

تحليل بازار رمزارزها

گزارشی از روند انجام پروژه بوت کمپ علوم داده کوئرا



Quera
GROUP 8

تحلیگران پروژه

یوسف عبدالکریمی

علیرضا محمدزاده
حسین زندی نژاد
حسین نادری
مهرداد عادلی
مهرشاد فلاح

منتور و راهنما

طراحی پایگاه داده
تحلیل‌های آماری
مصورسازی و طراحی گزارش
مستندسازی و اسکریپینگ
استخراج داده‌ها و اطلاعات رمزارزها

Quera



بخش اول: استخراج داده

((ابتدا با استفاده از لینک‌ها و بهره‌گیری از **سِلنیوم^۱ و درایور فایرفاکس**، اطلاعات مورد نظر مانند اسم رمزارز، سمبل، لینک صفحه اصلی، لینک تاریخچه هر رمزارز، ارزش بازار، حجم مبادله ۲۴ ساعت اخیر و حجم در دسترس مشتریان را بدست آورده و این تعداد از اطلاعات در قالب چند لیست ذخیره می‌شوند.))

Scraping

در این بخش، با استفاده از Selenium، اطلاعات مربوط به ۲۰۰ ارز دیجیتال منتخب از وبسایت coinmarketcap.com در تاریخ خواسته شده استخراج می‌شود. همچنین در این مرحله، لینک‌های مرتبط با صفحات GitHub و تگ‌های مربوط به هر رمزارز نیز استخراج می‌شوند.

Pandas

با اطلاعات استخراج شده از مرحله قبلی، یک DataFrame پانداس ایجاد می‌شود برای مدیریت داده‌ها که در آن اطلاعات هر رمزارز در ردیف‌های مختلف قرار دارد. این DataFrame نیازی به مدیریت داده‌های خالی ندارد زیرا هر ردیف یکتا و منحصر به فرد می‌باشد.

Output

اطلاعات استخراج شده در DataFrame به یک فایل CSV ذخیره می‌شوند تا برای تحلیل‌های بعدی قابل استفاده باشند. همچنین، داده‌های تاریخ برای هر ارز از صفحات دانلود شده و به صورت فایل‌های CSV در دایرکتوری مشخص شده ذخیره می‌گردد.

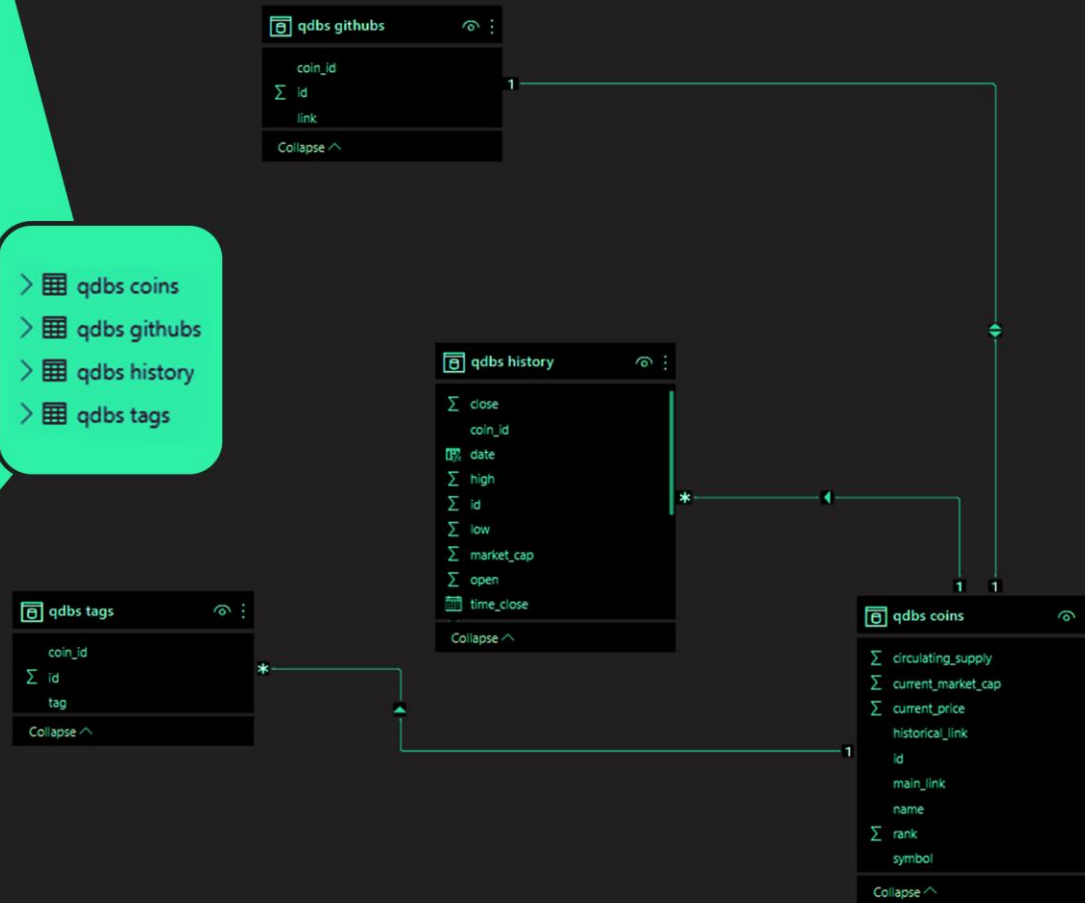
۱. "Selenium" یک ابزار اتوماسیون مرورگر است که برای تعامل با وبسایت‌ها و استخراج داده‌ها به کار می‌رود. این ابزار می‌تواند مرورگر وب را کنترل کرده و عملیاتی مانند کلیک کردن بر روی المان‌ها و تعامل با صفحات وب را انجام دهد، ممکن می‌سازد.

Database Design



بخش دوم: طراحی دیتابیس

ابتدا دیتاست‌ها^۱ رو برای وارد کردن در دیتابیس آماده می‌کنیم. این کار شامل وارد کردن و Import فایل‌های ساخته شده توسط اسکریپت، ترکیب سطرهای دیتاست‌های تاریخچه، تغییر نام ستون‌ها، ساخت دیتاست تگ‌ها و گیت‌هاب می‌شود.



سپس برای موجودیت‌های موجود که شامل **Coin, Tags, History, GitHub** می‌شوند، schema جدول‌های مربوطه را در دیتابیس می‌سازیم. این کار شامل ساخت ستون‌ها و تعیین نوع داده، ساخت کلیدها و ایجاد رابطه میان جداول می‌شود.

۱. "دیتانیت" مجموعه‌ای از داده‌ها است که معمولاً در جداول یا ماتریس‌ها ذخیره می‌شوند، هر سطر در دیتانیت به یک نمونه مربوط می‌شود و به تحلیل و استفاده در مطالعات و تحقیقات مختلف کمک می‌کند.

Statistical Analysis

بخش سوم: تحلیل‌های آماری

آمار توصیفی

در ابتدا با استفاده از sqlalchemy جداول به دست آمده در بخش قبلی را با کتابخانه Pandas به دیتافریم تبدیل کردیم. سپس با بررسی های اجمالی از دیتاست ها به ادامه کار پرداختیم. در این بخش، ابتدا نمودارهای پراکندگی ارزش بازار در مقابل حجم معاملاتی روزانه، نمودار توزیع حجم معاملات روزانه رمزارزهای قابل استخراج و نمودار میله ای ۱۰ رمزارز برتر از نظر تعداد افزایش قیمت در روزهای قمری بهار سال ۲۰۲۳ با استفاده از کتابخانه Plotly رسم شده است. سپس به بررسی رمزارزهایی که در سال اخیر تغییرات قیمت همسو داشته‌اند پرداخته شد. در این بخش، رمزارزهایی با ویژگی خواسته شده به همراه تعداد روزها با تغییرات قیمت همسو در جدولی نمایش داده شده است. در آخر نیز ماتریس همبستگی و نقشه حرارتی (heatmap) تغییرات قیمت برای شانزده رمزارز برتر از نظر ارزش بازار رسم شده است. این ماتریس کمک می‌کند تا روند تغییرات قیمت این رمزارزها در یک سال اخیر را در مقابل یکدیگر رصد کنیم.

آزمون فرض

در این بخش ابتدا با استفاده از روش‌های نمونه گیری، چهل رمزارز از میان داده های استخراج شده به صورت تصادفی انتخاب شده است. سپس میانگین حجم معاملاتی روزانه هر نمونه محاسبه شد و در نهایت بازه اطمینان ۹۸ درصدی برای حجم معاملاتی به دست آورده شد. ابتدا به بررسی روزهای متلاطم بازار برای پاسخ به سوال اول پرداخته شد. برای به دست آوردن شاخصی مناسب برای این آزمون کردن فرض خواسته شده، دو رویکرد ارائه شد. شاخص اول، تغییرات قیمت بر اساس تفاضل قیمت پایانی (close) و قیمت آغازین (open) و شاخص دوم به صورت تغییرات قیمت بر اساس تفاضل بالاترین قیمت (high) و پایین ترین قیمت (low)، در هر روز در نظر گرفته شدند. در هر دو حالت، فرض صفر که بیان کننده عدم وجود تفاوتی فاحش میان دو انتخاب ممکن برای روزهای کاری می باشد، با آزمون فرض t رد شد. سپس درستی ادعای "میانگین حجم معاملات روزانه، Bitcoin، Ethereum و USDt Tether به شدت بیشتر از میانگین حجم معاملات روزانه باقی رمزارزهاست" بررسی شد. این ادعا با آزمون فرض t رد شد. در آخر، اعتبار ادعای "ترمال بودن توزیع تغییرات قیمت رمزارزها در روزهای تعطیل (شنبه و یکشنبه)" با آزمون های فرض Anderson و Sminrov-Kolmogorov این ادعا رد شد.



Data Visualisation & Presentation

بخش چهارم: Power BI

مصورسازی

با استفاده از دیتابیس ساخته شده در مراحل قبل و استفاده از ابزار مصورسازی Power Bi ابتدا دیتاها را به منظور تحلیل بهتر به ستون‌های جزئی‌تر تقسیم کرده و در نهایت با استفاده از ویژوال‌های مختلف، چارت و نمودارهای موردنیاز برای ارائه این پروژه به مخاطب هدف را طراحی و اجرا می‌کنیم.

گزارش و ارائه

با بهره‌گیری از ابزارهای طراحی مثل Adobe Photoshop و PowerPoint یک گزارش حرفه‌ای از روند پیشرفت و انجام تسک‌های داده شده به همه اعضا گروه تحلیلگران به مخاطبین داده می‌شود. این گزارش باید دارای خلاصه‌ای از مراحل انجام هر بخش از پروژه باشد و نمای کلی از تاریخچه تکمیل وظایف بدهد.

REPORT SERVER

یک سرور برای نمایش داشبورد و نمودارهای طراحی شده در قدم قبلی تهیه و تنظیم می‌گردد و قابلیت نمایش اطلاعات در انواع دستگاه‌ها امکان پذیر می‌کند.

More Information...

بخش پنجم: اطلاعات بیشتر



روند کار اسکریپینگ گیت‌هاب

با استفاده از لینک‌های گیت‌هاب که به وسیله‌ی اسکریپر اصلی استخراج شده‌اند به گیت‌هاب هر کوین ریکوئست زده شده و اگر آدرس لینک متعلق به صفحه اول کوین نباشد، به وسیله توابع سلیوم به صفحه‌ی اول منتقل شده و از بخش people تعداد کل مشارکت‌کنندگان و لینک پروفایل مشارکت‌کنندگانی که در صفحه‌ی اول این بخش قرار دارند استخراج میشود. سپس همین روند تکرار شده و از صفحه‌ی اصلی هر کوین، زبان‌های برنامه‌نویسی‌ای که بیشترین استفاده را در ساختن آن کوین داشته‌اند استخراج می‌شود. سپس با استفاده از توابع سلیوم و قابلیت مرتب‌سازی و طبقه‌بندی ریپازیتوری‌ها، آنها را بر اساس تعداد استار و اینکه ریپازیتوری سورس باشند مرتب کرده و کل المان‌های ریپازیتوری را استخراج می‌کنیم. سپس از این المان‌های استخراج شده، نام ریپازیتوری، لینک آن، جزئیات آن اگر موجود باشد و همینطور لایسنس، تعداد فورک‌ها و استارهای ریپازیتوری را استخراج می‌کنیم. برای هر کوین که صفحه‌ی گیت‌هاب ندارد رشته‌ی «No Github Link» ذخیره می‌شود. همچنین داده‌ها بر اساس رتبه‌ی کوین شماره‌دهی می‌شوند. در آخر داده‌ها در سه فایل csv شامل مشارکت‌کنندگان، زبان‌های برنامه‌نویسی و ریپازیتوری‌ها ذخیره میشوند.

پایان!



Quera
GROUP 8