

COMPUTATIONAL FINANCE

# SUMMARY REPORT

**Paper :** Nowcasting Stock Implied Volatility with Twitter  
(<https://arxiv.org/pdf/2301.00248.pdf>)

---

Under the guidance of Prof. Diganta Mukherjee

By Hemant Banke (MD2107)

# Summary

## Introduction

This study uses random forests to forecast the end-of-day implied volatility of stock prices for the following day. It also assesses the efficacy of various predictor sources and highlights the importance of attention and sentiment variables mined from Twitter. It finds that equities in some industries are more predictable than others, maybe as a result of too much social media attention or insufficient option liquidity. It also identifies underlying market regimes in implied volatility that influence the performance of the approach.

The development of mass media makes it easier for information to spread quickly, perhaps boosting market pricing efficiency. The majority of the existing research has concentrated on the Twitter platform and its impact on the three key financial indicators of stock price, realized volatility, trading volume, or a combination of these. One of the more crucial metrics in the realm of derivatives is the market implied volatility, which is generated from option prices. It is a risk indicator that shows how much risk the market anticipates a specific asset will display in the future. Eg: VIX.

## Methodology

Instead of picking a small number of stocks at random for this study, the authors of the paper diversified the stock universe throughout the 11 typical US stock market sectors, picking 15 stocks each, resulting in a total of 165 stocks. Data used for each stock is : closing price, end-of-day 30-day IV using VIX method, two numerical features from textual Twitter corpus: end-of-day total tweet publication count and end-of-day average sentiment polarity (sentiment analysis using VADER). Two additional predictors were created for each feature to capture the temporal information contained in the original feature time series: the daily difference (first-order difference) and the difference between the daily value and its EMA of the last 10 trading days. The data was collected from January 1st, 2011 through March 1st, 2019.

The target variable is a binary variable which is 1 if the 30-day IV of the next day is higher than today, and 0 otherwise. A Random Forest model is built ,with 1000 trees and hyper parameters tuned from a grid of values, to predict the target based on the temporally ordered data. The tuning is performed using walk-forward validation, a cross-validation technique designed specifically for temporally ordered data, and performance is measured using ROC AUC.

The performance of the strategy is measured under various market regimes to assess how well it stands up to the financial market's change over time. Using a hidden Markov model, the paper quantifies market regimes as various states in mean IV. A distinct HMM model was trained on each stock's end-of-day IV time series, which covered the period from

January 1st 2007 to December 31st 2012, and was then used out-of-sample. Four different regimes were identified corresponding to low, medium, high, and very high mean IV.

## Results

### *Feature Importance*

Seven RF models were trained on data obtained from 7 possible subsets of feature sources for all 165 stocks and the performance was compared against stratified dummy classifiers, spanning an out-of-sample period from January 1st, 2013 till March 1st, 2019. Median AUC is used as the performance metric.

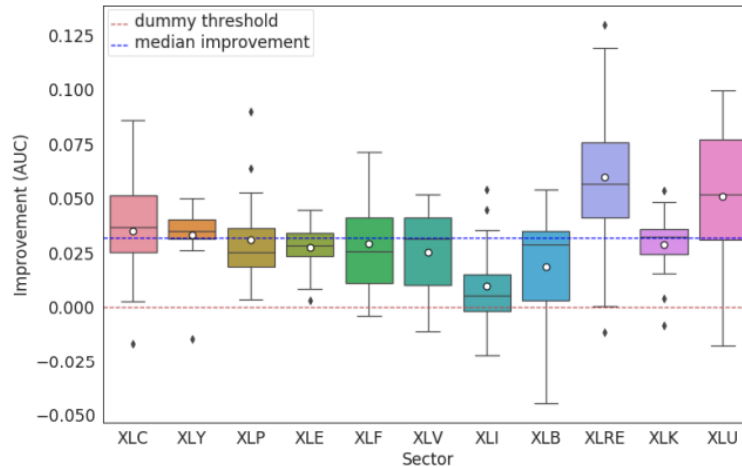
Scenario	Feature Source	Features
1	Stock Price	2
2	Stock Price, Tweets	8
3	Implied Volatility	3
4	Implied Volatility, Tweets	9
5	Tweets	6
6	Stock Price, Implied Volatility	5
7	Stock Price, Implied Volatility, Tweets	11

The outcomes offer proof that end-of-day IV movements can be predicted. Every scenario performs better than a random classifier, and scenarios with IV feature outperform the stratified dummy classifier. A superior median performance was always obtained when using characteristics extracted from tweets. This suggests that data from Twitter and future IV do interact in a predictive way. The best overall result was also obtained by utilizing all features, indicating that there are predictive patterns among all three feature sources.

### *Predictive Performance across Sectors*

The figure below shows Improvement in median AUC for RF model considering all features against stratified dummy classifier. The methodology is generally able to beat the stratified dummy classifier across all different sectors and the dummy classifier on 148 stocks. The approach does significantly better on XLRE (Real Estate) and XLU (Utilities). In contrast, XLI (Industrials) and XLB (Materials) seem to lack in performance.

Compared to other sectors, the equities in XLRE and XLU are substantially less liquid. Also, there is a modest negative association between option liquidity and predictive performance, suggesting that less liquid stocks are simpler to predict. The authors propose that the low liquidity of both XLRE and XLU in the options market may be followed by a less effective price discovery process, which may make these markets reflect new information more slowly and make them simpler to predict using the information at hand.



### ***Effect of Social Media***

There is a very strong correlation between attention on Twitter and liquidity. Suggesting that prediction is easier on stocks that are more popular on social media. Comparing the Improvement in AUC by considering Twitter features among others suggest most sectors seem to have a sizable number of stocks that benefit from using social media features. This can be explained by : social media inciting herd behavior and emotional reactions among investors, possibly driving inefficiency or the sentiment extraction technique on tweets is imperfect making stocks with more tweets more predictable because of more data.

### ***Performance across Implied Volatility Regimes***

	Frequency	Implied Volatility	Dummy AUC	Improvement
Low	392	18.6	50.75	+3.15
Medium	452	22.3	50.65	+3.83
High	411	26.7	52.12	+3.84
Very High	231	35.3	54.92	+1.47

The paper looks at performance variability of models trained on all features across four different market regimes in implied volatility, identified by using a HMM. The Table shows the median of all results across 165 stocks. The approach outperforms dummy classifiers in all regimes, the added value seems to be most significant in the low to high regimes. This can be due to correlation between predictive performance and frequency of data from the regime.