

به نام خدا

استاد درس: دکتر اکبری دانشگاه صنعتی امیرکبیر دانشکدهی ریاضی و علوم کامپیوتر مباحثی در علوم کامپیوتر گزارش پروژه ۷ بهمن ۹۹

پیشبینی قیمت سهام از شبکههای اجتماعی

چکیده

در این مقاله، یک روش برای پیشبینی روند کلی ارزش بیتکوین ارائه شده و سپس خروجی آن، روی یک اپلیکیشن موبایل پیاده شده است. در حال حاضر، بیتکوین بزرگترین رمزارز ٔ از لحاظ سهم از بازار است. با این حال، قیمت آن در یک بازهی روزانه یا حتی بلند مدت، بسیار پرنوسان است.

توییتر یک شبکهی مجازی بسیار محبوب است که گاهی به عنوان یک منبع خبری تاثیرگذار برای تصمیم گیری در خرید ارزهای دیجیتال مورد استفاده قرار می گیرد. تحلیل توییتها و فهمیدن نظر عموم مردم دربارهی روند قیمت بیت کوین، می تواند کمک شایانی به تصمیم گیری در مورد مبادلهی آن باشد.

پس از جمع آوری توییتهای مرتبط با بیت کوین و قیمت بیت کوین در بازه ی زمانی توییتها، این نتیجه حاصل شد که به دلیل تعدد متغیرها، استفاده از رگرسیون آراه درستی نیست و درنهایت با استفاده از مدلهای مدلهای Classification متفاوت و مقایسه ی خروجی هر مدل با ارزیابی، به نتیجه گیریهایی که به آن خواهیم پرداخت، دست یافتیم.

۱ مقدمه

امروزه عده ی زیادی بیت کوین را به عنوان طلای دیجیتال می شناسند. این شناسایی دلیل خوبی هم دارد. این رمزارز حجم معاملات بسیار بالایی داشته و با توجه به تعداد محدود آن و تقاضای روزافزون مردم، به راحتی می توان به تشابه آن با طلا پی برد.

اما با تمام اینها، نوسان شدید این رمزارز باعث شده تا پیشبینی زمان مناسب برای خرید یا فروش آن، بسیار دشوار شود.

میان فاکتورهای تاثیرگذار روی بیت کوین، رسانه تاثیر به سزایی روی قیمت آن می گذارد. به کمک شبکههای مجازی، هر خبر می تواند به سرعت پخش شده و قیمت بیت کوین را تحت تاثیر قرار دهد. همچنین بسیاری از پیش بینی کنندههای بیت کوین، نظر عمومی حال حاضر مردم را در پیش بینی خود، در نظر نمی گیرند. بنابراین ما بر آن شدیم تا روشی برای پیش بینی روند کلی ارزش بیت کوین ارائه کنیم.

Bitcoin 1

Crypto Currency ²

Regression 3

به این منظور، ابتدا توییتهای مربوط به بیت کوین به همراه قیمت بیت کوین در زمان توییتهای جمع آوری شده، به روشی که در بخش مجموعه دادهها به تفصیل بیان شده است، جمع آوری و اطلاعات لازم از آن استخراج و پیش پردازش شده است.

همچنین برای ارتباط میان احساسات[†] توییت و جهت حرکت قیمت بیت کوین، از مدلهای Classifier استفاده شده که اطلاعات بیشتر در این باره، در بخش مدل قابل مشاهده است.

۲ مجموعه داده

۲.۱ توستها

توییتها در دو دسته جمع آوری شدند؛ توییتهای قدیمی و جدید.

توییتهای قدیمی مربوط به سال ۲۰۱۷ هستند که از مجموعه دادهی موجود در لینک //thtps:// هستند که از مجموعه دادهی موجود در لینک //thub.com/teoYQ/Bitcoin-Twitter-sentiment-analysis/tree/master و github.com/teoYQ/Bitcoin-Twitter-sentiment-analysis/tree/master و توییتهای جدید، مربوط سال ۲۰۲۱ با استفاده از API توییتر استخراج شدهاند.

معیار انتخاب توییتها در ابتدا، تنها وجود کلمهی "bitcoin" و انواع مخلف نوشتاری آن مانند "BTC" در متن و برچسبهای آن بود. (bitcoin OR Bitcoin OR BitCOIN OR btc OR Btc OR BTC) پس از استخراج تعدادی از توییتها، متوجه شدیم ریتوییتها به عنوان یک توییت جدا به دیتاست وارد شده و موجب بروز اختلال در مدلها می شدند که با اعمال محدودیت ریتوییت نشده بودن، این مشکل برطرف شد. (iang:en)

از جمله اطلاعاتی که از این توییتها استخراج شد نیز، می توان به زمان توییت، تعداد علاقه مندی 9 ، ریتوییت، پاسخ و نقل قول 4 اشاره کرد. (tweet.fields=public_metrics,created_at)

https://api.twitter.com/2/tweets/search/recent?query=(bitcoin OR Bitcoin OR BITCOIN OR btc OR Btc OR BTC) -is:retweet lang:en&tweet.fields=public_metrics,created_at

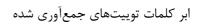
Sentiment 4

Retweet 5

Favorite 6

Reply 7

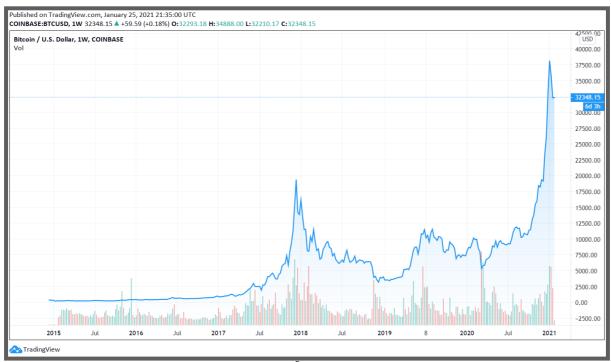
Quote 8





۲.۲ قیمت بیت کوین

برای استخراج قیمت بیت کوین در بازه ی زمانی توییتهای استخراج شده، از سایت https:// سایت Binance به آدرس // api.binance.com/api/v3/ticker/price?symbol=BTCUSD



نمودار قیمت بیت کوین به دلار آمریکا از سال ۲۰۱۵ تا الان

۲.۳ پیش پردازش

در ابتدا برای پیشپردازش توییتها، لینکها، برچسبها و تمام کلمات غیر انگلیسی حذف شده تا تنها حروف انگلیسی باقی بماند. همچنین کلمات توقف^۹ نیز از توییتها حذف شده و تمام کلمات با ریشهی آنها جایگزین شد.

همچنین توییتهای دارای کلمات advertisement ،ad ،lucky و block chain تبلیغ و هرزهنامه ۱۰ شناخته شده و از مجموعه دادهها حذف شدند.

Stop Words 9

Spam 10

۳ مدل

در ابتدا تلاش شد تا مسئله، به روش رگرسیون حل شود. به همین منظور، یک رگرسیون چندجملهای ۱۱ نوشته شد اما جوابهای قابل قبولی از آن به دست نیامد و قابل استفاده نبود.

سپس تلاش شد تا مسئله به صورت طبقهبندی ۱۲ حل شود. به همین منظور از classifierهای زیر با نتایج قابل مشاهده، استفاده شد.

KNN مدل ۳.۱

برای تنظیم ابرپارامتر k، با استفاده از روش بازو۱۳، عدد بهینه، ۴ به دست آمد.

```
[n [12]: knn ()
1 [[ 819 1011]
[ 613 1257]]
2 [[ 819 1011]
[ 613 1257]]
3 [[ 819 1011]
[ 613 1257]]
4 [[ 819 1011]
[ 613 1257]]
5 [[ 819 1011]
[ 613 1257]]
5 [[ 819 1011]
[ 613 1257]]
7 [[ 819 1011]
[ 613 1257]]
8 [[ 819 1011]
[ 613 1257]]
9 [[ 819 1011]
[ 613 1257]]
10 [[ 819 1011]
[ 613 1257]]
11 [[ 819 1011]
[ 613 1257]]
11 [[ 819 1011]
[ 613 1257]]
12 [[ 819 1011]
[ 613 1257]]
12 [[ 819 1011]
[ 613 1257]]
13 [[ 819 1011]
[ 613 1257]]
14 [[ 819 1011]
[ 613 1257]]
15 [[ 819 1011]
[ 613 1257]]
16 [[ 819 1011]
[ 613 1257]]
17 [[ 819 1011]
[ 613 1257]]
18 [[ 819 1011]
[ 613 1257]]
```

ماتریس سردرگمی به ازای kهای متفاوت

نتایج به دست آمده:

Accuracy: 56.1% F1: 0.49

Precision: 0.57 Recall: 0.44

۳.۲ مدل Naïve Bayes

با استفاده از روش گاوسی نتایج زیر به دست آمد:

Accuracy: 52.4% F1: 0.60

Precision: 0.75 Recall: 0.51

Polynomial Regression 11

Classification 12

Elbow 13

۳.۳ مدل Random Forest

در این روش، ابرپارامتر، یعنی تعداد درختهای مورد استفاده، برای اعداد ۱ تا ۱۰۰ آزموده و در نهایت مقدار بهینه ۱۱ به دست آمد و نتایج زیر از آن کسب شد:

Accuracy: 59% F1: 0.66

Precision: 0.78 Recall: 0.56

۳.۴ مدل Decision Tree

نتایج زیر به دست آمد:

Accuracy: 59.5% F1: 0.65

Precision: 0.8 Recall: 0.56

۴ خروجی نهایی

در نهایت مدل به دست آمده روی یک سرور قرار گرفت و در یک اپلیکیشن موبایل، با وارد کردن متن یک توییت، با استفاده از API ساخته شده، نتیجهی پیشبینی مدل، مبنی بر بالا یا پایینرونده بودن قیمت بیت کوین، به کاربر نشان داده می شود.

همچنین این اپلیکیشن، از قابلیتهای دیگری همچون اعلام قیمت لحظهای بیتکوین، نمایش نمودار قیمت بیتکوین و آخرین توییتهای مرتبط با بیتکوین نیز برخوردار است.

در آینده می توان، قابلیت پیش بینی قیمت با توجه به آخرین توییتهای موجود در توییتر را نیز به این اپلیکیشن اضافه کرد.

۵ نتیجه

یکی از چالشهای اصلی، جمعآوری مجموعه دادهی مناسب بود که با توجه به آنچه در بخش مجموعه داده توضیح داده شد، مجموعه دادهی مناسب برای آموزش مدلها، جمعآوری و سپس پردازش شد.

برای جمعآوری توییتها از API توییتر، key نیاز بود که برای به دست آوردن آن، یک فرم پر شد که در آن مواردی همچون هدف استفاده از این API و موارد دیگر ذکر شد که پس از جند روز، این دسترسی داده شد.

به منظور بهبود یادگیری شبکه، بازههای زمانیای که بیتکوین هم دارای روند نزولی و هم دارای روند صعودی باشد نیاز بود که با تحلیل بازههای زمانی متفاوت، یک بازهی مناسب یافت شد و مجموعه دادهی همان بازه، مورد استفاده قرار گرفت.

همچنین در این پروژه، زمان زیادی صرف آزمایش مدلهای مختلف شدو با این حال، در بهترین مدل، دقت ۶۰٪ به دست آمد که با توجه به زمان مجدود بهینهسازی مدلها، بهتر از چیزی بود که گمان میرفت.

۴ اعضای گروه

٥ اميرحسين امامي ٩٧١٢٠٠٢

۰ سارا بابایی ۹۷۱۳۰۰۵

٥ ارشيا مجيدي ٩٧١٣٠٤١