

# Bitcoin Tweet Sentiment Analysis

...

By Natalya Doris  
Flatiron School, Data Science  
Phase 4

# Business Problem

In this project, we help Crypto Consultancy firm X build a model that can accurately **classify** Tweets about Bitcoin as having either positive or negative **sentiment**. Firm X would like to use the model to classify **unlabeled** Tweets and **understand** any characteristics distinguishing the positive and negative Tweets.



# Breaking it Down

## Explore the Data

- Understand **trends** and key **takeaways** to bring to the client
- **Clean** the text data so that it is ready to be used in modeling

## Build a Model

### Test and tune:

- Three vectorizers
- Four classification models

## Evaluate the Model

**Score** each model, selecting a final **best-performing** model

# **The Data:**

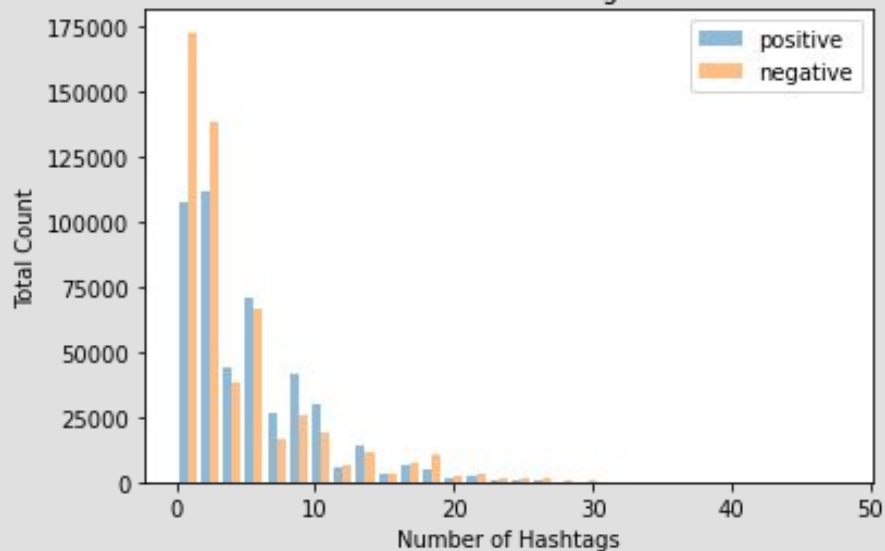
**One million Tweets referencing  
Bitcoin, spanning a six-month  
period from February 2021 to  
August 2021**



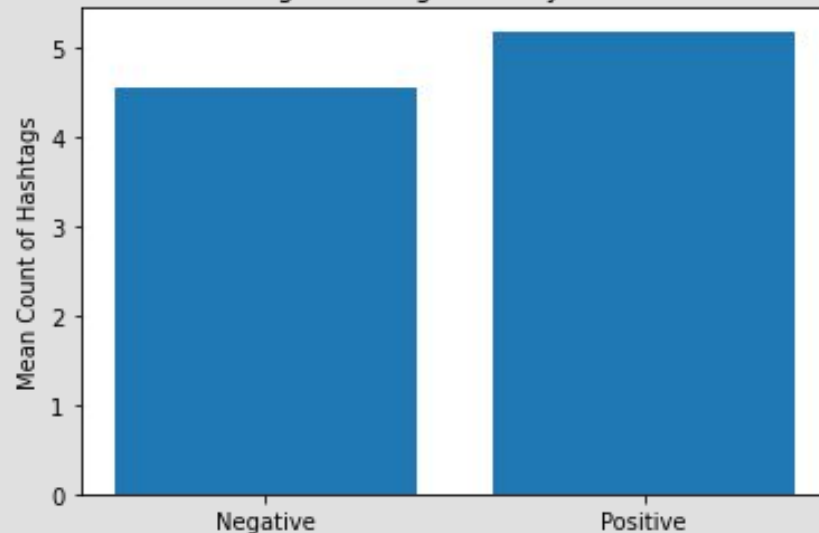


# Hashtags

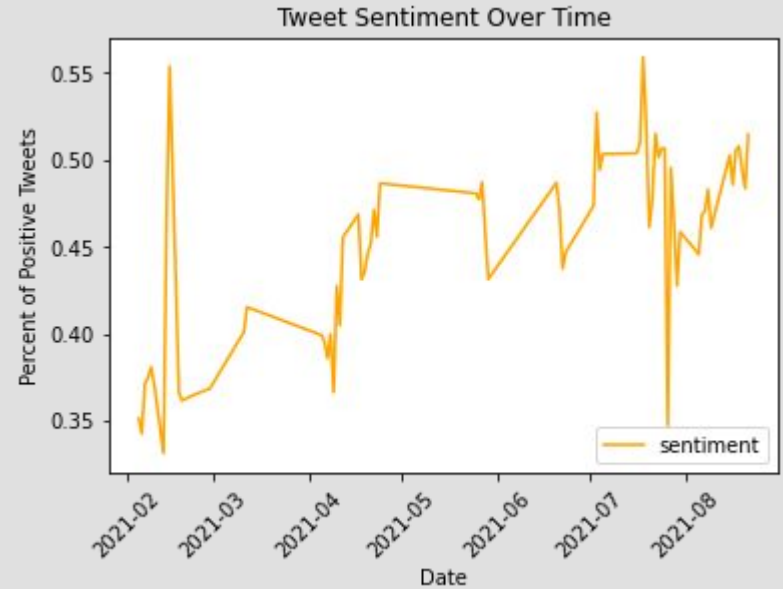
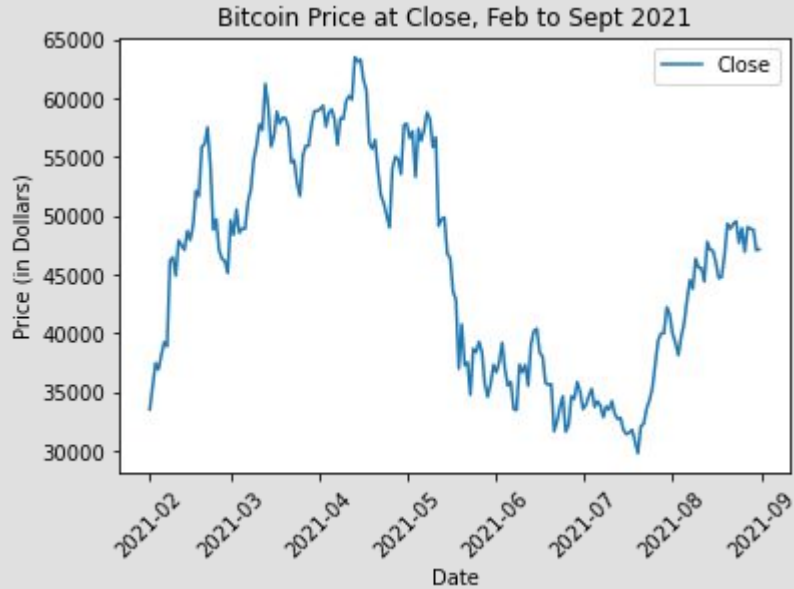
Distribution of Hashtag Count



Average Hashtag Count by Sentiment

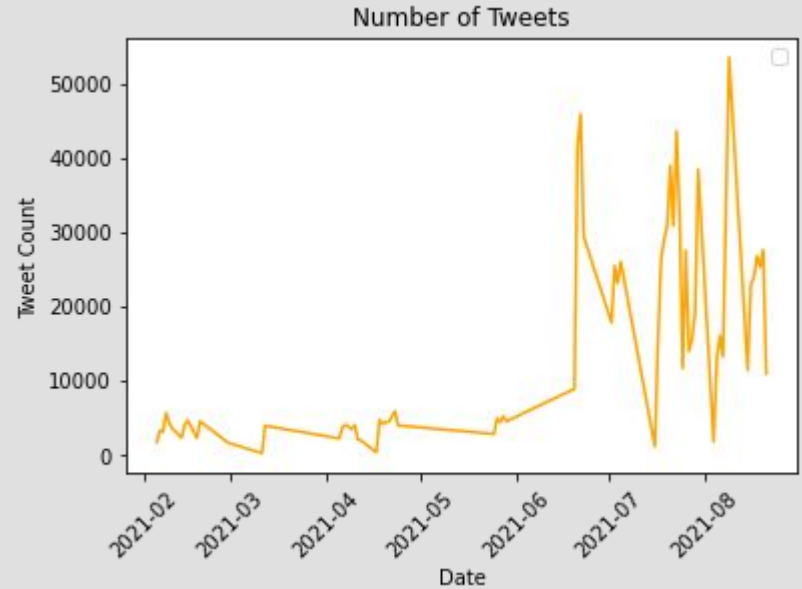
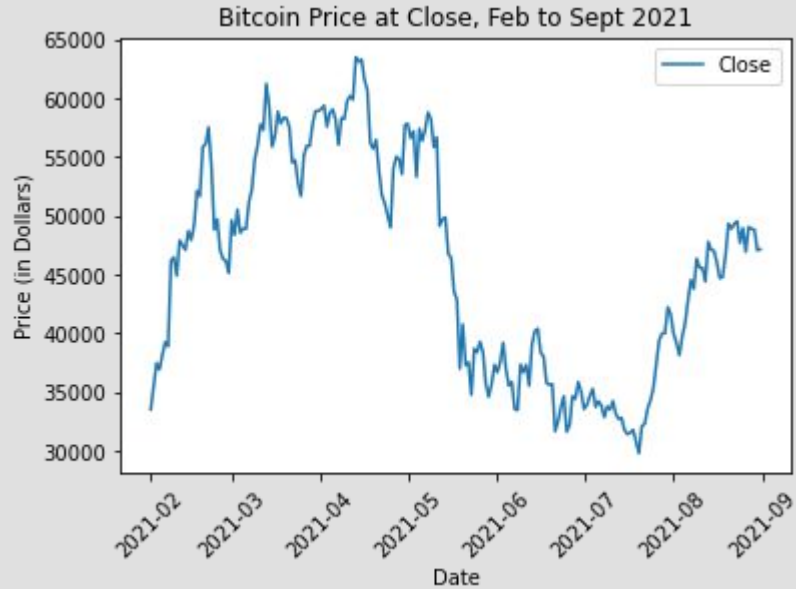


# #Trending





# #Trending



# Modeling

|                 | f1_train | f1_test  | accuracy_train | accuracy_test | roc_auc  | pr_auc   | model                | vectorizer                                |
|-----------------|----------|----------|----------------|---------------|----------|----------|----------------------|---|
| <b>baseline</b> | 0.575926 | 0.575227 | 0.619592       | 0.619370      | 0.615580 | 0.547573 | MultinomialNB()      | TfidfVectorizer(max_features=10)          |
| <b>base_tf</b>  | 0.873472 | 0.853652 | 0.883644       | 0.865230      | 0.863466 | 0.809588 | MultinomialNB()      | TfidfVectorizer()                         |
| <b>base_cv</b>  | 0.967555 | 0.858882 | 0.969429       | 0.869114      | 0.867714 | 0.813043 | MultinomialNB()      | CountVectorizer(ngram_range=(2, 2),\n ... |
| <b>base_cv2</b> | 0.958314 | 0.867798 | 0.960734       | 0.875402      | 0.874822 | 0.817498 | MultinomialNB()      | CountVectorizer(ngram_range=(2, 2))       |
| <b>logreg</b>   | 0.991206 | 0.917310 | 0.991740       | 0.925155      | 0.922720 | 0.901446 | LogisticRegression() | CountVectorizer(ngram_range=(2, 2))       |
| <b>svc</b>      | 0.999980 | 0.926771 | 0.999981       | 0.933035      | 0.931121 | 0.909578 | LinearSVC()          | CountVectorizer(ngram_range=(2, 2))       |

**Final model is Linear SV with Count Vectorizer, achieving ~93% accuracy and F1 score!**

# Conclusions

## Conversation Matters

- Period of higher Twitter volume associated with period of lower Bitcoin price

## Sentiment Impactful

- Spike in positive sentiment occurred around same time as rise in Bitcoin price

## #Positive

- Positive sentiment Tweets tend to have more hashtags on average

## Model Accuracy

- Final model can classify unlabeled Tweets as positive or negative with ~90% accuracy

# Next Steps

- Pull more recent Tweets on Bitcoin via Twitter API
  - Run final model on unseen data, i.e. new Tweets
  - Use model-labeled Tweets to conduct Time Series Analysis, with the aim of understanding the predictive power of Tweet sentiment on the price of BTC
  - Use deep learning models to classify Tweets, comparing these results to previous model
-