



# Analyzing The Trends In Stock Volatility Through The Application Of Abduction Methods

---

KOUSTUBH RAO

# Who Am I?

---

Yermal Koustubh Rao (200100176) is a final year undergraduate pursuing his bachelors in the Computer Science and Engineering department at Indian Institute of Technology, Bombay (IITB), graduating with a Honors degree in 2024.

This is the final project presentation for the course CS 781: Formal Methods in Machine Learning, offered by Prof. Supratik Chakraborty at the Indian Institute of Technology, Bombay (IITB) in the Autumn Semester (2023 - 2024).

# Introduction

---

Machine Learning in the recent times has proven itself to be indispensable. Due to its easy scaling as well as the fast improvements in GPU technology, larger models are trained with larger datasets leading to smart and fast models.

To this end, machine learning is also used extensively in the Fintech field to predict instrument performance like stocks, options etc. LSTMS, CNNs and RNNs have proven to perform really well in predicting prices of instruments in the recent past. But it is a hard task to give strict confirmation certificates for desired properties of ML models.

For example, a model performing really well up until now may end up crashing at the next moment leading to losses in the order of millions. Hence, machine learning models aren't used to the expected levels in the finance sector .

We observed this phenomenon throughout the course where proving strict properties of learning models is indeed a difficult task. The approaches that we learned were very loose in terms of the approximations made or would require heavy calculations to be made for large models. This is inefficient as the learning models are continuously learning with time. To this end, in this report we see how we can get strict confirmations on some of the properties of a machine learning model which predicts stock volatility.

# Problem Statement

---

To find and understand the minimal set of days from the recent past which sufficiently explains the volatility trend of a stock for the day. Volatility trend here refers to the close open difference in the stock price for a given day. The day ends on a positive trend when closing price is greater than the opening price and negative trend otherwise.

# Motivation

---

There are several factors which affects the performance of a stock for a given day. Events like elections, budgets, government policies, calamities etc. can greatly affect the performance during market hours.

It is a difficult task to judge which events in the recent past actually affect's a company in a visible way, Some events may lead to a bearish performance and some events may lead to bullish performance. Some might have little to no effect in regards to a company's performance.

It is exactly this question upon which I wanted to spend time and use the tools and theory taught during the course, to get into deeper.

# Motivation

---

This minimal set, gives us a deeper understanding of how certain events are related to a firm and how this gets portrayed in the market. To limit the scope of this vast field, we commit this research only towards companies from the finance industries. Although one might expect to find several events which affect the finance sector as a whole in common, certain events/days might be particular for that bank/firm. For example commercial banks might be affected by changes in tariff rules but private banks might not. Hedge funds might be immune to many of these events, at least in theory due to the hedging of their risk, which might be detrimental to the common bank firms otherwise.

This understanding can be leveraged in the future for better understanding and hence prediction of the relation between events and their play in market. In effect preparing oneself better and hence saving firms from losses, giving a more stable economy.

# Tools Used

---

This analysis is made along the lines of the following papers, Abduction-Based Explanations for Machine Learning Models by Alexey Ignatiev, Nina Narodytska, Joao Marques-Silva and Using MaxSAT for Efficient Explanations of Tree Ensembles by Alexey Ignatiev, Yacine Izza, Peter J. Stuckey, Joao Marques-Silva. An exact tool for the earlier paper, which was the main basis upon which the following project was undertaken was not available on the internet. Hence, the tool 'XReason' for the latter paper was used. This tool was modified extensively to come up with a working tool for the earlier paper. The authors of the earlier paper was contacted, and upon request a broken prototype of the earlier paper was obtained. A lot of wiring and restructuring was done to come up with a working tool which was the main ingredient of the paper. The source code for 'Xreason' can be found at <https://github.com/alexeyignatiev/xreason>, the code base for 'XPlainer', which is the actual tool developed for the first paper can be requested from the author.

# Modifications

---

The tools obtained, both 'XPlainer' and 'XReason' were hardcoded for the MNIST dataset, hence augmentations had to be made to accommodate the stock data needed for analysis in this project. To this end, certain configuration files were written which consisted configuration details like, loss functions, target feature, weights, epochs, batch size etc. The models.py was updated to handle deeper layers, the older files had only 2 hidden layers which was not sufficient for the dataset in our case. The Neural Network to SMT encoding was present in the 'XPlainer' repo, but the minimum hitting set algorithm was present in 'XReason', this was stitched.

The loss functions were updated to reshape the model predictions to match the shape of the target predictions. The pipeline was hardcoded for MNIST and other datasets which the original authors used in their research, so these had to be updated to intake 'stock\_data' dataset, relevant for this project. These updates were tiny but many in number.



# Data Preparation

---

The top 99 firms from the financial/banking sector was chosen for the analysis listed on NYSE/NASDAQ in the United States of America (USA). Yahoo Finance's yfinance api, was used to obtain the entire stock history of these 99 firms which were then stored in csv format for each firm.

This data contained various information for a trading day like opening price, closing price, high, low, traded volume etc. To limit the scope of the analysis, we were only interested in the closing price opening price difference. Hence a finer dataset was made which had only one point for a trading day, 1 if closing price was greater than the opening price else 0. This indicator would roughly tell about the stock volatility trend for that day.

# Model Architecture

---

It was found that a Feed Forward Neural Network (FFNN) was sufficient for the analysis, and there was no need for more complex architectures like CNN. Intuitively this should make sense as local contexts and patterns observed don't affect the final prediction which lies at the bottom right corner of the image. Of course this could have identified structures in the data but nevertheless unimportant for the task at hand.

Deep FFNN: with five fully connected layers of sizes  $n$ ,  $2n$ ,  $2n$ ,  $n$  and  $1$  where  $n$  is the size of the input. ReLU activation was used for each layer with a sigmoid activation at the final layer which predicts the class, positive or negative trend.

# Training

---

A training module was implemented to train the above architectures. Configuration files dictated the location for the dataset, location to store the model, number of epochs, model type, learning rate, loss function, optimizer etc. Since we needed the analysis to be particular for some day, the dataset contained 99 points for each of the firm. Experiments were done on the various choices for the length of the stock history used to predict the trend for a particular day. Due to the small size of the data, training was pretty fast. Further the hyper-parameters and the training duration was chosen so as to achieve a 100\% accuracy on the training data so that the prediction is accurate for all data points. The trained model was stored as a pickle file.

# Encoding and Hitting

---

The MILP encoding of the trained neural network was obtained using the configuration file and the stored pickle file of the model parameters. The encoding obtained was in the form of a text file. And this text file was taken as input to the module which performed the main minimum/minimal hitting set algorithms to find the hitting sets for the input data points. Z3/PySMT Solvers were used.

# Hitting Set Algorithm

---

---

**Algorithm 2:** Computing a smallest size explanation

---

**Input:**  $\mathcal{F}$  under  $\mathcal{M}$ , initial cube  $C$ , prediction  $\mathcal{E}$

**Output:** Cardinality-minimal explanation  $C_{\mathcal{M}}$

```
1 begin
2    $\Gamma \leftarrow \emptyset$ 
3   while true do
4      $h \leftarrow \text{MinimumHS}(\Gamma)$ 
5     if  $\text{Entails}(h, \mathcal{F} \rightarrow \mathcal{E}, \mathcal{M})$  then
6       return  $h$ 
7     else
8        $\mu \leftarrow \text{GetAssignment}()$ 
9        $C' \leftarrow \text{PickFalseLits}(C \setminus h, \mu)$ 
10       $\Gamma \leftarrow \Gamma \cup C'$ 
11 end
```

---

# Results

---

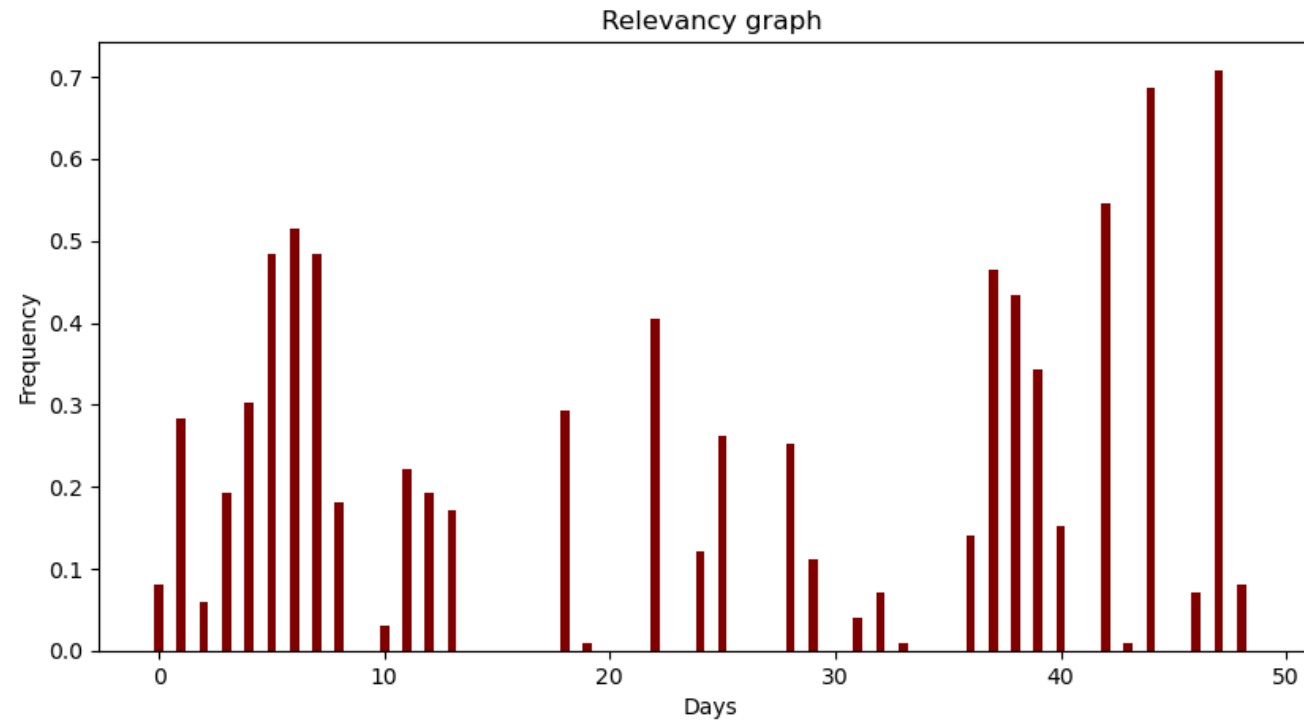
The above pipeline was run across 99 of the top financial firms/banks listed in NYSE/NASDAQ to predict the trend for 8<sup>th</sup> November 2023. A history of 49 days was used for the prediction and subsequently the minimum hitting set was calculated across all the firms.

The days which were the most important were 6<sup>th</sup> November (which was part of the minimum hitting set for 70/99 firms, 1<sup>st</sup> November (68/99) firms and 30<sup>th</sup> October (54/99) firms.

One can notice that the days 6<sup>th</sup> November and 1<sup>st</sup> November were a major deciding factor for the stock trend observed on 8<sup>th</sup> November for the financial firms.

# Importance and Relevance

---



# Result

---

One could have assumed that only the very recent past (say the past week) were sufficient to predict today's trend but the graph says otherwise. Points as far as a month ago were relevant for predicting the stock trend.

Surprisingly 7<sup>th</sup> November was not so relevant as expected. In fact, 7<sup>th</sup> November was part of the minimum hitting set of only 8 out of the 99 firms studied.

One can hypothesise that an important event might have occurred on 6<sup>th</sup> November relevant to the finance sector. Or rather an important event occurred on 3<sup>rd</sup> November (as 6<sup>th</sup> was a Monday) whose decision was portrayed by the stock market on the 6<sup>th</sup>.



# Finance Sector Index

---



# Finance Sector Index

---



# Results

---

Indeed, a spike can be observed on 3<sup>rd</sup> November and an observable change on the graph slope can be seen on 1<sup>st</sup> November.

A talented eye can deduce those events particular to these days responsible for having such an impact on the finance sector.

# Minimum Sets

---

The minimum set for Arch Capital is [1, 7, 22, 42, 44, 47]

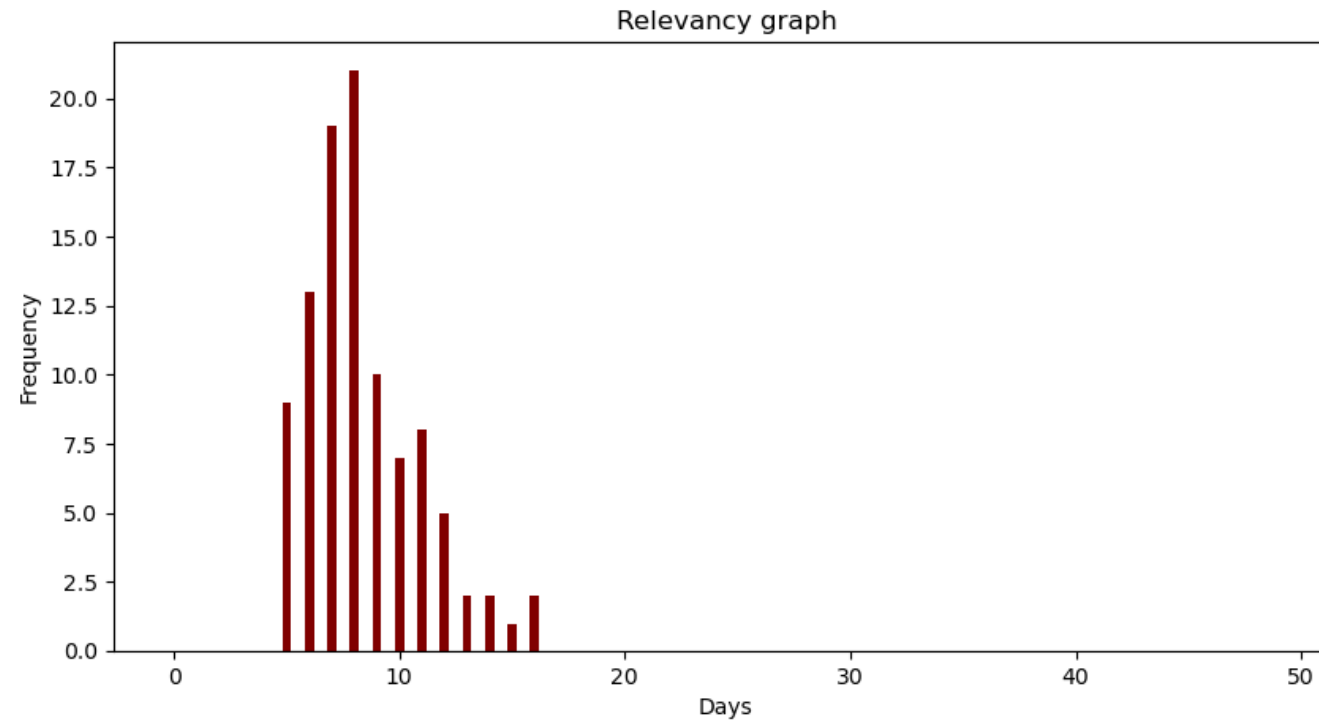
The minimum set for AFLAC Inc. is [7, 37, 39, 42, 44, 47]

The minimum set for American International is [3, 6, 11, 13, 18, 22, 28, 37, 39, 42, 44]

The length of the minimal set was 8 on average (8.4) ranging from 5 to 16, out of a total of 49 input features.

# Minimum Hitting Set Size

---



# Bibliography

---

<https://github.com/alexeyignatiev/xreason>

<https://ojs.aaai.org/index.php/AAAI/article/view/3964/3842>

<https://pytorch.org>

<https://www.nyse.com/index>

<https://finance.yahoo.com>

<https://www.geeksforgeeks.org>

Emails with Prof. Alexey Ignatiev and Nina Narodytska