095	0.090	0.080	0.059
<b>12:00PM-12:30PM</b>	0.065	0.080	0.086	0.091	0.091	0.087	0.080	0.058
<b>12:30PM-01:00PM</b>	0.061	0.076	0.084	0.089	0.089	0.086	0.078	0.057
<b>01:00PM-01:30PM</b>	0.061	0.076	0.082	0.086	0.088	0.085	0.076	0.057
<b>01:30PM-02:00PM</b>	0.062	0.076	0.083	0.088	0.089	0.085	0.076	0.058
<b>02:00PM-02:30PM</b>	0.067	0.081	0.087	0.091	0.091	0.085	0.076	0.053
<b>02:30PM-03:00PM</b>	0.066	0.079	0.086	0.089	0.088	0.082	0.072	0.052
<b>03:00PM-03:30PM</b>	0.072	0.084	0.090	0.090	0.088	0.079	0.067	0.043
<b>03:30PM-04:00PM</b>	0.087	0.093	0.091	0.084	0.080	0.067	0.055	0.041

Table 5.3: Mean transfer entropy estimates for non-overlapping 30-minute daily trading windows at each sample rate.

	1 SECOND	2 SECONDS	3 SECONDS	5 SECONDS	6 SECONDS	10 SECONDS	15 SECONDS	30 SECONDS	1 MINUTE	2 MINUTES
MEAN	0.253	0.305	0.328	0.341	0.342	0.331	0.309	0.273	0.268	0.294
STD	0.037	0.031	0.024	0.015	0.014	0.020	0.027	0.033	0.030	0.026
MIN	0.205	0.260	0.291	0.317	0.312	0.271	0.227	0.193	0.200	0.235
1%	0.207	0.260	0.294	0.318	0.315	0.272	0.233	0.200	0.201	0.239
25%	0.224	0.281	0.309	0.330	0.332	0.323	0.296	0.257	0.250	0.278
50%	0.251	0.302	0.326	0.338	0.339	0.336	0.313	0.280	0.268	0.296
75%	0.266	0.318	0.340	0.351	0.350	0.343	0.331	0.300	0.290	0.315
99%	0.367	0.385	0.383	0.379	0.376	0.367	0.345	0.325	0.321	0.337
MAX	0.373	0.389	0.384	0.380	0.383	0.369	0.345	0.326	0.324	0.337

Table 5.4: This table displays the weighted network density for measures of transfer entropy at each sample rate.

Sample Rate	1 second	2 seconds	3 seconds	5 seconds	6 seconds	10 seconds	15 seconds	30 seconds	1 minute	2 minutes
Intercept	0.039*** (0.001)	0.060*** (0.002)	0.077*** (0.002)	0.091*** (0.002)	0.094*** (0.002)	0.091*** (0.002)	0.085*** (0.002)	0.071*** (0.002)	0.055*** (0.001)	0.043*** (0.001)
Wtd_InDegree <sub>t-1</sub>	0.602*** (0.012)	0.484*** (0.012)	0.385*** (0.013)	0.296*** (0.014)	0.277*** (0.013)	0.251*** (0.011)	0.2222*** (0.010)	0.124*** (0.008)	0.039*** (0.006)	0.020*** (0.006)

	Observations	29,642	29,642	29,642	29,642	29,642	29,642	29,642	29,642	29,642
Adjusted R2	0.361	0.234	0.148	0.087	0.077	0.063	0.049	0.015	0.001	0.000

Table 5.5: Auto-correlation in weighted network in-degree (Wtd\_InDegree). Standard errors clustered by firm appear in parentheses below each coefficient. \*, \*\*, and \*\*\* indicate significance at the 10%, 5%, and 1% levels, respectively.

Sample Rate	1 second	2 seconds	3 seconds	5 seconds	6 seconds	10 seconds	15 seconds	30 seconds	1 minute	2 minutes
Intercept	0.015*** (0.004)	0.019*** (0.005)	0.023*** (0.006)	0.029*** (0.008)	0.033*** (0.008)	0.039*** (0.009)	0.037*** (0.009)	0.032*** (0.006)	0.027*** (0.004)	0.025*** (0.003)
Wtd_OutDegree <sub>t-1</sub>	0.847*** (0.038)	0.838*** (0.043)	0.819*** (0.047)	0.773*** (0.057)	0.748*** (0.061)	0.680*** (0.074)	0.659*** (0.076)	0.605*** (0.073)	0.535*** (0.069)	0.419*** (0.061)
Observations	29,642	29,642	29,642	29,642	29,642	29,642	29,642	29,642	29,642	29,642
Adjusted R2	0.702	0.689	0.658	0.587	0.551	0.456	0.43	0.366	0.287	0.177

Table 5.6: Auto-correlation in weighted network out-degree (Wtd\_OutDegree). Standard errors clustered by firm appear in parentheses below each coefficient. \*, \*\*, and \*\*\* indicate significance at the 10%, 5%, and 1% levels, respectively.

	Sample Rate	1sec	2sec	3sec	5sec	6sec	10sec	15sec	30sec	1min	2min
Model	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	
EA <sub>ijt</sub> * Morning <sub>jt</sub>	-0.013*** (0.002)	-0.024*** (0.003)	-0.030*** (0.003)	-0.032*** (0.003)	-0.033*** (0.003)	-0.031*** (0.003)	-0.026*** (0.003)	-0.021*** (0.002)	-0.015*** (0.002)	-0.014*** (0.002)	
EA <sub>ijt</sub> * Morning <sub>jt</sub> * Abs_Surprise <sub>ijt</sub>	1.672*** (0.660)	2.627*** (0.802)	2.995*** (0.853)	3.117*** (0.899)	3.094*** (0.922)	2.920*** (0.923)	2.510*** (0.858)	2.525*** (0.744)	2.379*** (0.660)	2.379*** (0.742)	
EA <sub>ijt</sub> * Afternoon <sub>jt</sub>	0.006*** (0.001)	0.006*** (0.001)	0.006*** (0.001)	0.007*** (0.001)	0.006*** (0.001)	0.007*** (0.001)	0.005*** (0.002)	0.005*** (0.002)	0.002 (0.002)	-0.002 (0.002)	
EA <sub>ijt</sub> * Afternoon <sub>jt</sub> * Abs_Surprise <sub>ijt</sub>	0.672*** (0.321)	0.882*** (0.408)	0.946*** (0.480)	0.784*** (0.476)	0.755 (0.533)	0.715 (0.594)	0.795 (0.571)	0.649 (0.596)	1.019* (0.595)	1.058 (0.823)	
EA <sub>ijt</sub> * Evening <sub>jt</sub>	0.008*** (0.001)	0.009*** (0.001)	0.008*** (0.001)	0.007*** (0.001)	0.007*** (0.001)	0.005*** (0.001)	0.003*** (0.001)	-0.001 (0.001)	-0.001 (0.001)	-0.006*** (0.001)	
EA <sub>ijt</sub> * Evening <sub>jt</sub> * Abs_Surprise <sub>ijt</sub>	-0.045 (0.366)	0.0001 (0.375)	0.031 (0.389)	0.174 (0.417)	-0.022 (0.436)	0.102 (0.493)	0.008 (0.499)	0.829* (0.587)	0.382 (0.587)	0.716 (0.711)	
Share_Turnover <sub>it</sub>	-0.171*** (0.066)	-0.304*** (0.074)	-0.360*** (0.075)	-0.395*** (0.076)	-0.400*** (0.072)	-0.394*** (0.070)	-0.393*** (0.063)	-0.345*** (0.057)	-0.246*** (0.050)	-0.242*** (0.050)	
Abs_RET <sub>it</sub>	0.065*** (0.012)	0.094*** (0.014)	0.120*** (0.015)	0.150*** (0.016)	0.159*** (0.016)	0.176*** (0.016)	0.177*** (0.016)	0.179*** (0.015)	0.120*** (0.014)	0.200*** (0.013)	
Morning <sub>jt</sub>	0.026*** (0.000)	0.029*** (0.000)	0.030*** (0.000)	0.029*** (0.000)	0.027*** (0.000)	0.023*** (0.000)	0.019*** (0.000)	0.015*** (0.000)	0.014*** (0.000)	0.014*** (0.000)	
Evening <sub>jt</sub>	0.012*** (0.000)	0.010*** (0.000)	0.007*** (0.000)	0.003*** (0.000)	-0.002 (0.000)	-0.005*** (0.000)	-0.009*** (0.000)	-0.010*** (0.000)	-0.007*** (0.000)	-0.005*** (0.000)	

Fixed Effects	Firm, Week									
Observations	89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322
Adjusted R <sup>2</sup>	0.761	0.774	0.760	0.722	0.703	0.646	0.620	0.600	0.537	0.417

Table 5.7: Tests for earnings surprise as a determinant of daily outgoing information transfers. Standard errors clustered by firm appear in parentheses below each coefficient. \*, \*\*, and \*\*\* indicate significance at the 10%, 5%, and 1% levels, respectively.

	Sample Rate	1sec	2sec	3sec	5sec	6sec	10sec	15sec	30sec	1min	2min
Model	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	
$EA_{ijt} * Morning_{jt}$	0.014*** (0.002)	0.013*** (0.002)	0.012*** (0.003)	0.0113*** (0.003)	0.013*** (0.004)	0.008** (0.004)	0.008** (0.004)	0.014*** (0.005)	0.014*** (0.006)	0.019*** (0.006)	0.006 (0.008)
$EA_{ijt} * Morning_{jt} * Abs\_Surprise_{ijt}$	1.619** (0.803)	1.531* (0.852)	1.066 (1.084)	1.134 (0.942)	0.121 (1.202)	0.213 (1.487)	1.907 (1.627)	-1.06 (2.530)	2.655 (2.390)	-0.816 (3.893)	
$EA_{ijt} * Afternoon_{jt}$	0.006*** (0.002)	0.010*** (0.002)	0.014*** (0.002)	0.018*** (0.003)	0.020*** (0.003)	0.019*** (0.003)	0.017*** (0.004)	0.012*** (0.004)	0.001 (0.005)	-0.003 (0.006)	
$EA_{ijt} * Afternoon_{jt} * Abs\_Surprise_{ijt}$	0.821 (0.635)	0.789 (0.710)	0.365 (1.067)	-0.287 (1.121)	0.922 (1.033)	0.18 (1.289)	1.932 (1.598)	2.072 (2.087)	2.737 (2.145)	4.522 (3.070)	
$EA_{ijt} * Evening_{jt}$	0.008*** (0.613)	0.013*** (0.679)	0.014*** (0.932)	0.012*** (1.078)	0.015*** (1.141)	0.017*** (1.335)	0.018*** (1.608)	0.009* (1.960)	0.013*** (2.684)	0.011 (3.591)	
$EA_{ijt} * Evening_{jt} * Abs\_Surprise_{ijt}$	0.826 (0.079)	0.1 (0.074)	0.731 (0.065)	2.201** (0.060)	-0.392 (0.052)	-0.72 (0.063)	0.7 (0.058)	5.270*** (0.075)	0.974 (0.075)	0.91 (2.684)	
$Share\_Turnover_{it}$	0.538*** (0.079)	0.462*** (0.074)	0.354*** (0.065)	0.148** (0.060)	0.064 <sub>t</sub> (0.052)	-0.093 (0.063)	-0.004 (0.058)	-0.005 (0.075)	-0.006 (0.075)	-0.008 (0.094)	
$Abs\_RET_{it}$	0.026*** (0.011)	0.007 (0.013)	-0.023 (0.015)	-0.060*** (0.017)	-0.081*** (0.018)	-0.126*** (0.021)	-0.138*** (0.024)	-0.143*** (0.030)	-0.196*** (0.035)	-0.159*** (0.050)	
$Morning_{jt}$	0.025*** (0.001)	0.029*** (0.001)	0.029*** (0.001)	0.028*** (0.001)	0.027*** (0.001)	0.023*** (0.001)	0.018*** (0.001)	0.015*** (0.001)	0.013*** (0.001)	0.014*** (0.001)	
$Evening_{jt}$	0.012*** (0.000)	0.009*** (0.000)	0.007*** (0.001)	0.003*** (0.001)	-0.0002 (0.001)	-0.006*** (0.001)	-0.009*** (0.001)	-0.010*** (0.001)	-0.008*** (0.001)	-0.006*** (0.001)	

	Fixed Effects	Firm, Week									
Observations		89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322	89,322
Adjusted R2	0.521	0.42	0.338	0.259	0.244	0.216	0.182	0.107	0.042	0.016	

Table 5.8: Tests for earnings surprise as a determinant of daily incoming information transfers . Standard errors clustered by firm appear in parentheses below each coefficient. \* , \*\* , and \*\*\* indicate significance at the 10%, 5%, and 1% levels, respectively.

$\alpha$ levels	1	0.5	0.4	0.375	0.35	0.325	0.3	0.2
<b>Jan 31st</b>	(498, 247506)	(498, 41802)	(450, 3082)	(361, 1353)	(257, 592)	(153, 240)	(85, 110)	(3, 2)
<b>Feb 1st</b>	(498, 247506)	(498, 41244)	(464, 3181)	(414, 1574)	(324, 804)	(219, 418)	(145, 219)	(8, 6)
<b>Mar 13th</b>	(494, 243542)	(494, 47139)	(487, 5010)	(459, 2568)	(402, 1335)	(306, 714)	(215, 396)	(34, 40)

Table 5.9: This table shows the number of nodes and edges in network backbones at various  $\alpha$  levels for January 31st (Jan 31st), February 1st (Feb 1st), and March 13th (Mar 13th). Each table cell has the number notes followed by the number of edges. For example, Jan 31st has a node count of 498 and an edge count of 247,506 when  $\alpha = 1$ .

# Chapter 6

## Conclusion

### 6.1 Summary and Conclusion

In Chapter 2 we introduced an agnostic framework called SML that integrates a query-like language to simplify the development of machine learning pipelines. We provided a high level overview of it's architecture and grammar. SML was applied to well-known machine learning problems and demonstrated how the amount of code one has to write significantly decreases when SML is used.

In the future, we plan to extend the connector in SML's architecture to support more machine learning libraries and additional languages. We plan to extend SML's web application to include additional functionality to improve the ease of use. Feature selection, model selection, and parameter optimization are additional areas to add to SML. In addition to improving SML it can also be tested in a comparative analysis against similar works.

In Chapter 3 we used random forest trees in a supervised learning paradigm to predict the annual direction of profitability for firms with minimal information. We generated out-of-sample predictions of directional changes (increases or decreases) in five profitability measures: return on equity (ROE), return on assets (ROA), return on net operating assets (RNOA), cash flow from operations (CFO), and free cash flow (FCF) from 2011-2016. We found that the classification accuracy for each measure outperformed the random walk models.

We can further improve the classification accuracies by selecting a model that yields better performance but offers less interpretability. Given that we focus on a minimal set of

features in this research, we can incorporate additional features for additional performance gains. With more observations, we can extend the methodology outlined in this paper to a regression setting and make predictions about the magnitudes of profitability.

In Chapter 4 we implemented software to estimate bi-variate transfer entropy (TE). Our method was able to scale better to larger data than existing implementations. On large data, our implementation PyIF is up to 1072 times faster utilizing GPUs and up to 181 times faster utilizing CPUs than existing implementations that estimate bi-variate TE. For future work, we plan to improve the existing code base to improve the computational performance of PyIF further. In addition to this, we plan to implement additional estimators to estimate bi-variate TE.

In Chapter 5 we introduced a new approach to examine information transfers around earnings announcements. We studied the network effect of earnings announcements by constructing daily networks of pairwise cross-firm information transfers. Our approach to construct this network relies on non-parametric estimates in equity prices with measures of transfer entropy drawn from information theory. We provide evidence that earnings information produced by a single firm flows to other firms across industries. Results from our tests showed that cross-firm links are substantially stronger for firms on days with releases of earnings information and for firms with more unexpected earnings news.

Further, we found that communities of firms in the network also form between firms with links that are not adequately captured by characteristics that focus on existing literature. We plan to develop community detection algorithms to work in very dense weighted networks such as our information transfer networks for future work. We also plan to explore the utility of information transfers for predictions.

# References

- Vic Anand, Robert Brunner, Kelechi Ikegwu, and Theodore Sougiannis. Predicting profitability using machine learning. *SSRN*, Oct 2019. URL [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3466478](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3466478).
- Matthew Anderson and Jeff McMullin. Detecting information flows in markets, Nov 2018.
- S. Arya and D. Mount. Ann: library for approximate nearest neighbor searching. 1998.
- John W. Backus. The syntax and semantics of the proposed international algebraic language of the zurich acm-gamm conference. In *IFIP Congress*, pages 125–131. Butterworths, London, 1959.
- Sp Baginski. Intraindustry information transfers associated with management forecasts of earnings. *Journal of Accounting Research*, 25(2):196–216, 1987. URL <https://EconPapers.repec.org/RePEc:bla:joares:v:25:y:1987:i:2:p:196-216>.
- K. Bakshi and K. Bakshi. Considerations for artificial intelligence and machine learning: Approaches and use cases. In *2018 IEEE Aerospace Conference*, pages 1–9, 2018. doi: 10.1109/AERO.2018.8396488.
- Ray Ball and P Brown. Empirical evaluation of accounting income numbers. *Journal of Accounting Research*, 6(2):159–178, 1968.
- Monica Billio, Mila Getmansky, Andrew Lo, and Loriana Pelizzon. Econometric measures of connectedness and systemic risk in the finance and insurance sectors. *Journal of Financial Economics*, 104(3):535–559, 2012. URL <https://EconPapers.repec.org/RePEc:eee:jfinec:v:104:y:2012:i:3:p:535-559>.
- Terry Bossomaier, Lionel Barnett, Michael Harré, and Joseph Lizier. *An Introduction to Transfer Entropy*. 01 2016. doi: 10.1007/978-3-319-43222-9.
- Mark T. Bradshaw, Michael S. Drake, James N. Myers, and Linda A. Myers. A re-examination of analysts’ superiority over time-series forecasts of annual earnings. *Review of Accounting Studies*, 2012. URL <https://doi.org/10.1007/s11142-012-9185-8>.
- L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, Monterey, CA, 1984.

Francois Brochet, Kalin Kolev, and Alina Lerman. Information transfer and conference calls. *Review of Accounting Studies*, 23(3):907–957, 2018. URL [https://EconPapers.repec.org/RePEc:spr:reaccs:v:23:y:2018:i:3:d:10.1007\\_s11142-018-9444-4](https://EconPapers.repec.org/RePEc:spr:reaccs:v:23:y:2018:i:3:d:10.1007_s11142-018-9444-4).

Robert Brunner, Kelechi Ikegwu, Matthew Anderson, and Jeff McMullin. Cross-firm information transfers during earnings season: A network approach, April 2021.

John Y. Campbell, Andrew W. Lo, and A.Craig MacKinlay. *The Econometrics of Financial Markets*. Princeton University Press, 1997. ISBN 9780691043012. URL <http://www.jstor.org/stable/j.ctt7skm5>.

Vitor Cerqueira, Luis Torgo, and Igor Mozetic. Evaluating time series forecasting models: An empirical study on performance estimation methods, 2019.

Prithwiraj Choudhury, Ryan T. Allen, and Michael G. Endres. Machine learning for pattern discovery in management research. *Strategic Management Journal*, 42(1):30–57, 2021. doi: <https://doi.org/10.1002/smj.3215>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/smj.3215>.

Aaron Clauset, M. E. J. Newman, and Cristopher Moore. Finding community structure in very large networks. *Phys. Rev. E*, 70:066111, Dec 2004. doi: 10.1103/PhysRevE.70.066111. URL <https://link.aps.org/doi/10.1103/PhysRevE.70.066111>.

Greg J. Clinch and Norman A. Sinclair. Intra-industry information releases: A recursive systems approach. *Journal of Accounting and Economics*, 9(1):89–106, 1987. ISSN 0165-4101. doi: [https://doi.org/10.1016/0165-4101\(87\)90018-8](https://doi.org/10.1016/0165-4101(87)90018-8). URL <https://www.sciencedirect.com/science/article/pii/0165410187900188>.

Andrew Collette. *Python and HDF5*. O'Reilly, 2013.

Francis X. Diebold and Kamil Yilmaz. On the Network Topology of Variance Decompositions: Measuring the Connectedness of Financial Firms. Koç University-TUSIAD Economic Research Forum Working Papers 1124, Koc University-TUSIAD Economic Research Forum, October 2011. URL <https://ideas.repec.org/p/koc/wpaper/1124.html>.

Pedro Domingos. A few useful things to know about machine learning. volume 55, pages 78–87, New York, NY, USA, October 2012. ACM. doi: 10.1145/2347736.2347755. URL <http://doi.acm.org/10.1145/2347736.2347755>.

Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Tobias Springenberg, Manuel Blum, and Frank Hutter. Auto-sklearn: Efficient and robust automated machine learning. pages 123–143.

George Foster. Intra-industry information transfers associated with earnings releases. *Journal of Accounting and Economics*, 3(3):201–232, 1981. ISSN 0165-4101. doi: [https://doi.org/10.1016/0165-4101\(81\)90003-3](https://doi.org/10.1016/0165-4101(81)90003-3). URL <https://www.sciencedirect.com/science/article/pii/0165410181900033>.

E. Frank, M. A. Hall, G. Holmes, R. Kirkby, B. Pfahringer, and I. H. Witten. *Weka: A machine learning workbench for data mining.*, pages 1305–1314. Springer, Berlin, 2005. URL <http://researchcommons.waikato.ac.nz/handle/10289/1497>.

Jerry C. Y. Han and J. Wild. Unexpected earnings and intraindustry information transfers - further evidence. *Journal of Accounting Research*, 28:211–219, 1990.

Weiqiang Hang and Timothy Banks. Machine learning applied to pack classification. *International Journal of Market Research*, 61(6):601–620, 2019. doi: 10.1177/1470785319841217.

Rebecca N. Hann, Heedong Kim, and Yue Zheng. Intra-industry information transfers: evidence from changes in implied volatility around earnings announcements. *Review of Accounting Studies*, 24(3):927–971, September 2019. doi: 10.1007/s11142-019-9487-1. URL [https://ideas.repec.org/a/spr/reaccs/v24y2019i3d10\\_1007\\_s11142-019-9487-1.html](https://ideas.repec.org/a/spr/reaccs/v24y2019i3d10_1007_s11142-019-9487-1.html).

Charles R. Harris, K. Jarrod Millman, Stéfan J van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585:357–362, 2020. doi: 10.1038/s41586-020-2649-2.

Kewei Hou, Mathijs A. van Dijk, and Yinglei Zhang. The implied cost of capital: A new approach. *Journal of Accounting and Economics*, 53(3):504 – 526, 2012. ISSN 0165-4101. doi: <https://doi.org/10.1016/j.jacceco.2011.12.001>. URL <http://www.sciencedirect.com/science/article/pii/S0165410111000966>.

K. M. Ikegwu, J. Trauger, J. McMullin, and R. J. Brunner. Pyif: A fast and light weight implementation to estimate bivariate transfer entropy for big data. In *2020 SoutheastCon*, pages 1–6, 2020. doi: 10.1109/SoutheastCon44009.2020.9249650.

Kelechi Ikegwu, Micheal Hao, Nerraj Asthana, and Robert Brunner. Standard machine learning language: A language agnostic framework for streamlining the development of machine learning pipelines. 2017. URL [https://github.com/lcdm-uiuc/Publications/blob/master/2017\\_Kelechi\\_Mike\\_Brunner/main.pdf](https://github.com/lcdm-uiuc/Publications/blob/master/2017_Kelechi_Mike_Brunner/main.pdf).

Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning: with Applications in R*. Springer, 2013. URL <https://faculty.marshall.usc.edu/gareth-james/ISL/>.

A. Kaiser and T. Schreiber. Information transfer in continuous processes. *Physica D, v.166, 43-62 (2002)*, 166, 06 2002. doi: 10.1016/S0167-2789(02)00432-3.

Shiraj Khan, Sharba Bandyopadhyay, Auroop Ganguly, Sunil Saigal, David Erickson, Vladimir Protopopescu, and George Ostroumov. Relative performance of mutual information estimation methods for quantifying the dependence among short and noisy data.

*Physical review. E, Statistical, nonlinear, and soft matter physics*, 76:026209, 09 2007. doi: 10.1103/PhysRevE.76.026209.

Brent Komer, James Bergstra, and Chris Eliasmith. Hyperopt-sklearn. pages 105–121.

M. Koning and C. Smith. *Decision Trees and Random Forests: A Visual Introduction for Beginners*. Amazon Digital Services LLC - Kdp Print Us, 2017. ISBN 9781549893759. URL [https://books.google.com/books?id=Hi\\_CtAEACAAJ](https://books.google.com/books?id=Hi_CtAEACAAJ).

Lars Kotthoff, Chris Thornton, Holger H. Hoos, Frank Hutter, and Kevin Leyton-Brown. Auto-weka: Automatic model selection and hyperparameter optimization in weka. pages 89–103.

Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Phys. Rev. E*, 69:066138, Jun 2004. doi: 10.1103/PhysRevE.69.066138.

Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. Numba: A llvm-based python jit compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*, pages 1–6, 2015.

M. Lichman. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>.

Michael Lindner, Raul Vicente, Viola Priesemann, and Michael Wibral. Tren-tool: A matlab open source toolbox to analyse information flow in time series data with transfer entropy. *BMC Neuroscience*, 12(1), 01 2011. doi: 10.1186/1471-2202-12-119. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3287134/>. Citations 56.

R. Marschinski and H. Kantz. Analysing the information flow between financial time series. *The European Physical Journal B - Condensed Matter and Complex Systems*, 30:275–281, 2002.

Steve Menard and Luis Nell. "jpype documentation". 2018. URL <https://jpype.readthedocs.io/en/latest/>.

S. Monahan. Financial statement analysis and earnings forecasting. 2018.

ANN Library: David Mount, Sunil Arya. Transfer Entropy Packge: Ghazaleh Haratinezhad Torbati, and Glenn Lawyer. *TransferEntropy: The Transfer Entropy Package*, 2015. URL <https://CRAN.R-project.org/package=TransferEntropy>. R package version 1.5.

M. E. J. Newman. *Networks: an introduction*. Oxford University Press, Oxford; New York, 2010. ISBN 9780199206650 0199206651. URL [http://www.amazon.com/Networks-An-Introduction-Mark-Newman/dp/0199206651/ref=sr\\_1\\_5?ie=UTF8&qid=1352896678&sr=8-5&keywords=complex+networks](http://www.amazon.com/Networks-An-Introduction-Mark-Newman/dp/0199206651/ref=sr_1_5?ie=UTF8&qid=1352896678&sr=8-5&keywords=complex+networks).

NVIDIA, Péter Vingermann, and Frank H.P. Fitzek. Cuda, release: 10.2.89, 2020. URL <https://developer.nvidia.com/cuda-toolkit>.

Chris Olsen and J. Richard Dietrich. Vertical information transfers: The association between retailers' sales announcements and suppliers' security returns. *Journal of Accounting Research*, 23:144–166, 1985. ISSN 00218456, 1475679X. URL <http://www.jstor.org/stable/2490695>.

Randal S. Olson, Nathan Bartley, Ryan J. Urbanowicz, and Jason H. Moore. Evaluation of a tree-based pipeline optimization tool for automating data science. *CoRR*, abs/1603.06212, 2016. URL <http://arxiv.org/abs/1603.06212>.

Jane A. Ou and Stephen H. Penman. Financial statement analysis and the prediction of stock returns. *Journal of Accounting and Economics*, 11(4):295 – 329, 1989. ISSN 0165-4101. doi: [https://doi.org/10.1016/0165-4101\(89\)90017-7](https://doi.org/10.1016/0165-4101(89)90017-7). URL <http://www.sciencedirect.com/science/article/pii/0165410189900177>.

Milan Paluš, Vladimír Komárek, Zbyněk Hrnčíř, and Katalin Štěrbová. Synchronization as adjustment of information rates: Detection from bivariate time series. *Phys. Rev. E*, 63:046211, Mar 2001. doi: 10.1103/PhysRevE.63.046211. URL <https://link.aps.org/doi/10.1103/PhysRevE.63.046211>.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Jason Pellerin. “nose. 2021. URL <https://nose.readthedocs.io/en/latest/>.

Grace Pownall and Gregory Waymire. Voluntary disclosure choice and earnings information transfer. *Journal of Accounting Research*, 27:85–105, 1989. ISSN 00218456, 1475679X. URL <http://www.jstor.org/stable/2491066>.

Kenneth R. Ahren and Jarrad Harford. The importance of industry links in merger waves. *The Journal of Finance*, 69(2):527–576, 2014. ISSN 00221082, 15406261. URL <http://www.jstor.org/stable/43611161>.

Armin Rigo and Maciej Fijalkowski. "cffi documentation". 2018. URL <https://cffi.readthedocs.io/en/latest/>.

N. Rizzolo and D. Roth. Learning based java for rapid development of nlp systems. In *LREC*, Valletta, Malta, 5 2010. URL <http://cogcomp.cs.illinois.edu/papers/RizzoloRo10.pdf>.

D. Roth. Learning based programming. *Innovations in Machine Learning: Theory and Applications*, 2005. URL <http://cogcomp.cs.illinois.edu/papers/Roth05.pdf>.

L. Sandoval. Structure of a global network of financial companies based on transfer entropy. *Entropy*, 16:4443–4482, 2014.

Samer Muthana Sarsam, Hosam Al-Samarraie, Ahmed Ibrahim Alzahrani, and Bianca Wright. Sarcasm detection using machine learning algorithms in twitter: A systematic review. *International Journal of Market Research*, 62(5):578–598, 2020. doi: 10.1177/1470785320921779.

T. Schreiber. Measuring information transfer. *Physical review letters*, 85(2):461–464, 2000.

M. Ángeles Serrano, Marián Boguñá, and Alessandro Vespignani. Extracting the multiscale backbone of complex weighted networks. *Proceedings of the National Academy of Sciences*, 106(16):6483–6488, 2009. ISSN 0027-8424. doi: 10.1073/pnas.0808904106. URL <https://www.pnas.org/content/106/16/6483>.

Behrendt Simon, Dimpfl Thomas, Peter Franziska J., and Zimmermann David J. Rtransferentropy — quantifying information flow between different time series using effective transfer entropy. *SoftwareX*, 10(100265):1–9, 2019. URL <https://doi.org/10.1016/j.softx.2019.100265>.

J. David Spiceland, James F. Sepe, and Mark Nelson. ntermediate accounting. chapter 1. McGraw-Hill Education, 2018.

James Stone. *Information Theory: A Tutorial Introduction*. 02 2015. ISBN 978-0956372857. doi: 10.13140/2.1.1633.8240.

Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.

CLARE WANG. Accounting standards harmonization and financial statement comparability: Evidence from transnational information transfer. *Journal of Accounting Research*, 52(4):955–992, 2014. doi: <https://doi.org/10.1111/1475-679X.12055>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1475-679X.12055>.

Paul L. Williams and Randall D. Beer. Generalized measures of information transfer. *ArXiv*, abs/1102.1507, 2011.

Patricia Wollstadt, Joseph T. Lizier, Raul Vicente, Conor Finn, Mario Martinez-Zarzuela, Pedro Mediano, Leonardo Novelli, and Michael Wibral. Idtxl: The information dynamics toolkit xl: a python package for the efficient analysis of multivariate information dynamics in networks. *Journal of Open Source Software*, 4(34):1081, 2019. doi: 10.21105/joss.01081. URL <https://doi.org/10.21105/joss.01081>.

Ying Yang. Research on the optimization of the supplier intelligent management system for cross-border e-commerce platforms based on machine learning. *Information Systems & e-Business Management*, 18(4):851 – 870, 2020. ISSN 16179846.

# Appendix A

## Standard Machine Learning Language Supplemental Code

### A.1 Iris Python Code

This shows the code required to replicate the same actions of the SML *Query* in Figure ??.

It's important to note that detailed documentation is publicly available. The purpose of this listing is to highlight the level of complexity relative to a SML query.

```
1 import pandas as pd
2 import numpy as np
3
4 from sklearn.preprocessing import label_binarize
5 import sklearn.cross_validation as cv
6 from sklearn.multiclass import OneVsRestClassifier
7 from sklearn.svm import SVC
8 from sklearn.metrics import roc_curve, auc
9
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 names = ['sepal length(cm)', 'sepal width(cm)', 'petal length(cm)', 'petal
13 width(cm)', 'species']
14 data = pd.read_csv('../data/iris.csv', names=names)
15 iris_classes = ['Iris-setosa', 'Iris-versicolor', 'Iris-virginica']
16 features = np.c_[data.drop('species', 1).values]
17 labels = label_binarize(data['species'], classes=iris_classes)
18 n_classes = labels.shape[1]
```

```

19
20 x_train , x_test , y_train , y_test = cv.train_test_split(features , labels ,
   test_size=0.25)
21 svm = OneVsRestClassifier(SVC(kernel='linear' , probability=True))
22 l = svm.fit(x_train , y_train)
23 predict_score = model.decision_function(x_test)
24 test_set_results = model.score(x_test , y_test) * 100
25 print ('SVM Prediction Accuracy = {0:6.2f}%'.format(test_set_results) )
26 fpr = dict()
27 tpr = dict()
28 roc_auc = dict()
29
30 for i in range(n_classes):
31 fpr[i] , tpr[i] , _ = roc_curve(y_test[:, i] , predict_score[:, i])
32 roc_auc[i] = auc(fpr[i] , tpr[i])
33 plt.rcParams['figure.figsize']=(12,12)
34 # Class Info
35 columns = [0,1,2,3]
36 cmap_class = ['Purples_r' , 'Greens_r' , 'Oranges_r' , 'Greys_r' ]
37 color_class1D = ['purple' , 'darkgreen' , 'orange' , 'grey' ]
38 column_headers = data.columns.values.tolist() # Grab headers from df
39 column_headers = [column_headers[x] for x in columns] # Map headers to indices
   selected
40
41 label = 'species'
42 fig , ax = plt.subplots(len(columns) , len(columns))
43 for ic , cc , cc1D in zip(iris_classes , cmap_class , color_class1D):
44   iris_class_data = data.loc[data.species == ic] # sep class
45
46 #Generate kde plot matrix for class
47 for col1 , i in enumerate(columns):
48   for col2 , j in enumerate(columns):

```

```

49     if i == j:
50         sns.kdeplot(iris_class_data[iris_class_data.columns[col1]], ax=
51                     ax[col1][col2], color=cc1D, shade=True, legend=False)
52     else:
53         sns.kdeplot(iris_class_data[iris_class_data.columns[col1]],
54                     iris_class_data[iris_class_data.columns[col2]], ax=ax[col1][
55                     col2], cmap=cc)
56
57 # Formatting
58
59     if j == 0:
60         ax[i,j].set_xticklabels([])
61         ax[i,j].set_ylabel(column_headers[i])
62         ax[i,j].set_xlabel(' ')
63         if i == len(columns)-1:
64             ax[i,j].set_xlabel(column_headers[j])
65         elif i == len(columns)-1:
66             ax[i,j].tick_params(axis='y', which='major', bottom='off')
67             ax[i,j].set_yticklabels([])
68             ax[i,j].set_xlabel(column_headers[j])
69             ax[i,j].set_ylabel(' ')
70         else:
71             ax[i,j].set_xticklabels([])
72             ax[i,j].set_xlabel(' ')
73             ax[i,j].set_yticklabels([])
74             ax[i,j].set_ylabel(' ')
75
76 plt.show()
77 plt.close()

```

## A.2 Auto-MPG Python Code

This shows the code required to replicate the same actions of the SML *Query* in Figure ??.

It's important to note that detailed documentation is publicly available in <sup>14</sup>, the purpose of this figure is to highlight the level of complexity relative to a SML query.

```
1 import pandas as pd
2 import numpy as np
3
4 import matplotlib.pyplot as plt
5 from sklearn import linear_model
6 from sklearn.cross_validation import train_test_split
7 from sklearn.learning_curve import learning_curve, validation_curve
8
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 plt.rcParams['figure.figsize']=(12,12)
13 sns.set()
14
15 names = ['mpg', 'cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model_year', 'origin', 'car_name']
16
17 #load dataset
18 data = pd.read_csv('../data/auto-mpg.csv', sep = '\s+', header = None, names = names)
19 data_clean=data.applymap(lambda x: np.nan if x == '?' else x).dropna()
20 X = data_clean[['cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model_year', "origin"]]
21 #Select target column
22 y = data_clean['mpg']
23 #Split data into training and testing sets
24 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8,
```

```

    test_size=0.2)

25

26 # Define and train linear regression model
27 estimator = linear_model.LinearRegression()# Generate Learning Curves
28 train_sizes, train_scores, test_scores = learning_curve(estimator, X_train,
29   y_train)
30 # Train Linear Regression Model
31 estimator.fit(X_train, y_train)# Generate Validation Curves
32 param_range = np.arange(0, 5)
33
34 v_train_scores, v_test_scores = validation_curve(estimator, X_test, y_test,
35   param_name='normalize', param_range=param_range)
36
37 score = estimator.score(X_test, y_test)
38 print('Accuracy :', score)
39
40 g = sns.PairGrid(data_clean, palette='PuOr_r')
41 g = g.map_diag(sns.kdeplot, shade=True) # can't add color arg...
42
43 g = g.map_upper(sns.kdeplot, cmap='PuOr_r')
44 g = g.map_lower(sns.kdeplot, cmap='PuOr_r')
45
46 plt.show()
47 plt.close()

48 color_pal = ['purple', 'dark green', 'orange', 'grey'] # For 1-D KDE
49 cmap_pal = ['PuOr_r'] # For 2-D KDE
50 classes = [] # May not have a class for categories
51
52 column_headers = data_clean.columns.values.tolist() # Grab headers from df
53 column_headers = [column_headers[x] for x in columns] # Map headers to indices
      selected

```

```

53 fig , ax = plt . subplots( len (columns) , len (columns))
54 if not classes :
55     for col1 , i in enumerate (columns) :
56         for col2 , j in enumerate (columns) :
57             if i == j :
58                 sns . kdeplot ( data _ clean [ data _ clean . columns [ col1 ]] , ax=ax [ col1 ][
59                     col2 ] , color=color _ pal [0] , shade=True , legend=False )
60             else :
61                 sns . kdeplot ( data _ clean [ data _ clean . columns [ col1 ]] , data _ clean [
62                     data _ clean . columns [ col2 ]] , ax=ax [ col1 ][ col2 ] , cmap=cmap _ pal
63                     [0])
64
64 # Formatting
65 if j == 0 :
66     ax [ i , j ]. set _ xticklabels ([ ])
67     ax [ i , j ]. set _ ylabel ( column _ headers [ i ])
68     ax [ i , j ]. set _ xlabel ( '' )
69     if i == len (columns)-1 :
70         ax [ i , j ]. set _ xlabel ( column _ headers [ j ])
71     elif i == len (columns)-1 :
72         ax [ i , j ]. tick _ params ( axis='y' , which='major' , bottom='off' )
73         ax [ i , j ]. set _ yticklabels ([ ])
74         ax [ i , j ]. set _ xlabel ( column _ headers [ j ])
75         ax [ i , j ]. set _ ylabel ( '' )
76     else :
77         ax [ i , j ]. set _ xticklabels ([ ])
78         ax [ i , j ]. set _ xlabel ( '' )
79         ax [ i , j ]. set _ yticklabels ([ ])
80         ax [ i , j ]. set _ ylabel ( '' )
81 plt . show ()
82 plt . close ()

```

```

82 plt.figure()
83 plt.xlabel("Validation examples")
84 plt.ylabel("Score")
85
86 v_train_scores_mean = np.mean(v_train_scores, axis=1)
87 v_train_scores_std = np.std(v_train_scores, axis=1)
88 v_test_scores_mean = np.mean(v_test_scores, axis=1)
89 v_test_scores_std = np.std(v_test_scores, axis=1)
90
91 plt.fill_between(param_range, v_train_scores_mean - v_train_scores_std,
92                   v_train_scores_mean + v_train_scores_std, alpha=0.1, color="orange")
92 plt.fill_between(param_range, v_test_scores_mean - v_test_scores_std,
93                   v_test_scores_mean + v_test_scores_std, alpha=0.1, color="purple") plt.plot
94 (param_range, v_train_scores_mean, 'o--', color="orange", label="Training
95 score")
96 plt.plot(param_range, v_test_scores_mean, 'o--', color="purple", label="Cross-
97 validation score")
98 plt.legend(loc="best")
99 plt.show()
100 plt.close()

```