

School of Management, Economics, Law, Social Schiences and International Affairs

The effect of Twitter activity on Bitcoin price

Documentation

Software Engineering for Economists (7,610,1.00)

Dimitrios Koumnakes - 10-613-370 Severin Kranz - 13-606-355 Joël Sonderegger - 11-495-488 Alen Stepic - 11-475-258 Chi Xu - XX-XXX

Fall Term 2017

Supervisor
Prof. Dr. Philipp ZAHN
Department of Economics

December 27, 2017

Abstract

(Insert text)

Contents

1	Introduction 1			
	1.1	Goal o	of the paper	1
	1.2	Metho	dology	1
	1.3	Scope		1
2 Data Collection				2
	2.1	Tweets	s Data	2
		2.1.1	Python Script	2
		2.1.2	Hardware Setup	4
	2.2	Bitcoir	n Price Data	4
		2.2.1	Execution	5
		2.2.2	Output	5
		2.2.3	API: Bitcoinaverage.com	5
3	Data Aggregation and Data Wrangling			
		3.0.1	Execution	6
		3.0.2	Output	7
4	Dat	a Anal	lysis	7

List of Figures

List of Abbreviation

API Application Programming Interface

CKEY Consumer Key

CSECRET Consumer Secret

ATOKEN Access Token

ASECRET Access Token Secret

JSON Java Script Object Notation

1 Introduction

In the academic environment accountability and reproducibility is important. However, the publishing process of papers and journals seem to be outdated. New ways of data collection and data processing exist by using computational economics. The usage of algorithms can increase effectiveness and efficiency. Hence, much lager data sets can be proceeded. However this creates also new problems regarding to traceability and reproducibility. Often cited academic paper exist, where the initial computation is not reproducible. Furthermore, some academic paper even contain computational errors. Replicating data or existing results do not provide any new knowledge at all. Nevertheless, the ability to reproduce increases trustworthiness and indicate the quality of the conducted work. This explains why reproduction is of great relevance.

1.1 Goal of the paper

The goal of this documentation is the provision of a description. This description should enable the reader to reproduce the results discussed in the separate paper. Thus, it contains an explanation how the input data have been gathered, stored, aggregated and analysed. In other words, the input data, the model core, the model parameters and the applied math program are explained.

1.2 Methodology

This documentation consists out of four chapters. The first chapter contains a short introduction and provides the reader with an overview about the topic. Furthermore it points out the relevance of documentation. The second chapter discusses the input data. This includes the process of gathering and storing twitter tweets as well as the gathering of the bitcoin price data. The third chapter discusses how the data is aggregated by pointing out the core model and its parameters. Finally the fifth chapter discusses how the analysis has been conducted.

1.3 Scope

The scope of the documentation is the provision of an overview about the different steps which have been conducted to obtain the results in the paper. It does not contain any discussions about the results of the separate paper. It is not a deep description of the code as the code itself as the code is documented separately. Nevertheless, important lines of code are discussed.

2 Data Collection

This section contains a detailed description of how the data for the sequential analysis is gathered and stored. This includes two subsections the (1) tweets data and the (2) bitcoin price data.

2.1 Tweets Data

Whit the python script real-time twitter data are streamed and stored. This happens with help of a raspberry pi.

2.1.1 Python Script

Twitter offers different Application Programming Interfaces (API) for collecting data. However, the time frame for gathering data on a freely base is limited to 7 days. On the other hand, python offers different twitter libraries. Such as the open-source package tweepy. This package has been used for streaming the twitter data as it simplify the script.

Installing Tweepy

Tweepy is installed very simply by running following commands in the command prompt.

pip install tweepy

If the previous downloaded python installation package does not contain the tweepy library, the tweepy package has to be downloaded. The package can be downloaded for free from the following link:

https://pypi.python.org/pypi/tweepy

Twitter Authentication

To access the twitter data, twitter requests an identification of the user. The identification is assured by different keys and access tokens. Those are the (1) consumer key (ckey), (2) consumer secret (csecret), (3) access token (atoken) and (4) access token secret (asecret).

To obtain the mentioned keys and tokens a twitter account is needed. Once, a twitter account exist a application has to be created. This has to be conducted by login with the twitter account credentials under the following link:

https://apps.twitter.com/

After the creation of a application, the keys and tokens can be extracted. Figure 1 illustrates how to retrieve the keys and tokens.

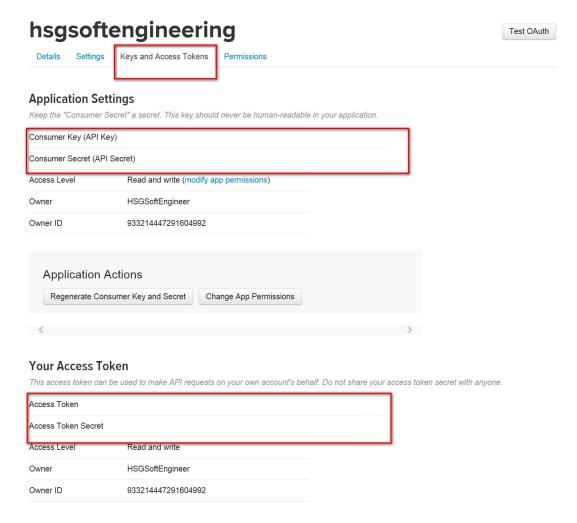


Figure 1: Twitter Keys and Access Tokens

Twitter Streaming API

By running the python script collectTwitterData.py real-time twitter data is pushed in a Java Script Object Notation (JSON) format. Tweets are just pushed in case a tweet contains the defined key word bitcoin.

\$ python collectTwitterData.py

From the JSON format the following parameters are decoded:

- created at: Timestamp of the created tweet
- text: text of the tweet

The time timestamps are in UTC time.

After a successfully decoding the parameters the data is appended to the twitterData.csv file in a new row and saved in the folder /data.

An excerpt of an example response looks like the following:

```
Г
   {
       "created_at": Tue Dec 19 20:40:48 +0000 2017,
       "text": "RT @WhiteBitcoin: New Bitcoin White. Effective, Flexible,
           Reliable",
   },
   {
       "created_at": Tue Dec 19 20:40:49 +0000 2017 ,
       "text": "As #bitcoin prices surge, #xrp remains affordable and FAST!",
   },
   {
       "created_at": Tue Dec 19 20:40:49 +0000 2017,
       "text": "Singapore issues bitcoin warning after price rise - #BTC #Bitcoin
           #Crypto",
   }
]
```

Furthermore to ensure a high quality data set an exception has been integrated into the collectTwitterData.py script. If an error occurs an email will be sent with the error message. Furthermore the logfile will be appended to /data/tweet_collection_error_log.csv.

2.1.2 Hardware Setup

As explained in the previous section the twitter API was used to filter tweets in realtime based on the buzzword "bitcoin". The python program was running for seven days without interruption (to be checked!!!!!!!!!!) on a raspberry pi (rpi) which was accessed with a secure shell (SSH) connection. This approach allowed us to ensure a continuous data stream while being able to access the program output regardless anywhere and anytime. Furthermore due to a low energy consumption of the rpi the method is also the most cost efficient one compared to a personal computer or NAS Server. The following section describes the different steps that have been taken in order to setup the rpi, ssh, git and python3.

The following steps were conducted by using a Raspberry Pi 2, a micro SD card and Python3.

(1) Setup of the Raspberry Pi

- Burn image 2017-11-29-raspbian-stretch.zip to micro SD card available from https://www.raspberrypi.org/downloads/raspbian/
- Boot image from micro SD card

- Open SSH session using PuTTY SSH client for Windows available from https://www.chiark.greenend.org.uk/s̃gtatham/putty/latest.html
- Login to server with credentials

(3) Update system to the newest software

- Update list of applications with the command: sudo apt-get update
- Update applications with the command: sudo apt-get upgrade

(4) Setup Git

- Install git using the command: sudo apt-get install git
- Clone git repository from https://github.com/joelsonderegger/twitterbitcoin.git

(5) Setup Pyhton3

- Install Python3 using the command: sudo apt-get install Python3
- Install the tweepy module for python using the command: pip install tweepy

(6) Collect Twitter Data

- Open a new screen using the command: screen -S 'name'
- start python program using the command: python3 CollectTwitterData.py
- Detach the screen to keep the program running after secure shell connection has been terminated using the commands: Ctrl+A, Ctrl+D
- Re-attach the screen to check whether the script is still running using the command: screen -r 'name'

(7) Download the collected data

- Install FileZilla available https://filezilla-project.org/
- Connect to rpi and navigate to data location
- Download data manually

In summary, this approach worked quite well and run smoothly after all bugs have been resolved from the file CollectTwitterData.py A clean setup is crucial to ensure continuous data collection over a longer period.

2.2 Bitcoin Price Data

We wrote a Python script which collects Bitcoin price data as there was no preexisting data set that satisfied our needs. The Bitcoin price is best expressed by the Bitcoin Price Index. The Bitcoin price index (BPI) is an index of the exchange rate between the Bitcoin (BTC) and the US dollar (USD) (**kristoufek2015main**). The objective of the script was to gather hourly Bitcoin Price Index data for at least the time period in which we gather the tweets data. We found the an API by bitcoinaverage.com which sufficed our needs. An API description follows later.

2.2.1 Execution

By executing the python script collectCryptocurrencyData.py hourly data for the Bitcoin Price Index is retrieved.

\$ python collectCryptocurrencyData.py

2.2.2 Output

After successfully running the python script CollectCryptocurrencyData.py the file bpi.csv is generated in the folder /data. It is important to note that every execution of the script overwrites any existing bpi.csv file.

The file bpi.csv contains historical Bitcoin Price Index data for one month on an hourly basis. Each data point consists of the following parameters:

• time: Timestamp on an hourly basis in UTC time

• average: Average price (in USD)

• high: Highest price (in USD)

• low: Lowest price (in USD)

• open: Opening price (in USD)

2.2.3 API: Bitcoinaverage.com

Bitcoinaverage.com offers a free API that provides real-time and historical price data for a range of crypto-currencies including Bitcoin. The following requests delivers data for an per hour monthly sliding window.

Request

The request to get the data for an per hour monthly sliding window looks as follows. This request require authentication that requires registration and the generation of an API key. The registration and generation of an API key is freely available on bitcoinaverage.com. The collectCryptocurrencyData.py already contains the necessary keys. This means that you need no register or generate keys to execute the script collectCryptocurrencyData.py.

https://apiv2.bitcoinaverage.com/indices/global/history/BTCUSD?period=monthly&?format=json

Response An excerpt of an example response looks like the following:

```
Г
   {
       "high": 8271.04,
       "average": 8247.83,
       "open": 8242.39,
       "low": 8217.72,
       "time": "2017-11-22 15:00:00"
   },
   {
       "high": 8246.82,
       "average": 8203.19,
       "open": 8203.81,
       "low": 8157.25,
       "time": "2017-11-22 14:00:00"
   },
   {
       "high": 8267.27,
       "average": 8238.62,
       "open": 8248.77,
       "low": 8198.54,
       "time": "2017-11-22 13:00:00"
   }
]
```

3 Data Aggregation and Data Wrangling

Before the collected data sets can be analyzed they need to be aggregated. In addition, some data wrangling is necessary to bring the data in a format which than can be analyzed.

3.0.1 Execution

By executing the python script aggregateTwitterBpi.py a aggregated data set that contains number of tweets and Bitcoin Price Index, both on an hourly basis, gets generated.

\$ python aggregateTwitterBpi.py

3.0.2 Output

After successfully running the python script aggregateTwitterBpi.py the file nr_of_tweets_bpi_closing is generated in the folder /data. It is important to note that every execution of the script overwrites any existing nr_of_tweets_bpi_closing_price.csv file.

The file aggregateTwitterBpi.csv contains historical tweets and Bitcoin Price Index data. Each data point consists of the following parameters:

- time: Timestamp on an hourly basis in UTC time
- nr_of_tweets:
- df_nr_of_tweets:
- log_df_nr_of_tweets:
- df_bpi_closing_price:
- log_df_bpi_closing_price:

4 Data Analysis

(Dimitri's Part)