



## روش کشف تقلب در پایگاه داده گرافی

الهام عبدالهی<sup>۱</sup>، فرشته آزادی پرند\*<sup>۱</sup>، علی اصغر صفائی<sup>۲</sup>

گروه رایانه، دانشکده علوم ریاضی رایانه، دانشگاه علامه طباطبائی<sup>۱</sup>

گروه انفورماتیک پزشکی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس، تهران، ایران<sup>۲</sup>

### چکیده

امروزه با توجه به گسترش ارتباطات بایستی حجم عظیمی از اطلاعات تولید، مبادله و در نهایت ذخیره‌سازی شود. سیستم‌های پایگاه داده‌ی رابطه‌ای، قابلیت ذخیره‌سازی این حجم از اطلاعات را ندارند. مفهوم NoSQL برای مواجهه با این مشکل در سال ۱۹۹۸ توسط کارلو استروزی ظهور پیدا کرد. از انواع پایگاه داده‌های NoSQL می‌توان به پایگاه داده‌های سندگرا، کلید مقدار، ستون محور و در نهایت پایگاه‌های گرافی اشاره کرد. در این کار تمرکز ما بر روی پایگاه داده‌های گرافی خواهد بود. پرکاربردترین پایگاه داده‌های گرافی، پایگاه داده neo4j است. پایگاه داده مذکور توسط شرکت‌های معروفی از جمله hp, Microsoft, adobe و نظایر آن استفاده می‌شود. از مهم‌ترین کاربردهای پایگاه داده مذکور می‌توان به مساله کشف تقلب و کلاهبرداری، کار با اینترنت اشیا و سیستم‌های توصیه‌گر، اشاره کرد. در این کار به بررسی تقلب افزایش قیمت سهم در یکی از شرکت‌های وابسته به سازمان بورس پرداخته خواهد شد. روند کار به این شکل است که عوامل دخیل در افزایش قیمت سهم شناسایی شده و داده‌های مربوطه در پایگاه داده گرافی neo4j ذخیره شده است. در نهایت زبان پرس و جوی cypher جهت بازیابی زنجیره‌های تقلب به کار گرفته شده است.

کلمات کلیدی: پایگاه داده غیررابطه‌ای گرافی Neo4j، زبان پرس و جوی Cypher، کشف تقلب، NoSQL



## تاریخچه مقاله:

تاریخ ارسال: ۹۷/۱۱/۱

تاریخ اصلاحات: ۹۸/۳/۱

تاریخ پذیرش: ۹۸/۴/۱

تاریخ انتشار: ۹۸/۵/۱۵

**Fraud Detection Method in NoSQL Graph Database**Elham Abdollahi<sup>1</sup>, Fereshteh-Azadi Parand<sup>\*1</sup>, Ali Asghar Safaei<sup>2</sup><sup>1</sup>Allameh Tabataba'i University, Tehran<sup>2</sup>Department of Medical Informatics, Faculty of Medical Sciences, Tarbiat Modares University, Tehran, Iran**Abstract**

Today, with the expansion of communications, a huge amount of production, exchange and ultimately data storage must be stored. Relational databases do not have the capacity to store this volume of information. The NoSQL concept emerged in 1998 by Carlo Strouzi to deal with this problem. NoSQL database types can include documentary databases, value keys, column-driven and eventually graphic databases. In this work, our focus will be on the graphic database. The most widely used graphic database is the neo4j database. The database is used by famous companies such as Cisco, HP, Microsoft adobe and the like. One of the most important uses of the database is the issue of fraud and fraud detection, work with the internet of things and recommendation systems. This work will investigate the fraud in the share price increase in one of the affiliated companies of the stock exchange. The market price is fluctuating momentarily. But this price change is not always based on the value of the relevant share. There are a number of factors involved in the story that have increased share prices without any added value. This paper examines the factors involved in increasing the share price and preserving relevant data in the neo4j database. Finally, the cypher query language is used to retrieve cheat chains.

**Keywords:***Graph database**Fraud detection**NoSQL database**Cypher query language*



## ۱ - مقدمه

وجود نکته این است که اطلاعات در طبیعت لزوماً به صورت جدولی نیست. فرض کنید که در یک رابطه دوستانه در یک شبکه اجتماعی بخواهیم دوستان دوستان را بازیابی کنیم. چنین پرس و جویی در پایگاه داده‌ی رابطه‌ی حجم عظیمی کد ایجاد کرده (پیچیدگی محاسباتی و فضای دارد) و انسجام را از بین می‌برد. لذا زمان انتقال به یک پایگاه داده‌ی غیررابطه‌ای فرارسیده است. از انواع پایگاه داده‌های غیررابطه‌ای می‌توان به پایگاه داده‌های سندمحور، کلید مقدار، ستون‌گرا و گرافی اشاره کرد که هر یک معماری و ویژگی عملیاتی متنوعی دارند و متناسب با منطق کسب و کار مورد استفاده قرار می‌گیرند.

در پایگاه داده‌های سندگرا اسناد به صورت جفت کلید-مقدار ذخیره می‌شوند. این مقادیر همانند فرمت‌هایی از جمله JSON و XML می‌توانند به صورت لیست یا نقشه باشند. اسناد می‌توانند با شناسه، ذخیره و بازیابی شوند. به عبارت بهتر یک پایگاه داده سندگرا، می‌تواند مشابه پایگاه داده کلیدمقدار عمل کند. در موارد کلی، پایگاه داده‌های سند محور به شاخص‌ها می‌پردازند تا دسترسی به اسناد را براساس هر یک از ویژگی‌های آنها سرعت بخشند. در این صورت با نوشتن شاخص داده، می‌توان به جزئیات داده‌ها دسترسی داشت.

در پایگاه داده‌های کلید-مقدار، مقادیر مبهم ذخیره شده، با کلید بازیابی می‌شوند. هر داده‌ای با درهم‌سازی شناسه، (کلید) ذخیره می‌شود. با فرض این که چند باکت و چند ماشین داشته باشیم، تابع درهم‌ساز طوری ساخته می‌شود که توزیع یکنواخت در باکت در دسترس را فراهم کند. به طوری که هیچ ماشین واحد به یک نقطه کانونی تبدیل نشود. با توجه به کلید، می‌توان از

امروزه اغلب سازمان‌ها درگیر مشکلاتی از قبیل تضاد منافع، اخذی اقتصادی، رشوه خواری، پولشویی، پرداخت‌های ساختگی و مسائلی از این قبیل هستند. این موارد همگی تقلب یا کلاهبرداری محسوب می‌شوند. طبق تعریف انجمن بازرسان رسمی تقلب<sup>۱</sup>، تصرف افراد در منافع سازمان جهت غنی‌سازی خود، تقلب است که زیان‌هایی با هزینه‌های گزاف را برای سازمان متحمل می‌کند. آنچه باعث تقلب می‌شود، بصورت اضلاع مثلث تقلب با عناوین فرصت، فشار و منطق تراشی دسته‌بندی شده است. اما برای بازیابی چنین تقلب‌هایی باید به بررسی اطلاعات ذخیره شده در پایگاه داده سازمان پرداخته شود. پایگاه داده‌ها در دونوع رابطه‌ای و غیر رابطه‌ای [۱] دسته بندی می‌شوند.

امروزه با گسترش ارتباطات، ذخیره‌سازی و بازیابی داده‌ها و اطلاعات در پایگاه داده‌های رابطه‌ای، بدلیل تاثیر عواملی چون حجم زیاد، تنوع، سرعت تغییر بالا، دقت اطلاعات و یک سری عوامل دیگر، امری هزینه بر است. مشکل پایگاه داده‌های رابطه‌ای این است که هنگام پیوست جداول در پرس‌وجوها مجموعه بزرگی از اطلاعات را بدون فیلتر کردن نمایش می‌دهند. در نتیجه سرعت کم و درمقابل پیچیدگی زمانی بالایی دارند. داده‌ها با گذشت زمان تغییر می‌کنند. تغییرات داخلی و خارجی یک سیستم و زمینه‌ای که در آن کار می‌شود، می‌تواند تاثیر قابل توجهی را در سرعت داشته باشد. عامل دیگر که باعث عدم عملکرد بهینه و مناسب پایگاه داده‌های رابطه‌ای شده است، تنوع در ساختار داده‌ها است. تنوع، به درجه نظم یا بی‌نظمی، تراکم یا عدم تراکم، اتصال یا عدم اتصال داده‌ها مربوط می‌شود. با این



های گرافی است. مناسبترین پایگاه داده گرافی جهت ذخیره سازی شبکه های اجتماعی، پایگاه داده neo4j است. یکی از کاربردهای پایگاه داده مذکور، کشف تقلب سازمان ها است [۳]. تقلب های تجاری در مواردی رخ می دهد که کارمند، مدیر، افسر و یا مالک سازمان به نفع خود و ضرر آن سازمان تقلب کند. سه نوع عمده تقلب تجاری عبارتند از: فساد، مصادره دارایی و بیانیه های دروغین. فساد می تواند در تضاد منافع که ریشه در طرح خرید و فروش دارد، رشوه خواری، پاداش غیرقانونی و اخاذی اقتصادی دیده شود. منظور از دارایی، نقدینگی و موجودی انبار و دارایی های دیگر است. مصادره دارایی به معنای سرقت از نقدینگی و پرداخت های ساختگی به روش های مختلف است. در صورتی که نسبت سود به درآمد بیش از حد و یا حتی بسیار کمتر بیان شود، این بیانیه ای دروغین است که کلاس دیگری از تقلب است. مدیران می بایست نتایج عملکرد خود را به مالکان واحد تجاری و سایر تأمین کنندگان وجوه مانند بانک ها گزارش کنند. باین حال، گزارش گیری مالی شامل طیف وسیعی از اعداد حسابداری برای نشان دادن سود، جریان نقدی و وضعیت مالی واحد تجاری است. بنابراین، مدیریت می تواند با استفاده نادرست از حسابداری ساختگی، وضعیت مالی شرکت را بهتر از واقعیت نشان داده و از اثرات آن در کوتاه مدت بهره برد [۴]. در این کار به بررسی تقلب های از نوع بیانیه دروغین پرداخته خواهد شد. تقلب مذکور توسط هر یک از عوامل مدیر، کارمند و مشتری قابل انجام است. این نوع تقلب ها که در اثر مخفی کردن واقعیات و یا حتی ارائه اطلاعات گمراه کننده در گزارش ها طی دوره های مختلف یک شرکت است، باعث ایجاد یک سری منافع

این آدرس برای ذخیره مقدار در باکت مربوط استفاده کرد.

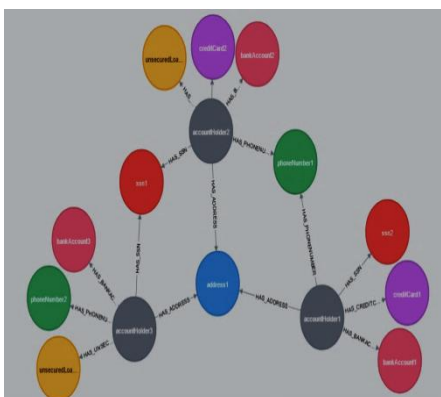
پایگاه داده ای ستون گرا، پایگاه داده ای است که مدل داده ای آن بر پایه جدولی با جمعیت کم است. یک ردیف می تواند شامل هر تعداد ستون دلخواه باشد. ساده ترین واحد ذخیره سازی ستون است که شامل یک جفت نام و مقدار می باشد. ترکیب هر تعداد ستون را ابرستون گویند. هر ابرستون یک نام مختص به خود را دارد. اگر چندین ستون کنار هم در یک ردیف ذخیره شوند، به آن خانواده ستون گویند. اگر چندین ابرستون در یک ردیف قرار گیرند، به آن خانواده ابر ستون گویند.

پایگاه داده های گراف به عنوان معروف ترین ابزار ذخیره سازی و پرس و جو در گراف ها هستند. پایگاه داده های گراف برای حل اطلاعات پیچیده و مشکلات تجاری حائز اهمیت هستند. گراف ها یک مدل طبیعی برای طیف وسیعی از حوزه ها هستند. به عنوان نمونه برای توصیف نحوه عبور کاربر از نقطه A به نقطه B پایگاه داده گرافی راه حل مناسبی است [۲]. انتخاب یک پایگاه داده غیررابطه ای و جایگزینی احتمالی آن با یک پایگاه داده رابطه ای، به دلیل عدم همسویی مناسب با نیازهای موجود یکی از اتفاقات مهمی است که این روزها در بسیاری از سازمان ها در حال وقوع است.

تقلب در سازمان های بیمه و بانک باعث از دست رفتن سالیانه میلیارد ها دلار می شود. صاحبان کسب و کار بایستی راه حل ها و یا سیستم های مناسبی را جهت کشف تقلب در نظر بگیرند. امروزه با توجه به گسترش اطلاعات و ارتباطات، روابط بین سازمان ها از طریق شبکه های اجتماعی صورت می گیرد. برای ذخیره سازی اطلاعاتی از نوع شبکه های اجتماعی نیاز به پایگاه داده



دارندگان حسابی هستند که به اطلاعات افراد دسترسی دارند.



شکل ۱. گراف دسترسی افراد به اطلاعات در سیستم بانکی

در ادامه گام‌هایی برای پرس‌وجوی مناسب جهت یافتن حلقه‌های تقلب، بررسی شده است. در ابتدا دارندگان حسابی که به اطلاعات دسترسی دارند، جست‌وجو می‌شوند. تعداد دارندگان حساب تحت عنوان اندازه حلقه تقلب مد نظر قرار می‌گیرد. دارندگان حسابی که از طریق داشتن `creditCard` و یا بدست آوردن `unsecuredLoan`، حق دسترسی به حساب‌های دیگر را کسب کرده‌اند، جست‌وجو شده و هریک با شناسه منحصر بفرد خود گردآوری می‌شوند. در صورتی که دسترسی از طریق ارتباط `HAS_CREDITCARD` صورت گرفته باشد، مجموع حدپایین حساب و در صورتی که دسترسی از طریق ارتباط `HAS_UNSECUREDLOAN` صورت گرفته باشد، مجموع حد تعادل حساب به عنوان ریسک مالی در نظر گرفته می‌شود.

مدلی برای کشف تقلب سازمان بیمه

یک مدل تقلب در سازمان بیمه حالتی است که متقلبان با صحنه‌سازی حوادث جعلی به سازمان مربوطه مراجعه کرده و ادعای خسارت می‌کنند. مدل

مالی برای متقلبان و زیان‌هایی سنگین برای سرمایه‌گذاران دیگر می‌شود.

از طریق طراحی پرس‌وجوهای مناسب و بررسی دقیق روابط بین موجودیت‌ها، می‌توان الگوی مناسبی را جهت کشف تقلب در نظر گرفت تا آنچه مطابق با این الگو است، از درون پایگاه داده گرافی بازیابی کرد. در ادامه کارهای مرتبط انجام شده، الگوریتم کشف تقلب در افزایش قیمت سهم یک شرکت بورسی بصورت فلوجارت بیان شده است. در نهایت نتایج کار را خواهیم داشت.

## ۲- کارهای مرتبط انجام شده

یکی از رویکردهای کشف تقلب که در پایگاه داده‌های غیررابطه‌ای مانند `neo4j` بررسی می‌شود، استفاده از مدل گرافی است. بدین شکل که با پیمایش گره‌ها و روابط بین آنها و نیز خواص گره‌های مذکور، می‌توان الگوهایی را برای شناسایی حلقه‌های تقلب در نظر گرفت. براساس این الگوها، می‌توان پرس‌وجوهای مناسبی را طراحی کرده و در نهایت حلقه تقلب مذکور را بازیابی کرد. مدل کشف تقلب در سیستم بانکی در [۴] الگویی که مربوط به کشف تقلب سیستم بانکی است، مطرح شده است. مدل‌سازی این مساله بصورتی است که دارندگان حساب بانکی که بصورت غیرمستقیم و از طریق دسترسی به اطلاعات مخفی دیگر دارندگان حساب، با یکدیگر در ارتباط هستند، به عنوان گره‌های متقلب شناسایی و بازیابی می‌شوند. ساختار کلی گراف عملکرد سیستم بانکی مطابق با شکل ۱ است. گره‌های طوسی رنگ

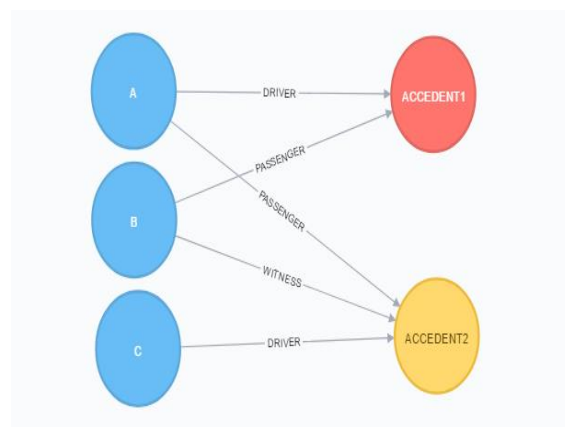


می‌تواند حلقه‌های جنایی را درگیر کند. متقلبین ماهر می‌توانند تعداد بسیار زیادی از هویت‌های ساختگی ایجاد کرده و از این حالت برای پیاده‌سازی طرح‌های بسیار بزرگ استفاده کنند.

معاملات برخط با شناسه‌هایی که در ادامه بررسی شده است، در نظر گرفته می‌شود: شناسه کاربر، آدرس ip، محل جغرافیایی، شماره کارت اعتباری و ردیابی کوکی‌ها. به طور معمول انتظار می‌رود روابط بین این شناسه‌ها بطور نسبی به حالت یک به یکی نزدیک باشد. بعضی از تغییرات به طور طبیعی برای حساب‌های مشترک، دستگاه‌ها و یا استفاده مشترک اعضای خانواده‌ای از یک کارت اعتباری و یا استفاده افراد از چندین کامپیوتر و مانند این موارد قابل تحمل هستند. با این حال، به محض اینکه روابط به سمت عددی فراتر از حالت معقول رفت، اغلب تقلب در حال وقوع است. ممکن است برای نمونه، تعداد زیادی از کاربران معاملاتی ناشی از ip یکسان داشته باشند. در تعداد بسیاری از معاملات، برای آدرس‌های مختلف، از یک کارت اعتباری استفاده شده باشد. تعداد زیادی از کارت‌های اعتباری همگی برای یک آدرس استفاده شده باشد. در هر یک از سناریوها، الگوهای درون گراف، با گذر از روابط بین قطعات مجزا، از طریق اطلاعاتی که می‌توانند تحت عنوان نشانه‌های قوی رویداد تقلب به کار گرفته شوند، قابل بازیابی هستند. عامل بزرگی که به عنوان دغدغه محسوب می‌شود، روابط بیشتر در میان شناسه‌ها است. گراف‌های بهم پیوسته شاخص‌های قدرتمندی برای انجام چنین تقلب‌هایی هستند.

پایگاه داده‌های گراف برای الگویابی در زمان واقعی طراحی شده‌اند. با انجام بررسی‌ها در محل و ارتباط آنها با محرک‌های رویداد، چنین طرح‌هایی قبل

تقلب به حالتی است که افرادی که در یک حادثه هستند، در حادثه دیگر نیز حضور داشته باشند. با این تفاوت که نقش افراد متفاوت باشد. به عنوان نمونه فردی با برچسپ A در حادثه ۱ به عنوان DRIVER و در حادثه ۲ به عنوان PASSENGER باشد. فرد B در حادثه ۱ نقش PASSENGER و در حادثه ۲ تحت عنوان WITNESS عمل کرده است. لازم به ذکر است در حادثه ۲ فرد دیگری (C) تحت عنوان DRIVER وارد عمل شده است. شکل ۲ نحوه عملکرد را نمایش می‌دهد.



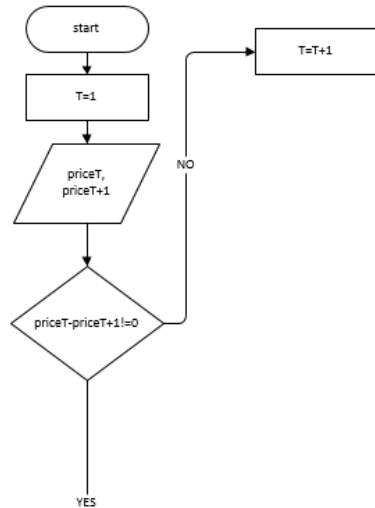
شکل ۲. نمونه ای از تقلب در سازمان بیمه

پرس وجوی مربوطه جهت یافتن حلقه‌های تقلب به این شکل صورت می‌گیرد که گره‌های با برچسپ A, B, C که به هر دو گره با برچسپ ACCIDENT1, 2 متصل شده‌اند، به عنوان گره‌های مقلب شناسایی و بازیابی می‌شوند. [۵]

#### مدل کشف تقلب تجارت الکترونیک

همانطور که زندگی ما بطور فزاینده‌ای به سمت دیجیتالی شدن پیش می‌رود، تعداد معاملات مالی برخط [۶]، رشد می‌یابد. متقلبین سرعت خود را با این روند وفق داده و روش‌های زیرکانه‌ای برای فریب سیستم پرداخت برخط، طراحی کرده‌اند. این فعالیت

می‌گیریم. طبق شکل ۴، ابتدا قیمت سهم در دو لحظه متوالی دریافت می‌شود. در صورت افزایش قیمت سهم وارد مرحله بررسی عامل ارز می‌شویم. در غیر این صورت قیمت سهم در دو لحظه متوالی دیگر دریافت می‌شود.

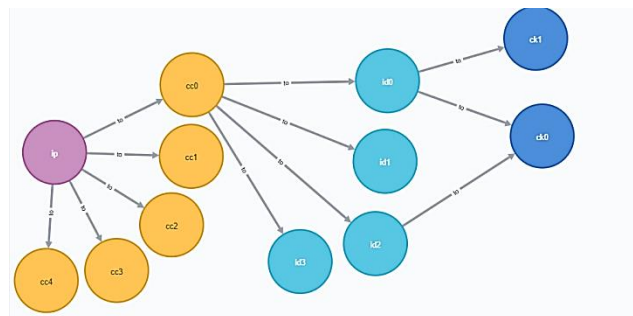


شکل ۴. بخشی از فلوچارت بررسی اختلاف قیمت سهم در دو لحظه متوالی

بایستی نوع شرکت از نظر وارد کننده یا صادر کننده بودن مشخص شود. در صورت افزایش قیمت ارز و صادرکننده بودن، افزایش قیمت سهم طبیعی است. در صورت افزایش قیمت ارز و وارد کننده بودن، افزایش قیمت سهم عامل دیگری دارد. پس باید عامل تصمیمات دولت بررسی شود. در صورت کاهش قیمت ارز و صادرکننده بودن، افزایش قیمت سهم عامل دیگری دارد و باید تصمیمات دولت بررسی شود. در صورت کاهش قیمت ارز و وارد کننده بودن، افزایش قیمت سهم طبیعی است. بخشی از فلوچارت مربوطه در شکل ۵ آورده شده است.

از اینکه بتوانند خسارت قابل توجهی ایجاد کنند، می‌توانند کشف شوند. محرک، می‌تواند رویدادهایی مانند ورود به سیستم، ثبت سفارش یا ثبت نام کارت اعتباری را شامل شود.

گراف شکل ۳، تراکنش آدرس ip با رویداد تقلبی که از ip رخ می‌دهد را نشان می‌دهد.  $cc_x$  نشانگر شماره کارت اعتباری،  $id_x$  نشانگر id کاربری برای انجام تراکنش‌های بر خط و  $ck_x$  به کوکی مشخص ذخیره شده در سیستم اشاره می‌کند. در شکل ۳، ip مربوطه، با استفاده از ۵ کارت اعتباری، چندین معامله را انجام می‌دهد که یکی از آنها ( $cc_0$ )، توسط چندین id استفاده شده است. دو کوکی  $ck_0, ck_1$  بین دو مشترک است.



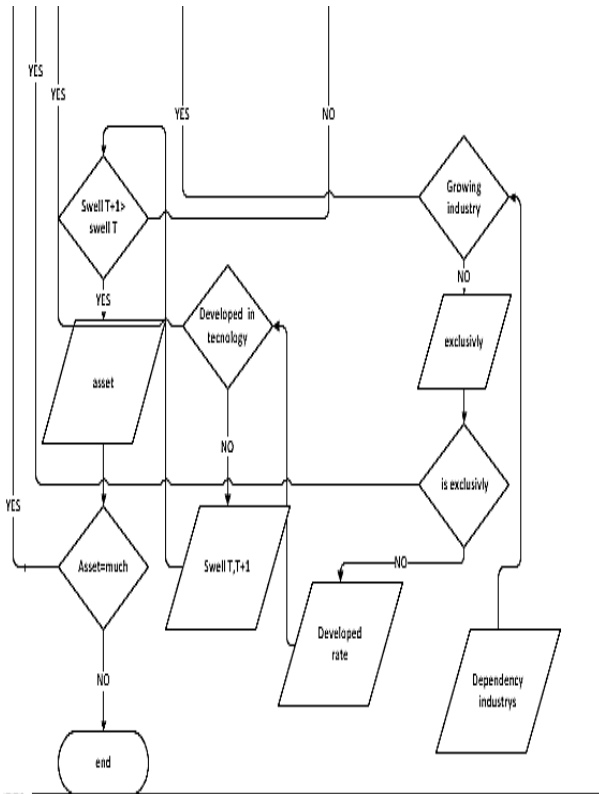
شکل ۳. نمونه‌ای از تقلب تجارت الکترونیک

### ۳- روش پیشنهادی

➤ الگوریتم کشف تقلب قیمت سهم زمانی که قیمت سهم تغییر پیدا می‌کند، بایستی یک سری عوامل بررسی شوند. از جمله این عوامل می‌توان به وضعیت ارز، تصمیمات دولت، روابط بین‌الملل، میزان عرضه و تقاضا، وابستگی صنعت به صنعت دیگر، میزان انحصاری بودن، نرخ تورم و میزان دارایی شرکت مورد نظر اشاره کرد. تک تک این عوامل را بصورت گره یک گراف در نظر

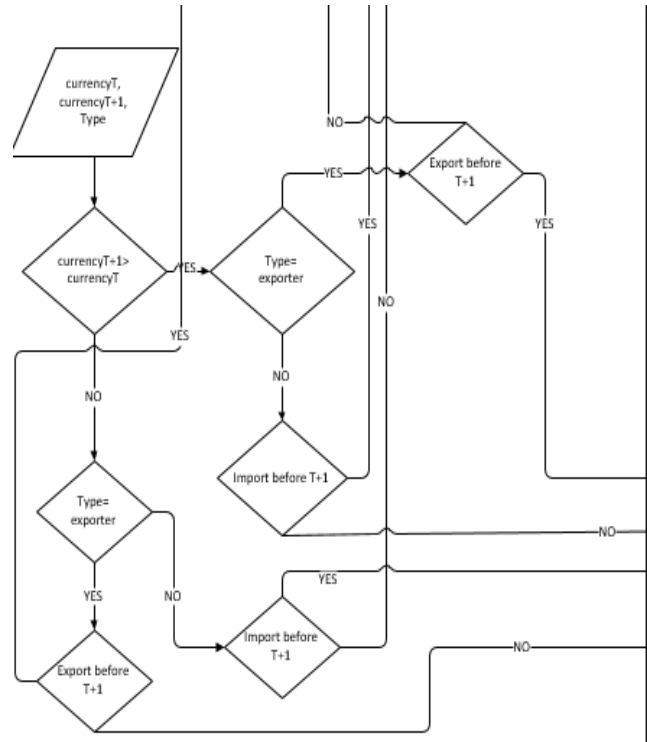


نحوه روابط به هر شکلی باشد باید میزان تقاضا بررسی شود. در صورت زیاد بودن تقاضا، افزایش قیمت سهم طبیعی است. پس قیمت در دو لحظه متوالی دیگر بررسی می‌شود. در غیر این صورت وضعیت صنایع وابسته بررسی می‌شود.



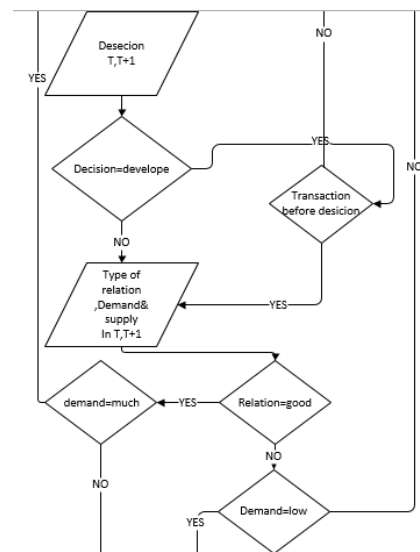
شکل ۷. بخشی از فلوچارت مربوط به سایر عوامل تاثیر گذار بر افزایش قیمت سهم

وضعیت صنایع وابسته بررسی می‌شود. در صورت روبه رشد بودن صنایع وابسته، افزایش قیمت سهم طبیعی است و قیمت در دو لحظه متوالی دیگر بررسی می‌شود. در غیر این صورت عامل انحصاری بودن شرکت، بررسی می‌شود. در صورت انحصاری بودن شرکت، افزایش قیمت سهم طبیعی است. در غیر این صورت عامل نرخ تورم بررسی می‌شود. در صورت افزایش نرخ تورم، افزایش قیمت سهم طبیعی است. در غیر این صورت عامل دارایی شرکت بررسی می‌شود. در صورتی که دارایی شرکت



شکل ۵. بخشی از فلوچارت بررسی وضعیت ارز

اگر دولت تصمیمی مبنی بر افزایش قیمت سهم گرفته باشد، افزایش قیمت سهم طبیعی است و قیمت در دو لحظه متوالی دیگر بررسی می‌شود. در غیر این صورت عامل روابط بین الملل و میزان عرضه و تقاضا بررسی می‌شود. بخشی از فلوچارت مربوطه در شکل ۶ آورده شده است.



شکل ۶. بخشی از فلوچارت روابط بین الملل





#### ۴ - نتایج

داده‌های این کار شامل اطلاعات مالی یک شرکت بورسی صادرکننده است. از جمله داده‌های مالی به کارگرفته شده می‌توان به قیمت سهم در زمان‌های مختلف، قیمت ارز، تغییرات نرخ تورم، زمان صادرکردن، میزان عرضه و تقاضا برای کالاهای شرکت، وضعیت صنایع وابسته به کالاهای شرکت، میزان انحصاری بودن، تصمیمات دولت، میزان پیشرفت و تکنولوژی و مجموع دارایی‌های شرکت که برای سال‌های ۱۳۹۴ و ۱۳۹۵، در نظر گرفته شده است، اشاره کرد. البته به دلیل عدم تاثیر، یک سری از عوامل مذکور در این نوع شرکت، نادیده گرفته شده است. از آن جمله می‌توان به تصمیمات دولت، وضعیت صنایع وابسته (که هر سال بدتر شده است، لذا نمی‌تواند تاثیری در افزایش قیمت سهم داشته باشد) و میزان انحصاری بودن اشاره کرد. چراکه شرکت مربوطه انحصاری نبوده و فضا رقابتی است و افزایش قیمت سهم با درجه بسیار کم‌تری رخ می‌دهد. البته امکان دسترسی به یک سری از اطلاعات به دلیل محرمانه بودن وجود نداشت. (از جمله زمان صادرات)

این روش برای محاسبه خطایی که در افزایش قیمت سهم رخ می‌دهد، استفاده می‌شود. شاید به نظر برسد که چه خطایی در کاهش قیمت سهم ممکن است باشد، اما باید دانست یک سری از شرکت‌ها برای فرار از مالیات قیمت‌های خود را پایین می‌آورند تا خود را زیان ده جلوه دهند. ابتدا داده‌های مربوطه موجود در فایل اکسل به فرمت CSV تبدیل شده و در پایگاه داده neo4j بارگزاری می‌شود. هر رکورد فایل بصورت یک گره ذخیره می‌شود. هر یک از داده‌های مربوطه بصورت

زیاد باشد، افزایش قیمت سهم طبیعی است در غیر اینصورت تقلب رخ داده است. بخشی از فلوچارت مربوطه در شکل ۷ آورده شده است.

پرس وجوی مربوطه جهت بازیابی حلقه تقلب در ادامه بررسی شده است. برای هر گره به ترتیب زمان یک id در نظر می‌گیریم. پرس وجوی لازم برای بازیابی گره آغازگر تقلب در ادامه بررسی می‌شود. روند کار به این شکل است که گره‌های متوالی بررسی می‌شوند. پس باید گره‌هایی که اختلاف id آنها یک واحد است با هم مقایسه شوند. حال به بررسی خصوصیات دو گره متوالی پرداخته خواهد شد. خاصیت قیمت گرهی با id بزرگتر باید بیشتر باشد، چراکه در این کار به بررسی تقلب افزایش قیمت سهم پرداخته شده است. وضعیت ارز بررسی می‌شود. قیمت ارز گرهی با id کمتر باید بیشتر از قیمت ارز گره دیگر باشد. خاصیت عرضه و تقاضا بررسی می‌شود. میزان عرضه و تقاضا گرهی با id کمتر به ترتیب باید کم‌تر و بیشتر از میزان عرضه و تقاضا گره با id بیشتر باشد. خاصیت تورم بررسی می‌شود. نرخ تورم گرهی با id کمتر باید بیشتر از نرخ تورم گرهی با id بیشتر باشد. در نهایت خاصیت دارایی بررسی می‌شود. میزان دارایی گرهی با id کمتر باید بیشتر از میزان دارایی گرهی با id بیشتر باشد. اگر تمامی خصوصیات گرهی با id کم‌تر برابر با خصوصیات گره id بیشتر باشد، دلیلی برای افزایش قیمت سهم وجود ندارد. گره مربوطه باید به عنوان تقلب بازیابی شود. برای نمایش بهتر گره متقلب در پایگاه داده گراف، گره مربوطه را با یک ارتباطی به نام fraud به خودش متصل کرده و تمامی گره‌هایی که ارتباط fraud داشته باشند، به عنوان گره آغازگر تقلب بازیابی می‌شوند.



۱۴۳ مورد تقلب واقعی همه کشف شده‌اند. برای پیاده سازی کار از سیستمی با پردازنده intel(r) core(TM) i7-5500u cpu@2.40GHz و RAM 6.00GB@2.40GHz استفاده شده است و پایگاه داده neo4j مورد استفاده نسخه ۳,۵,۰ آن است.

#### ۵- مراجع

- [1] Han, Jing, et al. "Survey on NoSQL database." 2011 6th international conference on pervasive computing and applications. IEEE, 2011.
- [2] L., Webber, J. & Eifrem, E. (2015). Graph databases: new opportunities for connected data. " O'Reilly Media, Inc." 1-64.
- [3] www.neo4j.com
- [4] Rathle, S., Rathle, G.& Rathle, Ph. (2014). Fraud detection: Discovering connections with graph databases." White Paper-Neo Technology-Graphs are Everywhere,1-11.
- [5] Laurent, C., Laurent, A. & Laurent, A. (2016) "Rogue behavior detection in NoSQL graph databases." Journal of Innovation in Digital Ecosystems 3(2), 70-82.
- [6] Rathle, S., Rathle, G.& Rathle, Ph. (2014). Fraud detection: Discovering connections with graph databases." White Paper-Neo Technology-Graphs are Everywhere,1-1



الهام عبداللہی عابد مدرک کارشناسی خود را در رشته علوم کامپیوتر در سال ۱۳۹۴ از دانشگاه دولتی محقق اردبیلی و مدرک کارشناسی ارشد خود را در همان رشته گرایش سیستم‌های کامپیوتری در سال ۱۳۹۸ از دانشگاه دولتی علامه طباطبائی اخذ کرده است.

خاصیت رکورد در نظر گرفته می‌شود. برای مشاهده میزان ثابت بودن قیمت‌ها، گره‌های با قیمت یکسان را با یک ارتباط به هم متصل می‌کنیم. که پرس‌وجوی مربوط به این کار به ترتیب در شکل های ۹ و ۱۰ نمایش داده شده است.

```
match(n),(m)
where
tointeger(n.id)+1= tointeger(m.id)
and
tointeger(n.price)=tointeger(m.price)
merge (n)-[:to]->(m)
```

شکل ۸. پرس و جوی مربوطه جهت بارگذاری داده

```
USING PERIODIC COMMIT 1000
load csv with headers from "file:///c:/test.csv" as line
create
(:share{id:line.id,date:line.date,time:line.time,
price:line.price,type:line.type,
currency:line.currency,swell:line.swell,
ir:line.ir,supply:line.supply,
demand:line.demand,darayee:line.darayee})
```

شکل ۹. پرس و جوی مربوط جهت ایجاد روابط

```
match(n),(m)
where
(tointeger(n.id)+1= tointeger(m.id) )
and (tointeger(n.price)<tointeger(m.price))
and (tointeger(n.currency)>tointeger(m.currency))
and (tointeger(n.supply)<tointeger(m.supply))
and (tointeger(n.demand)>tointeger(m.demand))
and (tointeger(n.swell)>tointeger(m.swell))
and (tointeger(n.darayee)>tointeger(m.darayee))
merge(m)-[:r:fraud]->(m)
return r as fraud
```

شکل ۱۰. پرس و جوی مربوط جهت بازیابی گره متقلب

زمان لازم برای بارگذاری ۲۴۸۱۶ رکورد و تنظیم ۲۴۸۱۵۶ خاصیت تنها ۴ ثانیه محاسبه شده است. در ضمن برای ادغام یک سری گره‌ها برای ایجاد ۲۴۳۱۸ رابطه ۲۳ ثانیه زمان محاسبه شده است. طبق آزمایش‌های انجام شده، ۴۹۸ مورد تغییر قیمت داشته ایم که شامل هردو حالت افزایش و کاهش در قیمت‌ها می‌باشد. البته در این کار تقلب‌های مربوط به افزایش قیمت سهم با میزان دقت ۰,۹۷ بازیابی می‌شود. تمامی



فرشته آزادی پرنده کارشناسی خود را از دانشگاه علم و صنعت، کارشناسی ارشد را از دانشگاه تربیت مدرس و دکتری خود را از دانشگاه علم و صنعت در رشته

مهندسی کامپیوتر نرم افزار اخذ نموده است.

ایشان از سال ۱۳۹۱ به عنوان هیئت علمی گروه رایانه در دانشگاه علامه طباطبایی مشغول به کار هستند. زمینه های پژوهشی مورد علاقه ایشان سیستم ها و پایگاه داده‌های توزیعی و تئوری بازی می باشد.



علی اصغر صفائی متولد سال ۱۳۵۷ کارشناسی و کارشناسی ارشد خود را در رشته مهندسی کامپیوتر گرایش نرم افزار به ترتیب در سال های ۱۳۸۰ و ۱۳۸۳ به پایان

رسانید. او در سال ۱۳۹۰ موفق به دریافت درجه دکترا در همین رشته و در حوزه پایگاه داده‌ها از دانشگاه علم و صنعت ایران شد. وی هم‌اکنون به‌عنوان عضو هیئت علمی گروه انفورماتیک پزشکی دانشگاه تربیت مدرس مشغول به فعالیت است. زمینه‌های پژوهشی ایشان شامل سامانه‌های پایگاه داده و جریان داده، داده‌های عظیم، پردازش موازی و بی‌درنگ پرس‌وجوها، حافظه‌های نهان معنایی، امنیت در سامانه‌های پایگاه داده و مدیریت اطلاعات در مقیاس وب است. از ایشان کتاب و مقالات متعددی نیز در همین زمینه‌ها به چاپ رسیده است.