

# Cryptocurrencies Price Prediction Using News and Social Networks Data

---

Arnaud Autef  
Catherine Gu  
Olivier Pham  
Charbel Saad  
Benoît Zhou



















---

# Overview

1. Problem description
2. Data Collection
3. Feature Engineering
4. Model & Evaluation
5. Conclusion & Takeaways

# Overview

## Top 100 Cryptocurrencies by Market Capitalization

Cryptocurrencies ▾		Exchanges ▾		Watchlist		USD ▾	Next 100 →	View All
#	Name	Market Cap	Price	Volume (24h)	Circulating Supply	Change (24h)	Price Graph (7d)	
1	 <b>Bitcoin</b>	\$141,491,248,675	\$7,976.12	\$23,153,242,863	17,739,362 BTC	-7.40%		...
2	 <b>Ethereum</b>	\$26,268,025,515	\$247.00	\$9,872,839,695	106,346,402 ETH	-7.68%		...
3	 <b>XRP</b>	\$17,376,386,250	\$0.411938	\$2,553,674,556	42,181,995,112 XRP *	-8.36%		...
4	 <b>Bitcoin Cash</b>	\$7,046,271,276	\$395.44	\$2,202,815,324	17,818,775 BCH	-9.79%		...
5	 <b>Litecoin</b>	\$6,438,460,335	\$103.73	\$4,019,675,150	62,068,951 LTC	-9.20%		...
6	 <b>EOS</b>	\$6,138,276,865	\$6.69	\$4,014,765,857	917,678,573 EOS *	-11.02%		...
7	 <b>Binance Coin</b>	\$4,302,042,406	\$30.47	\$419,528,781	141,175,490 BNB *	-5.81%		...
8	 <b>Bitcoin SV</b>	\$3,839,016,605	\$215.47	\$1,108,711,999	17,816,861 BSV	-1.95%		...
9	 <b>Tether</b>	\$3,129,805,000	\$0.997987	\$24,228,629,558	3,136,118,221 USDT *	-0.55%		...

---

## Problem description

- Low barrier to entry for trading cryptocurrencies
- Large tradable market by dollar value
  - >2200 cryptocurrencies
  - 19000 markets trading crypto
  - \$250bn market cap
- Is there a relationship between social network/news data and cryptocurrencies price?

---

# Research Roadmap

## Data Collection

- Market Data from Coinbase
- Google Trends search volume
- Twitter textual data



## Exploratory Analysis

- Google Trend volume // BTC volume
- BTC keyword lagged // BTC volume
- Keyword volume // volatility



## Directional prediction

- Sentiment Analysis
- Directional imbalances
- Condition on volatility



## Model & Evaluation

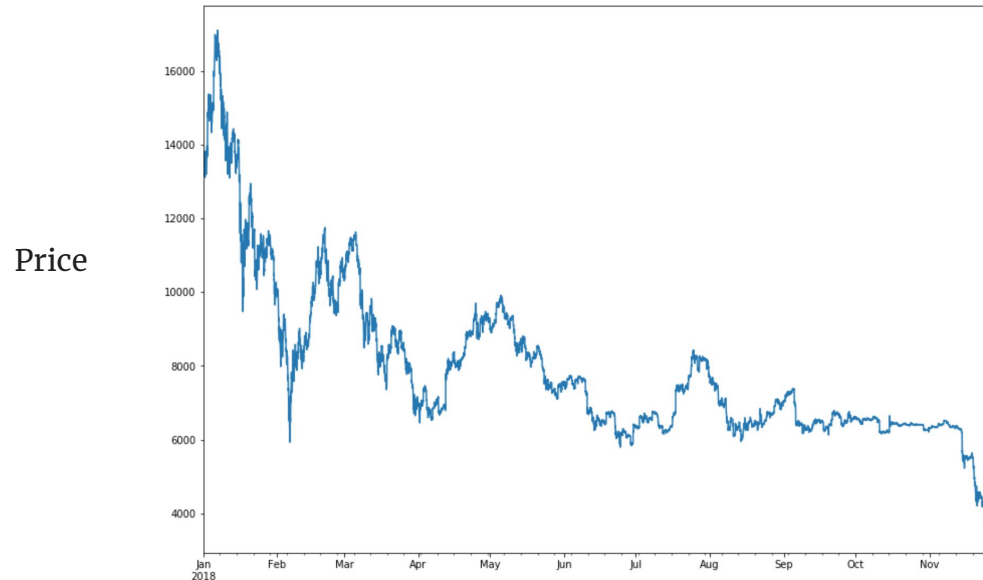
- Build a linear model
- Roll-forward Backtest

# Data Collection

---

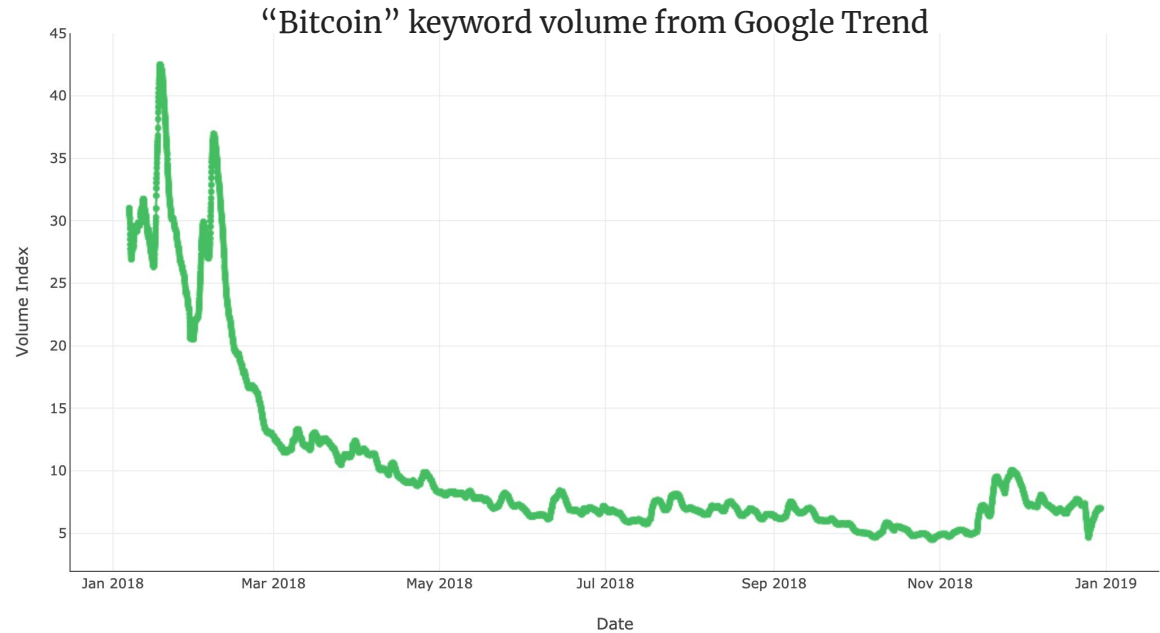
# Market Data Collection

- LO data from **Coinbase API** of currency pairs for 2018
- Different market behavior



# Google Trends Data

## Example: Google Trend on “*bitcoin*”

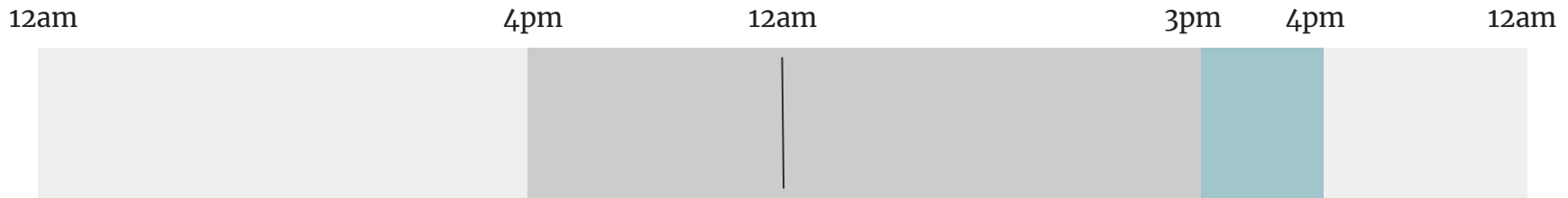




---

# Twitter Textual Data

- Used the “GetOldTweet” repository
- Returns all tweets in **chronological order** between 4pm to 4pm of the next day: **too much data**
- Selected only 1000 tweets between **3pm and 4pm**



# Feature Engineering

# Tweets:

- Some very informative tweets but...
- Not structured
- Noisy

Thread

**Guy Swann** ⚡ @TheCryptoconomy

!! #BCH /#bcash was hit by 51% attack from just 2 miners, [BTC.TOP](#) & [BTC.com](#) - & no one seems to be talking about it. 🤔

Thread 🗨️

1/ What I've gathered from loose details:  
First, there was an unintentional split with the recent #BCH "upgrade."

9:37 AM · May 24, 2019 · Twitter Web Client

900 Retweets 1.9K Likes

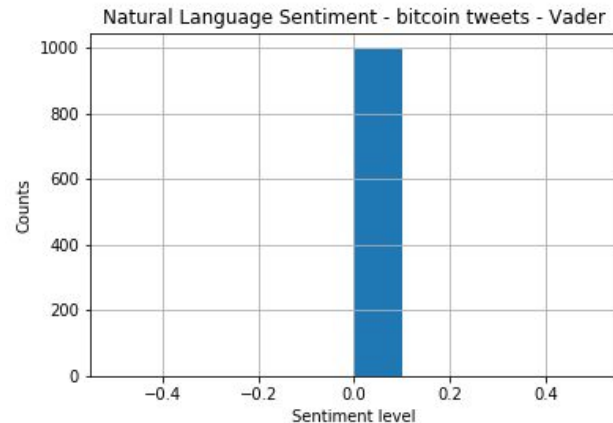
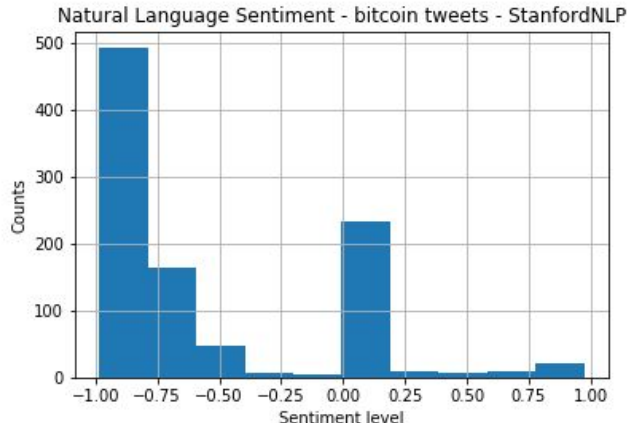
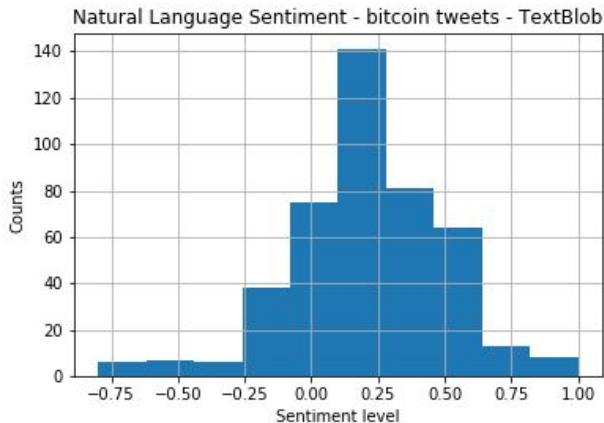
**Guy Swann** ⚡ @TheCryptoconomy · May 24  
Replying to @TheCryptoconomy

2/  
Since the original split in 2017, there has been a significant number of coins accidentally sent to "anyone can spend" addresses (due to tx compatibility of sigs, but no #SegWit on #BCH), or possibly they've been replayed from

# Sentiment Analysis for tweets

Compared 3 libraries on 1k tweets:

- *TextBlob* (naive Bayes classifier)
- *StanfordNLP*
- *Vader*



# Keyword volume

- Counted positive and negative keywords volume in Tweets and Google Trends volume
- *Imbalance =*

*# positive words - # negative words*

## Positive Keywords:

'conviction', 'bold', 'up', 'buy', 'bullish', 'bull', 'free money', 'long', 'rise', 'boom', 'bid', 'support', 'grow', 'get', 'make', 'earn', 'investment', 'invest', 'investing', 'invested', 'buying', 'bought', 'pump', 'like', 'skyrocket',

## Negative Keywords:

'scam', 'capitulation', 'down', 'fork', 'sell', 'short', 'bear', 'bearish', 'bubble', 'stop', 'crash', 'clamp', 'shut', 'freeze', 'fall', 'bust', 'trash', 'forbid', 'oppose', 'dash', 'sold', 'selling', 'collapse', 'plummet', 'plunge'

---

# Predictor Construction for Twitter

	Positive Tone	Negative Tone
Positive Imbalance	Buy	Sell
Negative Imbalance	Sell	Buy

**Keyword Score** = *Imbalance*  $\times$  *sign(sentiment score)*

---

# Twitter Predictors

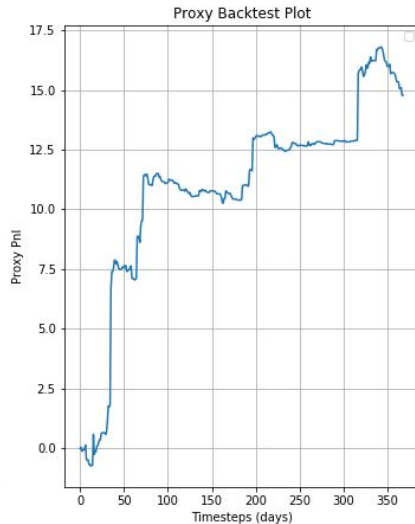
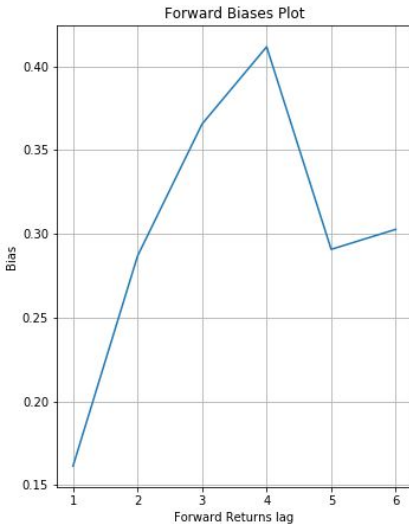
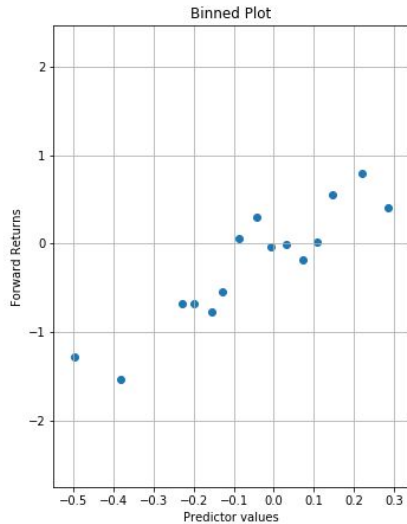
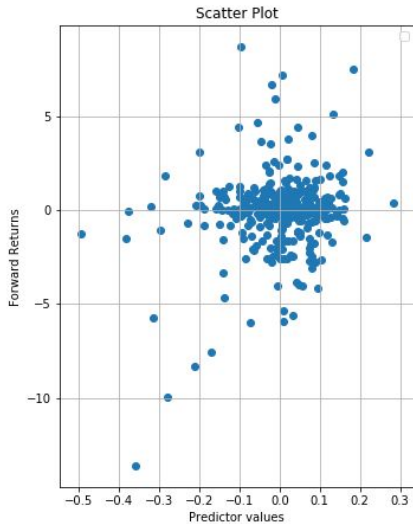
- **Buy** signal: Tweets with **positive Sentiment Analysis score** over the last hour.
- **Sell** signal: Tweets with **negative Sentiment Analysis score** over the last hour.
- Normalized signal: 
$$\frac{\text{Buy} - \text{Sell}}{\text{Buy} + \text{Sell}}$$

We expect signals to work best when using their discrete derivative, so we subtract them from their moving average.

# Plots & Metrics

$$\text{bias}(\text{lag}) = \text{corr} [\text{return}(t + \text{lag}), \text{predictor}(t)] \cdot \sigma(t + \text{lag})$$

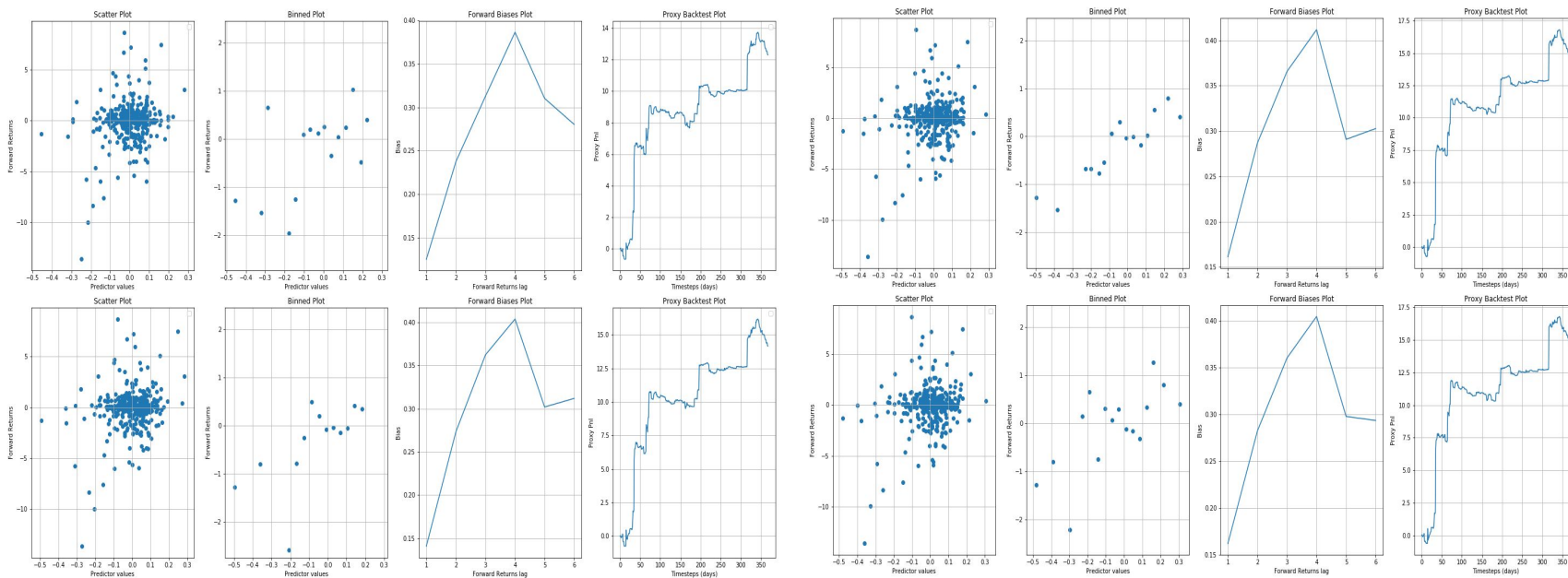
$$\text{lag}^* = \text{argmax}_{\text{lag}}(\text{bias}) = 4\text{hours}$$





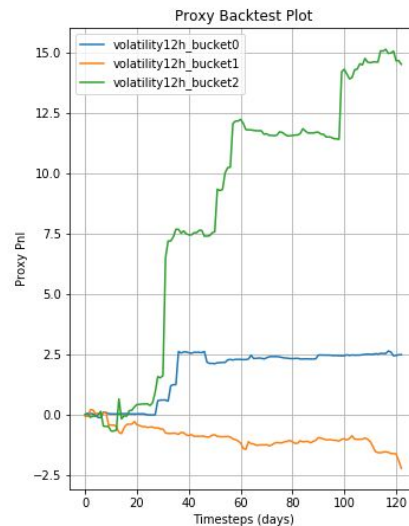
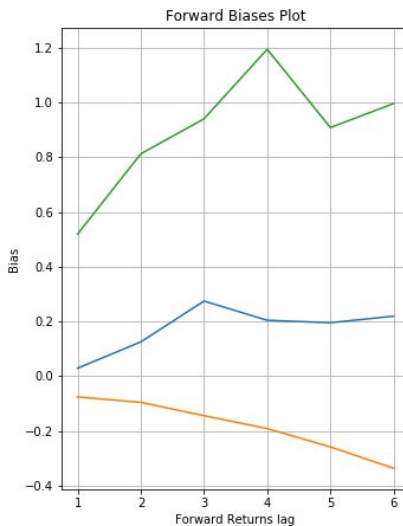
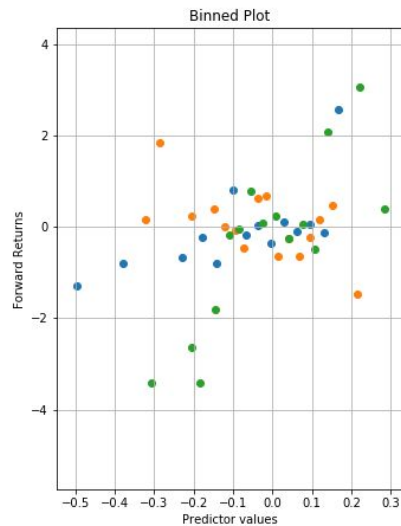
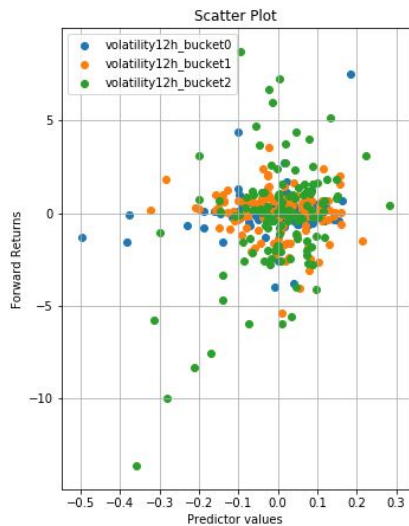
# Parameters selection

- Moving average: 3, 5, 7 & 9 days
- Returns horizon



# Final predictor

Normalized predictor + Conditioning on volatility regime



High

Medium

Low

---

# Predictor Construction for Google Trends



---

Signal for Google Trends:

Positive Volume – Negative Volume

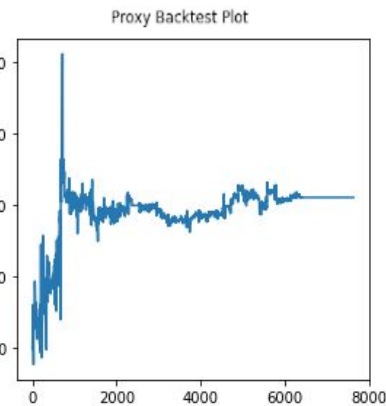
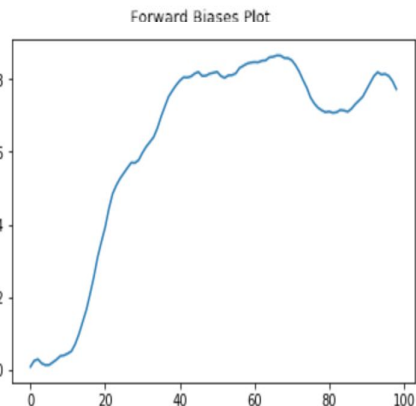
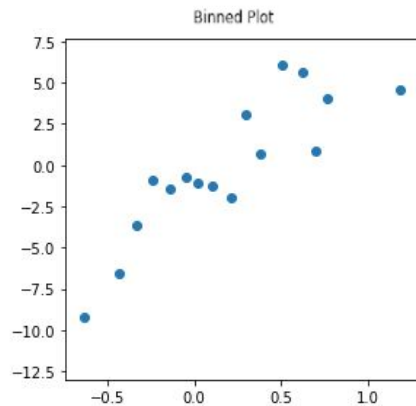
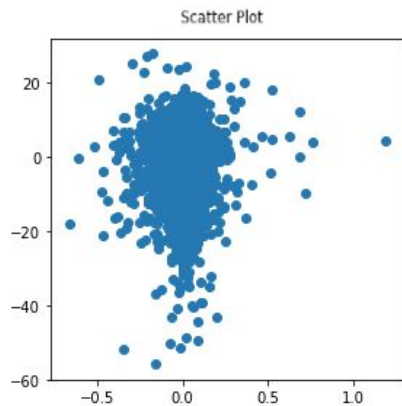
---

Total Volume

Centered with the ewm

# Google Trends Predictors

- Positive correlation
- Longer horizon (2 days)
- Limited prediction power after the beginning of the year



# Model & Evaluation

---

## Model

$$Y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

- No intercept
- $Y$  = BTC-USD return 4 hours ahead of 4pm
- $X_i$  = sentiment feature conditioned on past 12 hour volatility regime: low/median/high
- Moving average of 7 days
- In-sample  $R^2 = 0.041$

# Evaluation

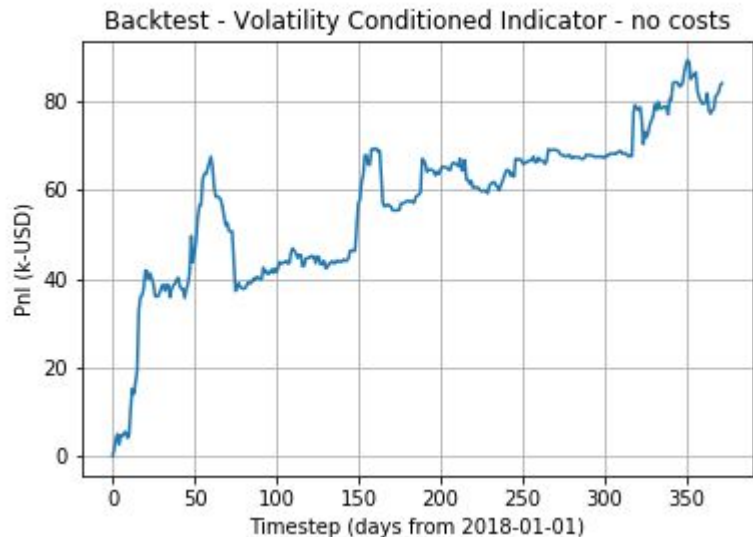
Roll-forward backtest of two weeks span

Week 1	Week 2	Week 3	Week 4				
Week 1	Week 2	Week 3	Week 4	Week 5	Week 6		
Week 1	Week 2	Week 3	Week 4	Week 5	Week 6	Week 7	Week 8

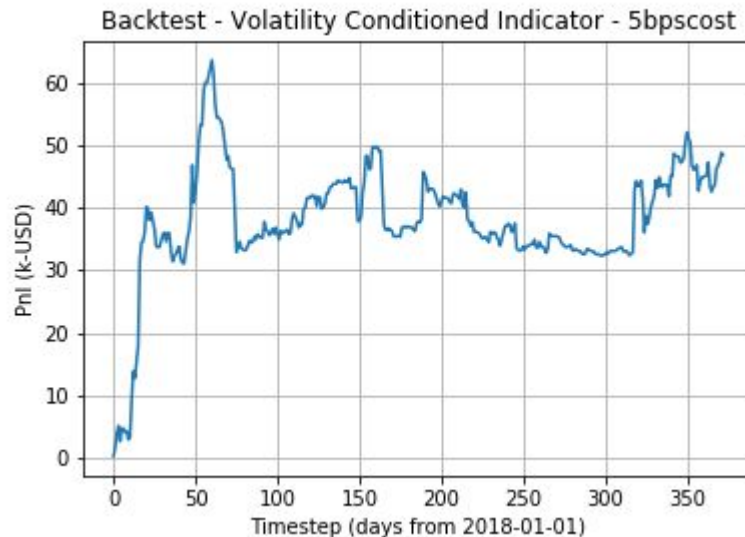
- Threshold based trading strategy
- Threshold  $\Gamma$  = trading cost
- Linear trading costs

# Evaluation

## Roll-forward backtest of two weeks span



no trading cost



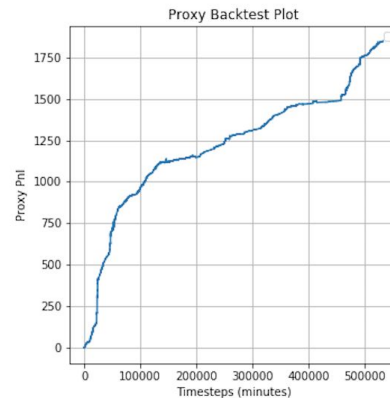
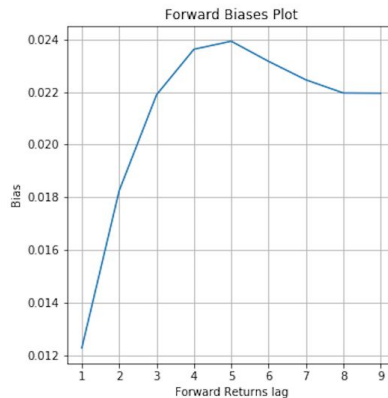
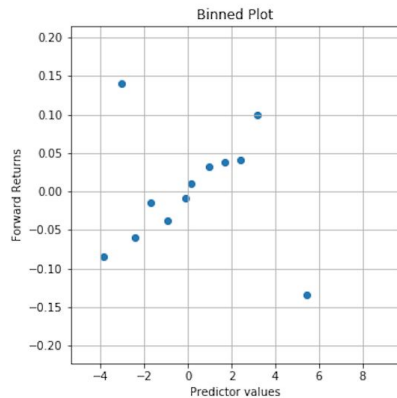
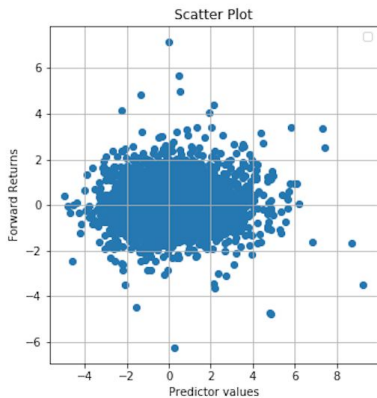
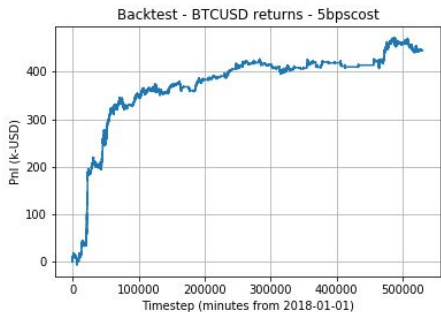
trading cost = 5 bps

Sharpe ratio with no trading cost:  $S = 2.44$



# Market Data based model

- Extension: Market data for higher frequency strategies
- Lead lag Bitcoin → Litecoin 5 min
- Sharpe ratio of 38.5



# Conclusion & Takeaways

- 
- Choice of sentiment analysis package is important
  - Engineering challenges: quality of sentiment data, feature construction
  - Filtering methods (e.g. volatility conditioning) produces more reliable results and align with the fundamentals observed in finance
  - Portfolio construction: accounting for trading cost, introducing alpha threshold to determine position turnover, imitating real trading environment
  - Extension: combining our signals to produce lower risk reliable strategies

# Questions