

Capstone Project

Sai Indukuri
April 15, 2018

I. Definition

Project Overview

Cryptocurrency is a digital currency in which encryption techniques are used to regulate the generation of units of currency and verify the transfer of funds, operating independently of a central bank. There is lot of buzz about cryptocurrency and everyday in the news we hear volatility in the prices of cryptocurrencies such as Bitcoin.

Bitcoin is actual implementation of decentralization issued under the consent of participants not the central bank. Hence variants like purchasing power or interest rate parity does not impact Bitcoin. Bitcoin price is largely impacted by demand and supply. 16.5 million Bitcoins are mined since it was created and the Bitcoins are capped at 21 millions part of the reason for scarcity and rush. Hence to predict Bitcoin price we need to predict Demand and Supply.

Blockchain is the technology used to create Bitcoin. "Block Chain information" along with "Bitcoin" trading information includes features that can determine demand and supply and hence the price. Therefore, the project uses historic blockchain information and Bitcoin trading information to predict future Bitcoin prices. This project uses data from <https://www.kaggle.com/sudalairajkumar/cryptocurrencypricehistory> which is publicly available. The data consist of "Bitcoin_Dataset" which consist of historic Blockchain information and Bitcoin prices for last 4 years.

Bitcoin dataset consist of several important blockchain information that might help us understand rise and fall continuum , most importantly, Total bitcoins, Trade volume, block size, avg block size, Orphaned blocks, transactions, hashrate, Mining difficulty and Miners revenue. We have a balanced dataset with around 47% of records with rising prices and 43% of records with falling prices and 10% of records with no price change.

Problem Statement

The extremely nonlinear nature of the crypto market data makes it very difficult to design a system that can predict the future direction of the crypto prices in Bitcoin with sufficient accuracy. Goal of the project is to predict price movements(that's either up or down) of cryptocurrencies such as Bitcoin. As problem involves predicting and classifying Bitcoin prices to rise , fall or remain the same this is a multi class classification problem.

The main goal of speculating a cryptocurrency price is to make profitable trades. To make profitable trades one doesn't need to know accurate value of the cryptocurrency price, but one merely needs to predict whether prices rise or fall or remain the same. In these lines we will train our models to classify subsequent day's closing price to be likely higher, lower than the last or the same as last, based on past days price data.

The solution involves

1. Preprocessing the data
2. Selecting features that can most help find the bitcoin trend
3. Testing the data on selected features on varying time series data using the models SVM, Random Forest and Logistic regression
4. Select the best model based on the models performance measured by Accuracy, Precision and FBeta score.
5. Fine tune the model to improve accuracy.

The final solution will be able to predict the bitcoin price trend(Up, Down, Constant) with accuracy better than that of Naive prediction.

Metrics

With Bitcoin predictions it is very important to Accurately predict price trends. Hence using accuracy as a metric for evaluating a particular model's performance would be appropriate. Additionally, Identifying price rise when price actually drops is detrimental as the investors will lose money. Therefore, a model's ability to precisely predict price rise is more important than the model's ability to recall . We can use F-beta score as a metric that considers both precision and recall:

In particular, when $\beta=0.5$, more emphasis is placed on precision. This is called the $F_{0.50.5}$ score (or F-score for simplicity).

In addition to FBeta score, we will also use Precision to measure the performance. Below are the formulas for each of the metric.

Accuracy = True positives/Total observation

Precision = True positives/(True positives+False positives)

FScore = $(1+\beta*\beta) * (\text{precision}.\text{recall}) / ((\beta*\beta*\text{precision})+\text{recall})$

II. Analysis

Data Exploration

Bitcoin is actual implementation of decentralization issued under the consent of participants not the central bank. Hence variants like purchasing power or interest rate parity does not impact Bitcoin. Bitcoin price is largely impacted by demand and supply. 16.5 million Bitcoins are mined since it was created and the Bitcoins are capped at 21 millions part of the reason for scarcity and rush.

Hence to predict Bitcoin price we need to predict Demand and Supply. Blockchain is the technology used to create Bitcoin. "Block Chain information" along with "Bitcoin" trading information includes features that can determine demand and supply and hence the price. Therefore, the project uses historic blockchain information and Bitcoin trading information to predict future Bitcoin prices. This project uses data from <https://www.kaggle.com/sudalairajkumar/cryptocurrencypricehistory> which is publicly available. The data consist of "Bitcoin_Dataset" which consist of historic Blockchain information and Bitcoin prices for 8 years from 2009 to 2017.

Bitcoin dataset consist of several important blockchain information that might help us understand rise and fall continuum , most importantly, Total bitcoins, Trade volume, block size, avg block size, Orphaned blocks, transactions, hashrate, Mining difficulty and Miners revenue.

We have an Unbalanced dataset with around 47% of records with rising prices and 43% of records with falling prices and 10% of records with no price change. Also the dataset contained few blank values and needs to be cleaned.

Below is a sample of dataset with features and values.

| | |
|------------------------------|-----------------|
| Date | 2009-11-10 |
| btc_market_price | 0 |
| btc_total_bitcoins | 1339450 |
| btc_market_cap | 0 |
| btc_trade_volume | 0 |
| btc_blocks_size | 0 |
| btc_avg_block_size | 0.0002154225352 |
| btc_n_orphaned_blocks | 0 |
| btc_n_transactions_per_block | 1 |

| | |
|---|----------|
| btc_median_confirmation_time | 0 |
| btc_hash_rate | 3.53E-06 |
| btc_difficulty | 1 |
| btc_miners_revenue | 0 |
| btc_transaction_fees | 0 |
| btc_cost_per_transaction_percent | 0 |
| btc_cost_per_transaction | 0 |
| btc_n_unique_addresses | 71 |
| btc_n_transactions | 71 |
| btc_n_transactions_total | 26958 |
| btc_n_transactions_excluding_popular | 71 |
| btc_n_transactions_excluding_chains_longer_than_100 | 71 |
| btc_output_volume | 3550 |
| btc_estimated_transaction_volume | 0 |
| btc_estimated_transaction_volume_usd | 0 |

Predominantly features can be classified in to

1. Demand features: Features that indicate bitcoin demand..

| |
|---|
| btc_market_price |
| btc_total_bitcoins |
| btc_market_cap |
| btc_trade_volume |
| btc_transaction_fees |
| btc_n_transactions |
| btc_n_transactions_total |
| btc_n_transactions_excluding_popular |
| btc_n_transactions_excluding_chains_longer_than_100 |
| btc_output_volume |
| btc_estimated_transaction_volume |
| btc_estimated_transaction_volume_usd |
| btc_cost_per_transaction_percent |
| btc_cost_per_transaction |

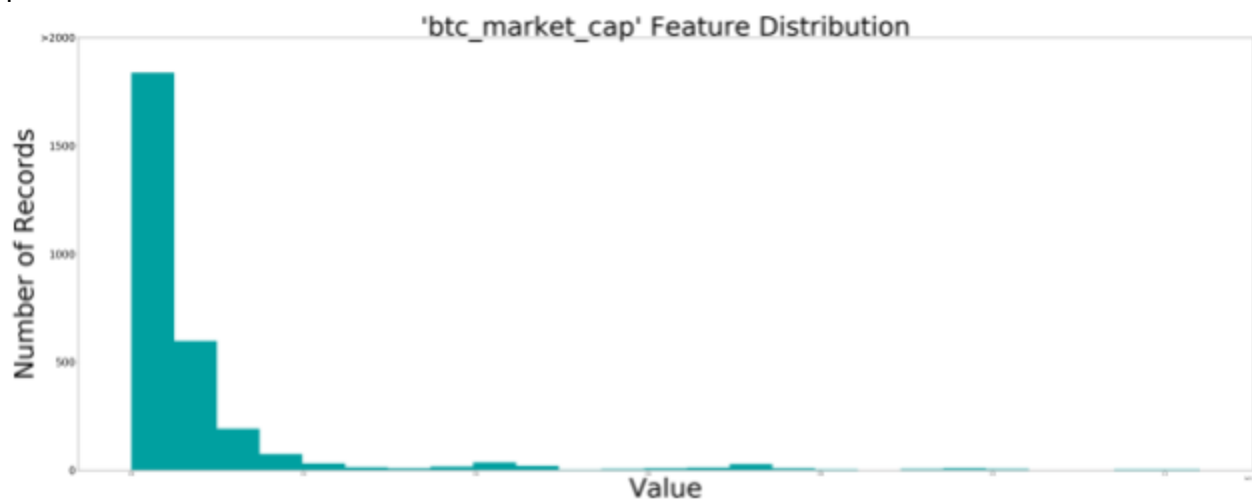
2. Supply features : Features that indicate bitcoin supply.

| |
|------------------------------|
| btc_blocks_size |
| btc_avg_block_size |
| btc_n_orphaned_blocks |
| btc_n_transactions_per_block |
| btc_median_confirmation_time |
| btc_hash_rate |
| btc_difficulty |
| btc_miners_revenue |
| btc_n_unique_addresses |

Bitcoin price can be predicted by analysing changes in Demand and supply continuum and above features have a role in defining the demand and supply.

Exploratory Visualization

Most of the features in the dataset are continuous features. Visualizations will help determine the skewness(if any) in the data and based on the below visualization most of the features appears extremely skewed. Below is the visualization of one of the features that shows skewness of the feature. Visualization of rest of the features can be found within the Bitcoin prediction code.



Skewness is a relevant characteristic extracted from the visualization as skewed data needs to be log transformed to ensure accurate performance of the algorithm.

Algorithms and Techniques

I choose to train the data with Support Vector Machines(SVM) and Random Forest algorithms. Below is the brief definition, strengths and weaknesses of these algorithms, reason I picked the algorithm and techniques applied on the algorithm..

Support Vector Machines(SVM):

Support vector machine uses a kernel mechanism that maximizes distance between closest member of separate classes.

Strengths: Unlike Logistic regression, SVM can model non linear decision boundary and have many kernels to choose from; SVMs are fairly robust against overfitting especially in high dimension space.

Weaknesses: Isn't suited to large dataset as training time might be high; Less effective with noisier dataset with overlapping classes.

Support vector machines may be a good fit for the problem since the dataset is small(~2000)

Support vector machine performs the classification by finding the hyperplane that differentiate the three classes(Up, Down, NoChg) very well. Bitcoin prediction is a multi class problem as we need algorithm to classify multiple classes. Within Support vector machines SVC, NuSVC and LinearSVC are the classifiers capable of performing multi class classification. As bitcoin analysis is complex, I used SVC instead of LinearSVC since linear SVC can only create linear kernel. Between SVC and NUSVC I chose SVC for its simplicity and also I do not want to control number of support vectors.

SVM by definition is well suited for binary classification. In order to perform multi-class classification, the problem needs to be transformed into a set of binary classification problems.

There are two approaches to do this:

One vs. Rest Approach (OvR): This strategy involves training a single classifier per class, with the samples of that class as positive samples and all other samples as negatives. This strategy requires the base classifiers to produce a real-valued confidence score for its decision, rather than just a class label; discrete class labels alone can lead to ambiguities, where multiple classes are predicted for a single sample.

One vs. One Approach (OvO): In the one-vs.-one (OvO) strategy, one trains $K(K - 1)/2$ binary classifiers for a K-way multi-class problem; each receives the samples of a pair of classes from the original training set, and must learn to distinguish these two classes. At prediction time, a voting

scheme is applied: all $K(K - 1)/2$ classifiers are applied to an unseen sample and the class that got the highest number of “+1” predictions gets predicted by the combined classifier.

SVC implements “One vs One” approach for multi class classification. Bitcoin is a multiclass problem with 3 classes. SVC algorithm would train 3 classifiers($3/(3-1)/2$) where each classifier receives samples with pair of classes(Up/Down,Up/NoChange,Down/NoChange) from the training set.As the classifier classifies the data a class with maximum number of votes is chosen.

Ensemble Methods:Random Forest :

Real World application: Most real world applications use some sort of ensemble learning as ensemble learning is a way of combining the output of multiple models.

Strengths: The major strength of Random Forest is that it is very difficult for the model to overfit, because the model only improves as you increase the complexity of the model

*Weakness:*The major weaknesses to do with Random Forest are related to increased storage and computation. Since RandomForest require multiple classifiers, there is a need for more storage for the classifiers which can be costly. In a similar vein, all of the classifiers need to be processed, rather than just one, which means that the run time of the algorithm is higher

RandomForest may be a good fit since we need to test multiple hypothesis.

I used RandomForestClassifier to train the bitcoin dataset. RandomForestClassifier would train as follows.

1. I used max_features as auto. Hence, the algorithm randomly selects $\sqrt{\text{total features}}$. As discussed in implementation we finally selected 14 features and hence approximately 4 ($\sqrt{14}$) features are randomly selected.
2. Among the 4 features , calculates the node based on best split
3. Splits the node in to daughter nodes using the best split
4. Repeats 1 to 3 steps until all leaves are pure
5. Builds forest by repeating steps 1 to 4 10 times by creating 10 trees

Once the RandomForestClassifier is trained it performs predictions as follows

1. Takes the test features and use the rules of each randomly created decision tree to predict the outcome and stores the predicted outcome
2. Calculate the votes for each predicted outcome
3. Consider the high voted predicted target as the final prediction from the random forest algorithm.

Benchmark

Looking at the distribution of classes (price rise(46.58%), drop(42.53%) and no change(10.89%)) there are very few records for 'No change' and a higher percentage of records with price drop, This can greatly affect **accuracy**, since we could simply guess the prices will go down all the time. Making such a statement would be called **naive**, since we have not considered any information to substantiate the claim. We will use a naive algorithm that would always predict price go down as Benchmark.

Following are the benchmark results for the Naive algorithm.

| Metrics | Accuracy | Precision | F Score |
|-------------------|----------|-----------|---------|
| Benchmark Results | 0.4253 | 0.4253 | 0.4806 |

III. Methodology

Data Preprocessing

Before data can be used as input for machine learning algorithms, it often must be cleaned, formatted, and restructured — this is typically known as **preprocessing**. Below are the preprocessing steps performed :

Cleaning:

As with any other realtime dataset bitcoin dataset came with its share of defects.

- There are few blank values within `btc_trade_volume`
There are 21 instances where `btc_trade_volume` is blank. `Btc_trade_volume` cannot be blank especially when other features have values. Trading volume information is not available during the dates and hence they were blanks. Missing values are handled by replacing them with mean of the `btc_trade_volume`.
- The records are not sorted on dates.
It is very important to order the records by dates with the time series problems like bitcoin estimation. The records in the input dataset are not sorted and needs to be sorted.

Transforming:

A dataset may sometimes contain at least one continuous feature whose values tend to lie near a single number, but will also have a non-trivial number of vastly larger or smaller values than that single number. Algorithms can be sensitive to such distributions of values and can underperform if

the range is not properly normalized. Bitcoin dataset has around 16 features that are continuous and that need to be transformed. Log transformation is applied on all these features.

Normalizing:

In addition to performing transformations on features that are highly skewed, it is often good practice to perform some type of scaling on numerical features. Applying a scaling to the data does not change the shape of each feature's distribution however, normalization ensures that each feature is treated equally when applying supervised learners. Note that once scaling is applied, observing the data in its raw form will no longer have the same original meaning, as exemplified below. Bitcoin dataset has 22 numerical features that are scaled and normalized.

Implementation

Implementation can be divided into 3 steps

Feature selection:

We can determine relevance of each of the input features by training a supervised random forest regression learner on the data. I have used 'Mean decrease impurity' within random forest regressor to identify the feature importance and feature importance is determined based on the features that have most information gain/entropy. Features are then ordered based on Feature importance and features are selected that would provide 80% of information Gain. Below are the 17 features shortlisted based on feature importance out of total 24 features.

| |
|---|
| btc_total_bitcoins |
| btc_trade_volume |
| btc_cost_per_transaction |
| btc_avg_block_size |
| btc_difficulty |
| btc_estimated_transaction_volume_usd |
| btc_n_unique_addresses |
| btc_transaction_fees |
| btc_n_transactions_per_block |
| btc_miners_revenue |
| btc_n_transactions_total |
| btc_n_transactions_excluding_chains_longer_than_100 |
| btc_median_confirmation_time |

| |
|-----------------------|
| btc_n_orphaned_blocks |
| btc_blocks_size |
| btc_hash_rate |
| btc_output_volume |

Creating a training and predicting pipeline:

To properly evaluate the performance of each model , it's important that we create a training and predicting pipeline that allows us to quickly and effectively train models using various sizes of training data and perform predictions on the testing data. We cannot use k-fold cross validation or any other general train test splits on Bitcoin dataset since bitcoin data is a time series data. We cannot randomly split timesplit data into groups. Instead, we must split data on temporal order in which values are observed.

We can either do a single split, by selecting an arbitrary split point in the bitcoin dataset and create two new dataset, one train other test, or do a multi split where in we repeat process of split multiple times. I chose multi split over single split since multi split provides more robust estimate of the expected performance of the SVC and RandomForestClassifiers.

I used TimeSeriesSplit to perform multi split and the algorithm splits the data into k fold split with an incremental superset of training data in every split. I chose to use 3 fold split since the dataset is small(~2000 records) and more than 3 splits may result in very small train-test dataset. So we have 3 sets of training and prediction pipeline each with 730 testing samples and 730,1460 and 2190(incremental superset of data) training samples.

Below is the Accuracy, Precision and FBeta score for each training sample on SVC and Random Forest

| SVC Classifier | Accuracy | Precision | FBeta Score |
|-----------------------|----------|-----------|-------------|
| 730 training samples | 0.02 | 0.00 | 0.00 |
| 1460 training samples | 0.50 | 0.25 | 0.28 |
| 2190 training samples | 0.4 | 0.16 | 0.18 |

| Random Forest Classifier | Accuracy | Precision | FBeta Score |
|--------------------------|----------|-----------|-------------|
| 730 training samples | 0.56 | 0.54 | 0.54 |

| | | | |
|-----------------------|------|------|------|
| 1460 training samples | 0.48 | 0.48 | 0.47 |
| 2190 training samples | 0.40 | 0.16 | 0.18 |

There are no major complications during the coding process except for the expected errors such as ill defined precision, which probably might have occurred as the algorithms might have been trained on a sample without positive samples. This is largely due to the imbalanced dataset.

Training and Testing the Model :

SVC and Random tree algorithms are time series trained and tested with the training data that is split as discussed in the previous section. Algorithms are trained with varying sample size 730,1460 and 2190 and tested on 730 samples each time.

Refinement

I used grid search (‘GridSearchCV’) to refine the Random forest algorithm.

Below are the parameters and scoring used to refine the algorithm.

Parameters

- Varying size of estimators(100,200,300) are used to allow random forest to test on varying number of trees.
- Max_depth is the maximum depth of the tree and controls the depth to which nodes need to be expanded. I used 150,155 and 160 as max_depth.
- Max_features of ‘sqrt’ will allow us to select sqrt of all features while searching for a best split. This will limit number of features considered for split and thus avoid overfitting.
- Min samples leaf with minimum values of (1,5,10,50,100,200,500) would help us try GridSearchCV with varying samples. Since the Bitcoin dataset is imbalanced it makes more sense to *use lower* values to avoid overfitting.

Scoring

When the dataset is imbalanced Accuracy is misleading. Hence It makes more sense to score the algorithm on Precision due to the the imbalanced bitcoin dataset. Accordingly, Algorithm is tuned and scored based on precision score.

Below are the optimal hyper parameter values for the refined algorithm.

| | |
|--------------|-----|
| n_estimators | 100 |
| max_depth | 150 |

| | |
|------------------|------|
| max_features | auto |
| min_samples_leaf | 5 |

IV. Results

Model Evaluation and Validation

Below are the results of the final model run with random state 100.

| Metric - Final Model | Value |
|----------------------|-------|
| Accuracy | 0.91 |
| Precision | 0.92 |
| FScore | 0.91 |

The final model seems to be reasonable with a very good accuracy, precision and FScore.

I ran the final model on multiple random states(1,20,100 and 500) to validate the final model and below are the validation results

| Metric | Value |
|----------------|-------|
| Mean Accuracy | 0.95 |
| Mean Precision | 0.96 |
| Mean FScore | 0.95 |

As expected, mean validation results came close to the final validation confirming the robustness of the final model.

Justification

Below are the results of the optimized model compared with unoptimized and Benchmark.

| | Benchmark | Unoptimized | Optimized |
|-----------|-----------|-------------|-----------|
| Accuracy | 0.4253 | 0.56 | 0.95 |
| Precision | 0.4253 | 0.54 | 0.96 |
| Fscore | 0.4806 | 0.54 | 0.95 |

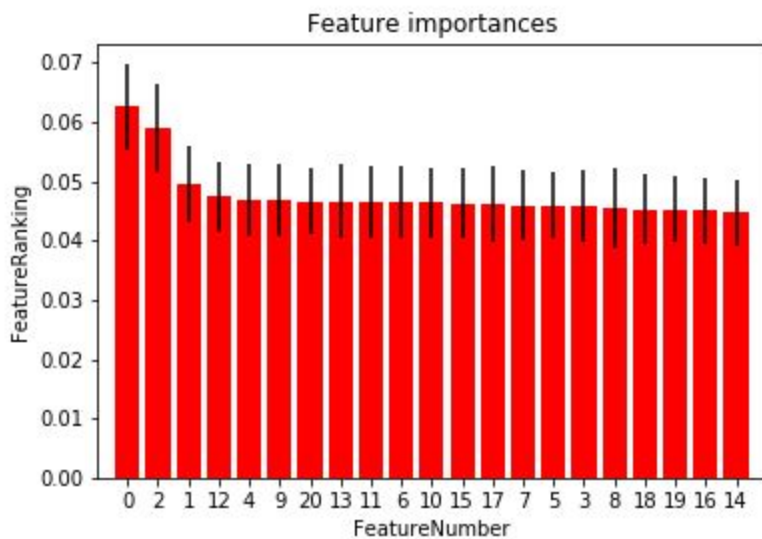
Clearly unoptimized model performed better than benchmark and model performance has significantly improved after optimizing in all the metrics the model is evaluated.

The optimized solution is satisfactory and has performed to the expectations. The optimized model is built by tuning the parameters as discussed in previous section and is significant enough to solve the problem.

V. Conclusion

Free form Visualization

Below graph on feature importance depicts the importance of features in the dataset.



X axis represents feature number and Y axis represent feature ranking. I did not put the feature names on x axis to keep the graph clean. Please find below the mapping of feature number, feature names and ranking for reference. I sorted the below table on feature ranking.

| Feature Number | Feature Name | Feature Ranking |
|----------------|--------------------------------------|-----------------|
| 0 | btc_total_bitcoins | 0.062526 |
| 2 | btc_trade_volume | 0.05903 |
| 1 | btc_market_cap | 0.0495 |
| 12 | btc_cost_per_transaction | 0.047326 |
| 4 | btc_avg_block_size | 0.046775 |
| 9 | btc_difficulty | 0.046676 |
| 20 | btc_estimated_transaction_volume_usd | 0.046612 |
| 13 | btc_n_unique_addresses | 0.046551 |
| 11 | btc_transaction_fees | 0.046457 |

| | | |
|----|---|----------|
| 6 | btc_n_transactions_per_block | 0.046446 |
| 10 | btc_miners_revenue | 0.046392 |
| 15 | btc_n_transactions_total | 0.046242 |
| 17 | btc_n_transactions_excluding_chains_longer_than_100 | 0.046207 |
| 7 | btc_median_confirmation_time | 0.045877 |
| 5 | btc_n_orphaned_blocks | 0.045862 |
| 3 | btc_blocks_size | 0.045779 |
| 8 | btc_hash_rate | 0.045441 |
| 18 | btc_output_volume | 0.045258 |
| 19 | btc_estimated_transaction_volume | 0.045248 |
| 16 | btc_n_transactions_excluding_popular | 0.045058 |
| 14 | btc_n_transactions | 0.044735 |

Reflection

Predicting Bitcoin price is not only interesting but also challenging. Challenging because the problem statement needs lot of research to understand the variables in the bitcoin price trend. As any real time problems Bitcoin dataset is not clean and the dataset is cleaned by replacing empty data with mean of the features and by formatting/sorting the dataset on the dates.

After Data cleansing the data needed to be transformed as most of the features are continuous and are highly skewed. In addition there are varying units among the feature which needed us to scale all the features.

Once data is cleansed, transformed and normalized we have selected the features that matter most by training the data on random forest algorithm and choosing the features that provided 80% of the information gain. The features are later trained on SVC and Random forest ensemble algorithms on a time series split.

Unoptimized model performed comparatively better than benchmark model and the model has been tuned well with accuracy and precision over 0.9. Final model and solution has fit well to my expectations and can be used in general setting to solve any cryptocurrency trend estimation problems.

Improvement

I see further scope for improvement in the solution in the following areas.

1. Although bitcoin prices are largely affected by demand and supply continuum, prices can also be predicted from the sentiment analysis. The algorithm can further be improved by applying sentiment analysis to predict price movement.
2. Since Bitcoin is relatively new the data available is limited and very small. To have a larger dataset algorithm can further be improved to predict on lower frequency time series (5 mins, 10 mins, 30 mins) instead of a day.

I believe improving the algorithm in above two areas might further improve the performance.