



پردیس علوم  
دانشکده ریاضی، آمار و علوم کامپیوتر

# کاربرد های داده کاوی در مهندسی نرم افزار

نگارنده

زهرا امیدی

استاد راهنما

مهندس ابراهیم نقیب زاده مشایخ

پایان نامه برای دریافت درجه کارشناسی  
در رشته علوم کامپیوتر

تیر ماه ۱۳۹۷

## چکیده

دسترسی بیشتر به داده های ایجاد شده به عنوان بخشی از فرآیند تولید نرم افزار به ما این امکان را می دهد تا تکنیک های جدید تجزیه و تحلیل را برای این داده ها به کار گیریم و از نتایج آن ها برای هدایت بهینه سازی فرآیند استفاده کنیم. در این نوشتار، منابع اطلاعاتی مختلفی را معرفی می کنیم و اصول و تکنیک های داده کاوی را در مورد داده های مهندسی نرم افزار مورد بحث قرار می دهیم. داده هایی که می توانند استخراج شوند، توسط بیشتر بخش های فرآیند تولید ایجاد می شوند: استخراج نیازها، تجزیه و تحلیل تولید، آزمایش، اشکال زدایی و نگهداری. بر اساس این رده بندی، ما روش های کاوشی را که مورد استفاده قرار گرفته اند، مورد بررسی قرار می دهیم و آن ها را با توجه به بخش های متناظر فرآیند تولید و وظایفی که دارند، رده بندی می کنیم. بنابراین این مطالعه، محققان را با یک مرور کلی از تکنیک های داده کاوی که به داده های مهندسی نرم افزار مربوط می شوند، آشنا می کند و به محققان در انتخاب تکنیک های داده کاوی مناسب برای کار خود کمک می کند.

# پیش‌گفتار

پیشرفت‌های اخیر تکنولوژی مدیریت داده‌ها به همراه ابزارها و روش‌های جمع‌آوری کارا، ذخیره‌سازی و نمایه‌سازی داده‌ها را فراهم می‌کنند. حجم زیادی از داده‌ها در تولید نرم‌افزار تولید می‌شوند که سازمان‌های نرم‌افزاری آن‌ها را به امید استخراج اطلاعات مفید و به دست آوردن درک بهتری از فرآیندها و محصولاتشان جمع‌آوری می‌کنند. همچنین مقدار قابل توجهی از فراداده‌های پروژه‌های نرم‌افزاری متن‌باز (OSS) در مخازن نرم‌افزار جمع‌آوری و نگهداری می‌شوند. به عنوان مثال، پروژه *FLOSSmole*<sup>1</sup> داده‌های چند پروژه‌های نرم‌افزاری متن‌باز را باهم تلفیق می‌کند و آزادانه آن‌ها را در قالب‌های مختلف ارائه می‌دهد. پس یک فراوانی افزایشی در داده‌های ذخیره‌شده در مخزن نرم‌افزار وجود دارد که می‌تواند نقش مهمی در بهبود بهره‌وری و کیفیت نرم‌افزار داشته باشد. به ویژه، داده‌ها در تولید نرم‌افزار می‌توانند به نسخه‌های برنامه، ردیاب‌های اجرا، گزارش‌های خطا / اشکال و پکیج‌های متن‌باز اشاره کنند. لیست‌های پستی، انجمن‌های گفتگو و خبرنامه‌ها نیز اطلاعات مفیدی در مورد یک قسمت از نرم‌افزار ارائه می‌دهند. این رشد انفجاری توانایی مهندسان نرم‌افزار برای جمع‌آوری و ذخیره اطلاعات مهندسی نرم‌افزار، نیاز به ابزاری جدید، مقیاس‌پذیر و کارآمد برای تجزیه و تحلیل داده‌ها را ایجاد کرده است.

داده‌های تک‌تک‌هایی را برای تجزیه و تحلیل و استخراج الگوهای جدید و جالب از داده‌ها فراهم می‌کند. به طور رسمی، داده‌های کاوی به عنوان فرآیند استخراج اطلاعات ناشناخته و مفید بالقوه از مجموعه داده‌ها، تعریف شده است. بنابراین استخراج اطلاعات مهندسی نرم‌افزار به تازگی توجه محققان

---

<http://flossmole.org/><sup>1</sup>

را به عنوان یک وسیله امیدوار کننده برای رسیدن به هدف بهبود نرم افزار جلب کرده است. الگوهای استخراج شده دانش می توانند مهندسی نرم افزار را در پیش بینی، برنامه ریزی و درک جنبه های مختلف پروژه یاری کنند تا بتوانند به طور مؤثرتری از فعالیت های توسعه و مدیریت پروژه حمایت کنند. چالش های متعددی در کاوش مخازن نرم افزاری وجود دارد:

- داده های پیچیده : مهندسان نرم افزار معمولاً از داده های با انواع منحصر بفردی برای انجام کار مهندسی نرم افزار استفاده می کنند. با این حال، حوزه مهندسی نرم افزار می تواند پیچیده باشد و در نتیجه، وظایف مهندسی نرم افزار به طور فزاینده ای نیاز به استخراج انواع داده های مرتبط چندگانه برای دستیابی به نتایج مؤثرتری دارد. علاوه بر این، مواردی وجود دارند که اطلاعات مهم، نه تنها به فقره های منحصر بفرد داده، بلکه به ارتباط بین آن ها وابسته است. بنابراین نیاز به تکنیک هایی که تجزیه و تحلیل انواع داده های پیچیده و همچنین ارتباط بین داده های مهندسی نرم افزار را در نظر می گیرند، در مهندسی نرم افزار، قوی تر از همیشه است.

- داده های بزرگ : حجم بسیار زیادی از اطلاعات در مخازن مهندسی نرم افزار جمع آوری و ذخیره می شود. در اغلب موارد مهندسان نرم افزار تنها چند مخزن محلی را برای انجام وظایف خود تجزیه و تحلیل می کنند. در حالی که، ممکن است مقدار بسیار کمی از داده های مربوطه در مخازن محلی مهندسی نرم افزار برای کمک به استخراج الگوهای مطلوب (به عنوان مثال استخراج الگوهای روش های API) وجود داشته باشد. اگر ما بتوانیم مخازن نرم افزاری در سطح اینترنت را کاوش کنیم، می توان این مشکل را حل کرد. همچنین تکنیک های کاوش را می توان برای بسیاری از مخازن نرم افزاری در یک سازمان یا در سراسر سازمان و یا حتی در کل جهان متن باز به کار برد. بنابراین نیاز به تکنیک هایی که قادر به کاوش کارآمد داده های با مقیاس های بزرگ هستند، بوجود می آید.

در این نوشتار، ما یک مرور کلی از رویکردهایی که هدفشان متصل کردن بخش های تحقیقاتی داده کاوی و مهندسی نرم افزار به منظور ایجاد تکنیک

های کارآمدتر برای پردازش نرم افزار است، ارائه می دهیم. در فصل اول، مقدمه ای بر مفاهیم و رویکردهای اصلی داده کاوی ارائه می دهیم، در حالی که در فصل دوم، انواع مختلفی از اطلاعات مهندسی نرم افزار را که می توانیم کاوش کنیم، توصیف می کنیم. و در فصل ۳ به بررسی روش هایی که از تکنیک های داده کاوی در زمینه ی مهندسی نرم افزار استفاده می کنند، می پردازیم. به ویژه، این مطالعه بر اساس ویژگی های اصلی زیر است:

۱. روش داده کاوی استفاده شده (فصل ۱)

۲. داده های مهندسی نرم افزار که روش های داده کاوی روی آن ها اعمال می شود (فصل ۲)

۳. وظایف مهندسی نرم افزار که به کمک روش های داده کاوی می توانند انجام شوند (فصل ۳).

بنابراین، هدف این مطالعه معرفی مفاهیم و تکنیک های اساسی به محققان است که بتوانند برای به دست آوردن درک بهتر از فرآیندهای مهندسی نرم افزار، داده کاوی را به کار ببرند. به طور موازی، محققان می توانند تکنیک های اخیر را برای درک بهتر اطلاعات پروژه، شناسایی محدودیت های فرایندهای فعلی و تعریف روش هایی که وظایف مهندسی نرم افزار را تسهیل می کنند، به کار ببرند. به طور خلاصه، مفاهیم کلیدی این کار عبارتند از:

- رده بندی رویکردهای استفاده شده در کاوش داده های مهندسی نرم افزار، با توجه به ناحیه هایی از مهندسی نرم افزار که آنها به کار می روند.
- چارچوب تجزیه و تحلیل مبتنی بر ماتریس، پل ارتباطی مهندسی نرم افزار با روش های داده کاوی است.
- گردهم آوردن ناحیه های تحقیقاتی داده کاوی و مهندسی نرم افزار. تعدادی از رویکردهایی که از داده کاوی در وظایف مهندسی نرم افزار استفاده می کنند، روش های کاری جدیدی را برای پژوهشگران و متخصصان در مهندسی نرم افزار ارائه می دهند.

# فهرست مطالب

۱	داده کاوی و کشف دانش	۱
۶	اطلاعات مهندسی نرم افزار	۲
۶	مستندات	۱.۲
۷	اطلاعات مدیریت پیکربندی نرم افزار	۲.۲
۸	کد منبع	۳.۲
۸	کد کامپایل شده و ردیابی اجرا	۴.۲
۹	پایگاه داده های اشکال و ردیابی موضوع	۵.۲
۹	لیست های پستی	۶.۲
۱۰	داده کاوی برای مهندسی نرم افزار	۳
۱۱	ردیابی و استخراج نیازها	۱.۳
۱۱	رده بندی	۱.۱.۳
۱۲	خلاصه سازی اطلاعات	۲.۱.۳
۱۲	تجزیه و تحلیل تولید	۲.۳
۱۳	خوشه بندی	۱.۲.۳
۱۴	رده بندی	۲.۲.۳
۱۶	قوانین انجمنی و کاوش الگوهای تکرار شونده	۳.۲.۳
۱۷	آزمایش کردن	۳.۳
۱۸	خوشه بندی	۱.۳.۳
۱۸	رده بندی	۲.۳.۳
۲۰	اشکال زدایی	۴.۳
۲۱	رده بندی شکست های نرم افزار	۱.۴.۳
۲۳	اصلاح رده بندی شکست ها با استفاده از دندروگرام ها	۲.۴.۳
۲۵	پالایش با استفاده از درختان رده بندی	۳.۴.۳

۲۶	قوانین انجمنی و کاوش الگوهای تکرار شونده .	۴.۴.۳	
۲۷	..... ننگه داری	۵.۳	
۲۷	..... خوشه بندی	۱.۵.۳	
۲۸	..... رده بندی	۲.۵.۳	
۲۹	قوانین انجمنی .	۳.۵.۳	
۳۱	..... استفاده مجدد از نرم افزار .	۶.۳	
۳۱	..... رده بندی	۱.۶.۳	
۳۳	قوانین انجمنی .	۲.۶.۳	
۳۵			۴ نتیجه گیری

## فصل ۱

# داده کاوی و کشف دانش

اهداف اصلی داده کاوی پیش‌بینی و توصیف است. پیش‌بینی برآورد کردن مقدار آینده یا پیش‌بینی مقدار متغیرهای هدف بر مبنای مطالعه‌ی سایر متغیرها است. توصیف به کشف الگوها به منظور کمک به ارائه‌ی داده‌ها به شیوه‌ای قابل فهم‌تر و قابل بهره‌برداری‌تر، تمرکز دارد. یک توصیف خوب، توضیح خوبی از رفتار داده‌ها ارائه می‌دهد. اهمیت نسبی پیش‌بینی و توصیف برای کاربردهای مختلف داده‌کاوی، متفاوت است. با این حال، در مورد کشف دانش، به نظر می‌رسد که توصیف مهم‌تر از پیش‌بینی باشد، بر خلاف برنامه‌های تشخیص الگو و یادگیری ماشین، که پیش‌بینی مهم‌تر است. داده‌کاوی، با توضیح و تجزیه و تحلیل عمیق مفاهیم و فرآیندهای نرم‌افزاری، به وظایف مهندسی نرم‌افزار کمک می‌کند. بر اساس تکنیک‌های داده‌کاوی، می‌توانیم روابط میان پروژه‌های نرم‌افزاری را استخراج کنیم. داده‌کاوی می‌تواند از اطلاعات استخراج شده برای ارزیابی پروژه‌های نرم‌افزاری و یا پیش‌بینی رفتار نرم‌افزار، بهره‌برد. تعدادی از روش‌های داده‌کاوی برای برآورده کردن نیازهای برنامه‌های مختلف پیشنهاد شده‌اند. همه آن‌ها مجموعه‌ای از عملکردهای داده‌کاوی را برای شناسایی و توصیف الگوهای جالب در دانش استخراج شده از مجموعه داده‌ها فراهم می‌کنند. در زیر ما به طور مختصر وظایف اصلی داده‌کاوی و نحوه استفاده از آن‌ها در مهندسی نرم‌افزار را شرح می‌دهیم.

- خوشه بندی - تکنیک یادگیری بی نظارت  
خوشه بندی یکی از کارهای مفید در زمینه داده کاوی است، که برای



کشف گروه‌ها و شناسایی توزیع‌ها و الگوهای جالب در داده‌های اساسی به کار گرفته می‌شود. مساله‌ی خوشه‌بندی در مورد بخش بندی یک مجموعه داده‌ی داده شده، به گروه‌ها (خوشه‌ها) است به طوری که نقاط داده‌ی یک خوشه بیشتر به یکدیگر شبیه باشند تا به نقاط خوشه‌های مختلف. در فرایند خوشه‌بندی، هیچ رده‌ی از پیش تعریف شده و یا نمونه‌ای که نوع روابط مطلوب و معتبر بین داده‌ها را پیش بینی کند، وجود ندارد. بنابراین این مساله به عنوان یک فرآیند یادگیری بدون نظارت شناخته می‌شود. خوشه‌بندی را می‌توان برای ایجاد یک دید کلی از توزیع داده‌های اساسی و همچنین شناسایی خودکار داده‌های دورافتاده یا منزوی، به کار برد. در مهندسی نرم‌افزار، خوشه‌بندی می‌تواند برای تعریف گروه‌های مولفه‌های مشابه بر اساس تعداد تغییرات و معیار تعداد سیکلمات (تعداد مسیرهای مستقل خطی در کد منبع برنامه) استفاده شود.

#### ● رده‌بندی - تکنیک یادگیری با نظارت

مسئله رده‌بندی به طور گسترده در زمینه‌های تحقیقاتی آمار، شناخت الگو و یادگیری ماشین به عنوان یک راه حل ممکن برای مسئله‌ی اکتساب دانش یا استخراج دانش مورد مطالعه قرار گرفته است. رده‌بندی یکی از وظایف اصلی در داده کاوی برای اختصاص دادن یک فقره داده به مجموعه‌ای از رده‌های از پیش تعریف شده، است. رده‌بندی را می‌توان به عنوان یک تابع (رده‌بند) تعریف کرد که یک فقره داده را به یکی از چند رده از پیش تعریف شده، نگاشت می‌کند. مجموعه‌ای از رده‌های از پیش تعریف شده و یک مجموعه‌ی آموزشی از نمونه‌های از پیش رده‌بندی شده، رده‌بندی را مشخص می‌کند. برعکس، فرایند خوشه‌بندی به نمونه‌ها و رده‌های از پیش تعریف شده وابسته نیست. هدف فرایند رده‌بندی، ایجاد یک مدل است که می‌تواند برای رده‌بندی فقره‌های آینده که رده‌بندی آنها ناشناخته است، مورد استفاده قرار گیرد. یکی از تکنیک‌های رده‌بندی که به طور گسترده استفاده می‌شود، ساخت درخت‌های تصمیم است. آن‌ها را می‌توان برای کشف قوانین رده‌بندی برای یک ویژگی منتخب از یک مجموعه داده، با تقسیم کردن سیستماتیک اطلاعات موجود در این مجموعه داده، مورد استفاده قرار داد. درخت‌های تصمیم نیز یکی از ابزارهایی هستند که برای ساخت مدل‌های رده‌بندی

در زمینه مهندسی نرم افزار انتخاب شده اند. شکل ۱.۱ یک درخت رده بندی را نشان می دهد که برای ارائه ی یک سازوکار شناسایی مولفه های نرم افزاری خطرناک بر اساس ویژگی های مولفه و سیستم آن، ساخته شده است. بنابراین، بر اساس درخت تصمیم داده شده، می توانیم قانون زیر را که به تصمیم گیری در مورد خطاها در یک مولفه کمک می کند، استخراج کنیم:

IF (# of data bindings > 10) AND (it is part of a non real-time system)  
THEN the module is unlikely to have errors

● قوانین انجمنی و کاوش الگوهای تکرار شونده

کاوش قوانین انجمنی علاقه مندی بسیاری را به خود جلب کرده است؛ چرا که قوانین، یک روش مختصر برای ارائه ی اطلاعات مفید بالقوه است که به راحتی توسط کاربران نهایی قابل درک است. قوانین انجمنی ”همبستگی“ اساسی بین ویژگی های یک مجموعه داده را نشان می دهد. این همبستگی ها به شکل زیر ارائه می شوند:

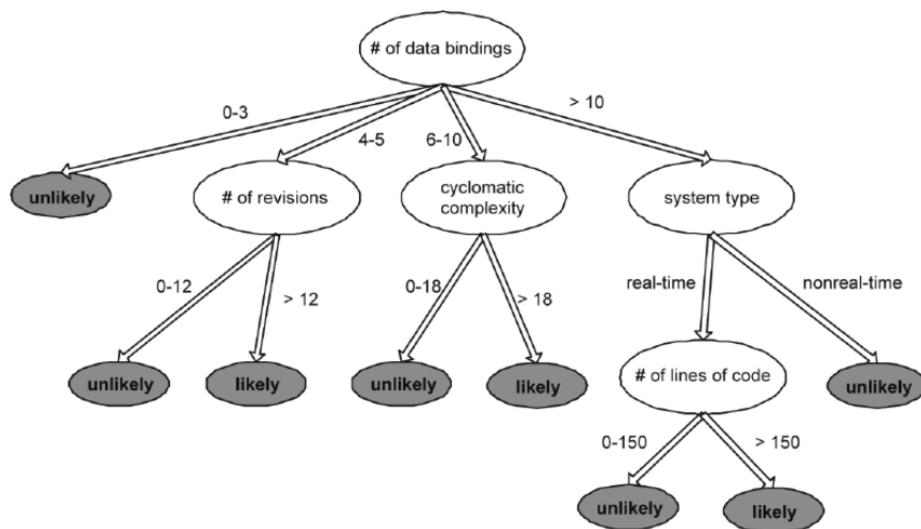
$$A \rightarrow B$$

که A و B به مجموعه ای از ویژگی ها در داده های اساسی اشاره دارند. بنابراین، آن ها برای استخراج اطلاعات بر اساس انطباق در یک مجموعه داده استفاده می شوند. به عنوان مثال، تجزیه و تحلیل کارنامه ی خطاهای سیستم که در مولفه های نرم افزاری کشف شده اند، می تواند روابط بین رویدادهای استنتاجی را بر اساس ویژگی های مولفه نرم افزاری و دسته بندی های خطاها استخراج کند. چنین قاعده ای به فرم زیر خواهد بود:

(large/small size, large/small complexity, number of revisions) → (interface error, missing or wrong functionality, algorithms or data structure error etc.)

● خلاصه سازی و تشریح داده ها.

تشریح داده ها خلاصه ای از ویژگی های داده برای یک رده خاص داده است. داده ها بر اساس نیازهای مشخص شده توسط کاربران گردآوری می شوند. این تکنیک ها می توانند برای کشف مجموعه ای از الگوها از مخازن مهندسی نرم افزار که ویژگی های خاصی را برآورده می کنند، استفاده شوند.



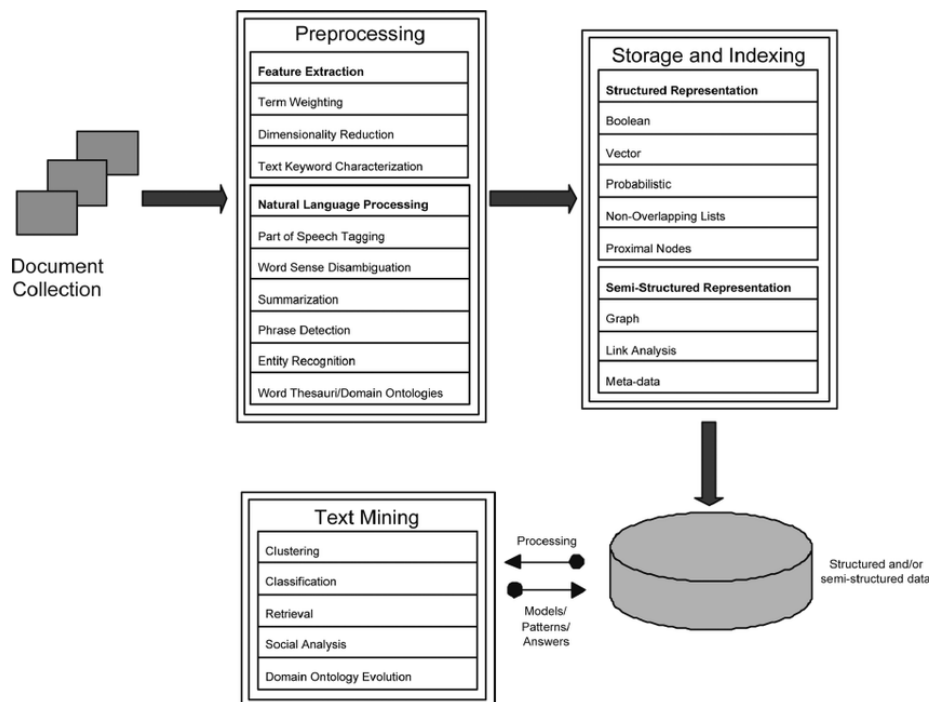
شکل ۱.۱: درخت رده بندی برای شناسایی مولفه های نرم افزار خطرناک.

● تشخیص تغییر و انحراف.

تشخیص تغییر و انحراف، بر کشف مهم ترین تغییرات در داده ها از مقادیری که قبلا اندازه گیری شده است، تمرکز دارد. بنابراین این تکنیک ها می توانند به شناسایی تغییرات کد منبع یا شناسایی تفاوت های بین الگوهای استخراج شده از مخازن مهندسی نرم افزار کمک کنند.

رویکردهای مختلفی برای انجام وظایف داده کاوی بالا و به کار بردن انواع مختلف داده ها توسعه داده شده اند. آنها تکنیک ها را از جنبه های مختلف مدیریت داده ها و تجزیه و تحلیل داده ها، از جمله شناخت الگو، یادگیری ماشین، آمار، بازیابی اطلاعات، تجزیه و تحلیل متن و مفهوم به کار می برند. متن کاوی (کاوش متن) به عنوان یک مورد خاص از داده کاوی معرفی می شود و به فرایند استخراج اطلاعات از متن اشاره دارد. مخازن مهندسی نرم افزار، شامل اطلاعات متنی مانند کد منبع، لیست های پستی، گزارش های اشکال و کارنامه ی اجرا است.

کاوش منابع متنی برای بسیاری از فعالیت های مهم در مهندسی نرم افزار ضروری است: ردیابی نیازها؛ بازیابی اجزای مخزن؛ شناسایی و پیش بینی ناکارایی های نرم افزار؛ نگهداری نرم افزار؛ آزمایش و غیره. روش های به کار گرفته شده در متن کاوی، بسته به کاربرد، معمولا نیاز به تبدیل متون به یک



شکل ۲.۱: پیش پردازش، ذخیره سازی و پردازش متن در متن کاوی.

نمایش ساختاری حد واسط دارند، که می‌تواند به عنوان مثال ذخیره‌سازی متون در یک سیستم مدیریت پایگاه داده، بر اساس یک شمایی خاص باشد. اگرچه در بسیاری از رویکردها، بهتر است یک فرم حد واسط نیمه ساخت یافته‌ی متن نیز نگهداری شود، برای مثال مثل نمایش مستندات در قالب یک گراف، که در آن تجزیه و تحلیل جامعه و تکنیک‌های گراف می‌توانند مورد استفاده قرار گیرند. مستقل از اهداف کاری، متن کاوی نیاز به تکنیک‌های پیش پردازشی دارد که معمولاً تحلیل کیفی و کمی ویژگی‌های مستندات را به همراه می‌آورد. در شکل ۲.۱، نمودار مهمترین مراحل تحلیل پیش پردازش و همچنین مهمترین تکنیک‌های متن کاوی نشان داده شده است.

## فصل ۲

# اطلاعات مهندسی نرم افزار

ماهیت داده‌هایی که توسط تکنیک‌های داده‌کاوی در مهندسی نرم‌افزار مورد استفاده قرار می‌گیرند، می‌تواند به عنوان روشی متمایز از روش‌های پایه‌ای عمل کند، زیرا بر پیش‌پردازش و همچنین پس‌تحلیل تأثیر می‌گذارد. در ادامه، ما منابع مختلفی از اطلاعات مهندسی نرم‌افزار را که داده‌کاوی برای آن‌ها به کار رفته است ارائه می‌دهیم. این ارائه همچنین تلاش می‌کند تا سختی تهیه داده‌ها را برای پردازش منعکس کند.

### ۱.۲ مستندات

داده‌های مستندات نرم‌افزاری از اهمیت زیادی برخوردار هستند اما در کنار آن از پیچیدگی بالایی نیز برای پردازش توسط تکنیک‌های داده‌کاوی برخوردارند. برنامه کاربردی، مدیریت سیستم و مستندات کد منبع، یک میانگیر بزرگ مستندات و متن رایگان برای تجزیه و تحلیل و کاوش نرم‌افزار تشکیل می‌دهند. از میان قسمت‌های اطلاعاتی متن، که می‌توان آن را با ارزش‌ترین مقدار برای استفاده در تکنیک‌های کاوش در نظر گرفت، می‌توان به توضیحات نرم‌افزار، پی‌کربندی استفاده و راه‌اندازی، راهنمای کاربر، مسائل مربوط به مدیریت فایل، ورود به سیستم، مجوز و مسائل مربوط به سازگاری اشاره کرد. علاوه بر مستندات خارجی نرم‌افزار، مستندات داخلی نیز ممکن است نقش منبع داده مهم را بازی کنند. اگرچه اکثر مستندات در چند دسته‌ی مختلف از اسناد،

مانند قالب سند قابل حمل، فایل های سیستم حروف چینی و متن، اچ تی ام ال قرار می گیرند. با توجه به تنوع انواع اسناد و داده های متنی مورد استفاده، لازم است که یک مولفه پیش پردازش برای اطلاعات اسناد، که قادر به استفاده از تجزیه کننده های تمام انواع اسناد ذکر شده در بالا است، وجود داشته باشد. یکی دیگر از منابع مهم اطلاعاتی که در مستندات نرم افزار قرار دارد، داده های چند رسانه ای است. دستورالعمل های شکل ها، و همچنین دستورالعمل های صوتی و تصویری، همگی می توانند به عنوان مقدار اطلاعات اضافه شده در نظر گرفته شوند. در چنین مواردی، تکنیک های کاوش چند رسانه ای باید برای به عهده گرفتن پیش پردازش و پس پردازش، و افزایش هزینه های کلی پردازش صورت گیرد.

## ۲.۲ اطلاعات مدیریت پیکربندی نرم افزار

داده های به وجود آمده از سیستم های مدیریت پیکربندی نرم افزار را ممکن است شامل کد نرم افزاری، مستندات، مدل های طراحی، حسابداری وضعیت، ردیابی نقص و همچنین داده های کنترل تجدید نظر (مستندات و نظرات همراه با نسخه های نرم افزاری در *CVS* تصویب شده) باشند. مستقل از سیستم کنترل نسخه پایه (متمرکز یا توزیع شده) مقدار داده های در دسترس *SCM* بزرگ است و بنابراین یک مطالعه دقیق و درک تمیز از دامنه نرم افزار مورد نیاز است، در نتیجه ارزشمندترین داده ها نگهداری می شوند. اکثر داده های *SCM* ها متن ساختار یافته است.

سیستم های مدیریت کد منبع غیر متمرکز<sup>۲</sup> در سال های گذشته رشد چشمگیری داشته است. *DSCM* ها می توانند به محققان مهندسی نرم افزار، داده های جدید و مفید ارائه دهند که آن ها را قادر می سازند تا فرآیندهای نرم افزاری را بهتر درک کنند. با این حال داده های *DSCM* باید با دقت کاوش شوند. از آنجا که تفاوت هایی در معنای اصطلاحات استفاده شده در سیستم های مدیریت کد منبع متمرکز و غیرمتمرکز وجود دارد، تله های بالقوه ای ممکن است در تجزیه و تحلیل این داده ها وجود داشته باشند.

---

Software configuration management systems (SCMs)<sup>۱</sup>  
Decentralized source code management (DSCM)<sup>۲</sup>

## ۳.۲ کد منبع

کد منبع برای داده‌کاوی در مهندسی نرم‌افزار می‌تواند منبع مهمی از داده‌ها باشد. کاربردهای متنوع داده‌کاوی در مهندسی نرم‌افزار، از کد منبع برای کمک به نگهداری نرم‌افزار، درک برنامه و تجزیه و تحلیل اجزای نرم‌افزار استفاده کرده‌اند. جزئیات این رویکردها در فصل ۳ مورد بحث قرار می‌گیرد. پیش‌پردازش اولیه برای کد منبع در دست، همیشه یک هشدار است، از آنجا که یک تجزیه‌کننده برای زبان کد منبع مربوطه باید در دسترس باشد. پس از تجزیه، کد منبع را می‌توان به عنوان متن ساختار یافته مشاهده کرد. جنبه‌های مرکزی استفاده از تکنیک‌های داده‌کاوی برای کد منبع شامل پیش‌بینی تغییرات آینده از طریق کاوش تاریخچه‌ی تغییرات، پیش‌بینی انتشار تغییرات، نقص‌های حاصل از تاریخچه نهان و همچنین پیش‌بینی تراکم نقص در فایل‌های کد منبع است.

## ۴.۲ کد کامپایل شده و ردیابی اجرا

کد کامپایل شده در قالب کد هدف، یکی از منابع داده جایگزین برای استفاده از تجزیه و تحلیل پایا در مهندسی نرم‌افزار است. همچنین کد کامپایل شده به عنوان یک منبع داده برای تکنیک‌های داده‌کاوی به منظور کمک به تشخیص نرم‌افزارهای مخرب مورد استفاده قرار گرفته است. علاوه بر این، اصول کاوش وب به طور گسترده در برنامه‌های قابل اجرای شی‌گرا به منظور کمک به درک برنامه برای روش‌های مهندسی معکوس استفاده شده است. هنگامی که مولفه‌ها و اجزای نرم‌افزاری مورد آزمایش قرار می‌گیرند، زنجیره‌ای از رویدادها رخ می‌دهد که در ردیابی اجرا<sup>۳</sup> ثبت می‌شوند. استخراج الگویی اجرا نیز در ردیابی اجرا تحت چارچوب تجزیه و تحلیل پویا به منظور کمک به استخراج عملکردهای سیستم‌های نرم‌افزاری استفاده شده است.

<sup>۳</sup>execution traces

## ۵.۲ پایگاه داده‌های اشکال و ردیابی موضوع

پایگاه داده‌های گزارش اشکال<sup>۴</sup> یا ردیابی موضوع<sup>۵</sup>، مهمترین گزارش موضوع در سیستم‌های نرم‌افزاری هستند. داده‌های ساخت یافته (چندتایی‌های پایگاه داده) شامل توضیحات یک موضوع، جزئیات گزارش دهنده و تاریخ / زمان هستند که سه نوع اطلاعات استاندارد<sup>۴</sup> می‌توانند در پایگاه داده‌های ردیابی موضوع یافت شوند. تکنیک‌های یادگیری ماشین در گذشته برای پیش‌بینی تخصیص صحیح برنامه‌نویسان به اشکال‌ها، پاک کردن پایگاه داده از ظهور یک خطای مشابه، یا حتی پیش‌بینی مولفه‌های نرم‌افزاری که در همان زمان تحت اثر اشکال‌های گزارش شده بوده‌اند، به طور موفقیت‌آمیزی استفاده شده است.

## ۶.۲ لیست‌های پستی

سیستم‌های نرم‌افزاری بزرگ، و به ویژه نرم‌افزارهای متن‌باز، لیست‌های پستی را به عنوان پل ارتباطی بین کاربران و برنامه‌نویسان ارائه می‌دهند. لیست‌های پستی، داده‌ی سخت به شمار می‌آیند، زیرا حاوی حجم زیادی از متن‌های آزاد هستند. گراف‌های پیام و نویسنده را به راحتی می‌توان از داده‌ها بیرون کشید، اما تجزیه و تحلیل محتوای آن‌ها سخت است زیرا از آن جایی که احتمالاً پیام‌ها از پاسخ‌هایی تشکیل شده‌اند، نیاز به بررسی بحث‌های اولیه و یا قبلی در لیست‌های پستی مطرح می‌شود. کاربردهای داده‌کاوی در لیست‌های پستی شامل تجزیه و تحلیل متن، خوشه‌بندی متن موضوعات مورد بحث، و تجزیه و تحلیل زبان‌شناسی پیام‌ها برای برجسته کردن پروفایل‌ها و ویژگی‌های شخصیتی برنامه‌نویسان می‌شود، اما محدود به این‌ها نیست.



## فصل ۳

# داده کاوی برای مهندسی نرم افزار

با توجه به توانایی داده کاوی برای مقابله با حجم زیاد داده‌ها و کارایی آن برای شناسایی الگوهای پنهان دانش، داده کاوی در تعدادی از کارهای تحقیقاتی به عنوان راهی برای پشتیبانی از نگهداری نرم افزار مقیاس صنعتی، اشکال زدایی، آزمایش ارائه شده است. نتایج کاوش می‌توانند به مهندسان نرم افزار برای پیشگیری از خرابی نرم افزار، استخراج و رده بندی اشکالات رایج، شناسایی روابط بین رده‌ها در یک کتابخانه، تجزیه و تحلیل داده‌های نقص، کشف الگوهای دوباره استفاده شده در کد منبع و به همین ترتیب خودکارسازی روند تولید کمک کنند. به طور کلی با استفاده از داده کاوی، متخصصین و محققان می‌توانند اطلاعات بالقوه‌ی مهندسی نرم افزار را بیاموزند و از نتایج کاوش برای مدیریت بهتر پروژه‌های خود و تولید سیستم‌های نرم افزاری با کیفیت بالاتر که در زمان و بودجه‌ی تعیین شده تحویل داده می‌شوند، استفاده کنند. در بخش‌های بعدی، در مورد ویژگی‌های اصلی رویکردهای کاوشی که در مهندسی نرم افزار مورد استفاده قرار گرفته‌اند و نحوه‌ی استفاده از آن‌ها در چرخه حیات مهندسی نرم افزار، بحث می‌کنیم. ما رویکردها را با توجه به وظایف مهندسی نرم افزاری که آن‌ها کمک می‌کنند و تکنیک‌های کاوشی که آن‌ها استفاده می‌کنند، رده بندی می‌کنیم.

جدول ۱.۳: رویکرد های کاوش استفاده شده در ردیابی و استخراج نیازها

نتایج تحلیل داده ها	داده ورودی	رویکرد کاوش
نیازها	مستندات	رده بندی
نیازها	SCM و لیست های پستی	بازیابی اطلاعات
	کارنامه CVS	خلاصه سازی اطلاعات

## ۱.۳ ردیابی و استخراج نیازها

در این بخش ما درباره‌ی چگونگی به کار بردن تکنیک‌های تجزیه و تحلیل داده‌ها برای استخراج یا ردیابی سیستم نیازها بحث می‌کنیم. تحقیقات انجام شده در زمینه‌ی تجزیه و تحلیل نیازها به داده کاوی با بالاترین سطح آن اشاره می‌کنند، از جمله فعالیت‌ها و متدلوژی‌های خاص و مرتبط از آمار، یادگیری ماشین و بازیابی اطلاعات. جدول ۱.۳ ویژگی‌های اصلی تکنیک‌های مورد بحث در زیر را خلاصه می‌کند.

### ۱.۱.۳ رده بندی

یک رویکرد جدید، بر بهبود بهره‌برداری از نیازهای سطح بالا و سطح پایین با استفاده از بازیابی اطلاعات متمرکز شده است. به طور خاص، آن‌ها جهان مستندات را به عنوان اجتماع عناصر طراحی و نیازهای فردی در نظر می‌گیرند و آن‌ها مساله ردیابی نیازها را به پیدا کردن شباهت‌های بین نمایش فضای برداری نیازهای سطح بالا و سطح پایین نگاشت می‌کنند، به این ترتیب آن را به یک وظیفه‌ی بازیابی اطلاعات کاهش می‌دهند. به عنوان گسترش این مطالعه، محققان بر کشف عوامل موثر بر رفتار یک تحلیلگر در هنگام کار با نتایج بدست آمده از ابزارهای داده کاوی در مهندسی نرم‌افزار تمرکز کرده‌اند. تمام تحقیق براساس فرضیه‌ی تایید شده‌ی دقت ردیابی‌های کاندیدای تولید شده توسط رایانه بر دقت ردیابی‌های تولید شده توسط تحلیلگر تاثیر می‌گذارد، است. این مطالعه نشان می‌دهد که چگونه عملکرد ابزارهایی که از طریق استفاده از بازیابی اطلاعات، نیازهای سطح بالا و سطح پایین را استخراج می‌کنند، بر زمان مصرف شده توسط یک تحلیلگر برای ارائه‌ی بازخورد و همچنین بر عملکرد آن تاثیر می‌گذارد. نتایج

نشان می‌دهند که سیستم-های داده‌کاوی فراخوان<sup>۱</sup> کمی در بازخورد زمان مصرف شده‌ی تحلیلگر دارند. به طور موازی، فراخوان بالا نیز منجر به تعداد زیادی مثبت کاذب<sup>۲</sup> می‌شود که باعث می‌شود تحلیلگر تعداد زیادی از نیازها را کاهش دهد و فراخوان را مبهم کند. به طور کلی نتایج گزارش شده نشان می‌دهند که تحلیلگر تمایل دارد که صحت<sup>۳</sup> و فراخوان را با هم در یک سطح نگه دارد.

### ۲.۱.۳ خلاصه سازی اطلاعات

از دیدگاه دیگر، متن کاوی در مهندسی نرم‌افزار برای اعتبارسنجی داده‌های لیست‌های پستی، کارنامه‌ی CVS و فایل‌های تغییر ورودی نرم‌افزار متن‌باز استفاده شده است. محققان مجموعه‌ای از ابزارها، یعنی *SoftChange*<sup>۴</sup>، را ایجاد کردند که اعتبارسنجی داده‌های منابع متنی نرم‌افزار متن‌باز ذکر شده در بالا را پیاده‌سازی می‌کنند. ابزارهای آن‌ها این انواع داده‌ها را از پروژه‌های متن‌باز، بازبینی، خلاصه‌سازی و اعتبارسنجی می‌کنند. بخشی از تجزیه و تحلیل آن‌ها می‌تواند فعال‌ترین برنامه‌نویسان یک پروژه متن‌باز را مشخص کند. آمار و دانش جمع‌آوری شده توسط تجزیه و تحلیل *SoftChange* به طور کامل مورد بهره‌برداری قرار نگرفته است، زیرا روش‌های پیش‌بینی دیگری می‌توانند با توجه به قطعات کدی که ممکن است در آینده تغییر کنند، یا تجزیه و تحلیل انجمنی بین اهمیت تغییرات و افراد، مورد استفاده قرار گیرند (یعنی آیا همه‌ی تغییرات ثبت شده توسط فعال‌ترین برنامه‌نویس، در مقیاس و در عمل به اندازه‌ی بقیه مهم بوده است؟).

### ۲.۳ تجزیه و تحلیل تولید

این بخش، یک مرور کلی از رویکردهای کاوش مورد استفاده برای کمک به روند تولید را فراهم می‌کند. ویژگی‌های اصلی این رویکردها در جدول ۲.۳ خلاصه شده است.

Recall<sup>۱</sup>

False negative<sup>۲</sup>

Precision<sup>۳</sup>

<http://sourcechange.sourceforge.net/><sup>۴</sup>

جدول ۲.۳: رویکرد های کاوش مورد استفاده در توسعه نرم افزار

نتایج تحلیل داده ها	داده ورودی	رویکرد کاوش
فرایندهای نرم افزار	کد منبع	خوشه بندی
ردیابی اشکال ها	SCM ، کد منبع	رده بندی، بازیابی متن
تصحیح نقص	نقص	کاوش الگوهای تکرار شونده
قوانین تلاش	داده	قوانین انجمنی
شرطهای نادیده گرفته شده	برنامه گراف وابستگی	قوانین انجمنی

### ۱.۲.۳ خوشه بندی

متن کاوی نیز در مهندسی نرم افزار برای کشف فرایندهای تولید استفاده شده است. فرایندهای نرم افزاری شامل رویدادهایی مانند روابط عوامل، ابزارها، منابع و فعالیت های سازماندهی شده توسط ساختارهای کنترل جریان است که بیان می دارد، این مجموعه از رویدادها به صورت سری، موازی، تکراری اجرا می شوند یا یکی از مجموعه ها به طور انتخابی اجرا می شود. کشف فرآیند نرم افزار، منابع تولید را (به عنوان مثال کد منبع، رونوشت های ارتباطات و غیره) به عنوان ورودی می گیرد و هدف آن کشف دنباله ای از رویدادهای مشخص کننده وظایفی است که منجر به تولید آنها شده است. محققان یک روش نوآورانه برای کشف فرایندهای نرم افزاری مخازن وب نرم افزارهای متن باز ارائه داده اند. روش آنها شامل تکنیک های استخراج متن، تحلیل موجودیت و تجزیه و تحلیل شبکه های اجتماعی است و بر اساس فرآیند رده بندی های موجودیت است. روش های خودکار اصلاح رده بندی با استفاده از وظایف متن کاوی می تواند مورد استفاده قرار گیرد، به طوری که این روش مستلزم وابستگی شدید به اقدامات، ابزار، منابع و عواملان رده بندی است. یک مثال می تواند استفاده از خوشه بندی متن روی منابع متنی نرم افزار باز و استخراج فقره های کاندیدای جدید برای رده بندی ناشی از برجسب های خوشه ای باشد.

آنها به عنوان متن ورودی، لیست های پستی توسعه دهنده Apache را استفاده کرده اند. تحلیل موجودیت ضروری است، زیرا بسیاری از افراد بیش از یک نام مستعار استفاده می کنند. پس از ساختن گراف اجتماعی که از

ارتباطات بین فرستنده و پاسخ‌دهنده ساخته شده است، آن‌ها یک تحلیل شبکه اجتماعی انجام دادند و به یافته‌های واقعا مهمی مانند رابطه‌ی قوی بین فعالیت ایمیل و فعالیت سطح کد منبع رسیدند. علاوه بر این، تجزیه و تحلیل شبکه‌های اجتماعی در این سطح، گره‌های مهم (افراد) را در بحث‌ها نشان می‌دهد. اگرچه تجزیه و تحلیل گراف و ارتباط در این روش دخیل بود، اما استفاده از تکنیک‌های رتبه‌بندی گره‌ها، مانند *PageRank* یا سایر تکنیک‌های پردازش گراف مانند *Spreading Activation* انجام نشدند.

### ۲.۲.۳ رده بندی

مخازن کد منبع، اطلاعاتی غنی را ذخیره می‌کنند که نه تنها برای مدیریت و ساخت کد منبع مفید است، بلکه یک گزارش دقیق نیز به دست می‌دهد که چگونه کد منبع در طول توسعه پیشرفت کرده است. اطلاعات مربوط به شواهدی از *refactoring* کد منبع در مخزن ذخیره می‌شود. همچنین همان‌طور که اشکال‌ها رفع می‌شوند، تغییرات انجام شده برای اصلاح مشکل ثبت می‌شوند. همان‌طور که *API* های جدید به کد منبع اضافه می‌شوند، یک راه مناسب استفاده از آن‌ها به صورت ضمنی در کد منبع، توضیح داده می‌شود. سپس، یکی از چالش‌ها، توسعه‌ی ابزارها و تکنیک‌هایی است که استخراج و استفاده از این اطلاعات مفید را به طور خودکار انجام می‌دهند.

در [۵۶] یک روش پیشنهاد شده است که از داده‌هایی استفاده می‌کند که رفع اشکال‌های کاوش شده از مخزن کد منبع برای بهبود تکنیک‌های تجزیه و تحلیل پایا برای پیدا کردن اشکال‌ها را توصیف می‌کند. این یک رویکرد دو مرحله‌ای است که از تاریخچه تغییر کد منبع یک پروژه نرم‌افزاری برای بهبود جستجوی اشکال‌ها استفاده می‌کند.

اولین گام در این فرآیند، شناسایی انواع اشکال‌هایی است که در نرم‌افزار رفع شده‌اند. هدف بررسی کردن داده‌های تاریخی ذخیره شده برای پروژه نرم‌افزاری است، تا بتوانیم بفهمیم که چه داده‌هایی وجود دارند و چگونه ممکن است در وظیفه‌ی یافتن اشکال‌ها مفید باشند. بسیاری از اشکال‌های موجود در تاریخچه‌ی *CVS*، نامزدهای خوبی برای شناسایی توسط تجزیه و تحلیل آماری، بررسی اشاره‌گر *NULL* و بررسی مقادیر بازگشتی تابع هستند. گام دوم ساختن یک یابنده‌ی اشکال به وسیله‌ی این یافته‌ها است. ایده

این است که یک بررسی کننده‌ی مقدار بازگشتی یک تابع، بر اساس این دانش که یک نوع خاص از اشکال چند بار در گذشته رفع شده است، توسعه یابد. به طور خلاصه، این بررسی کننده به دنبال مواردی می‌گردد که در آن‌ها مقدار بازگشتی یک تابع قبل از اینکه مورد آزمایش قرار بگیرد، در کد منبع استفاده شود. استفاده از یک مقدار بازگشتی می‌تواند به معنای پاس دادن آن به عنوان یک شناسه تابع به یک تابع، استفاده از آن به عنوان بخشی از یک محاسبه، ارزیابی مقدار آن اگر یک اشاره‌گر باشد یا مقداردهی دوباره‌ی آن قبل از آزمایش باشد. همچنین مواردی که مقادیر بازگشتی هیچ وقت توسط تابع فراخوانی ذخیره نمی‌شوند، بررسی می‌شوند. آزمایش کردن یک مقدار بازگشتی به معنی آن است که یک تصمیم کنترل جریان به مقدار بازگشتی اطمینان دارد.

بررسی کننده یک تحلیل جریان داده روی متغیر نگهداری مقدار بازگشتی انجام می‌دهد، فقط به منظور تعیین این که آیا این مقدار قبل از آزمایش مورد استفاده قرار گرفته است یا خیر. آن به سادگی، متغیر اصلی را که مقدار بازگشتی در آن ذخیره می‌شود، شناسایی می‌کند و استفاده بعدی از آن متغیر را تعیین می‌کند.

علاوه بر این، این بررسی کننده هشدارهایی را که در یکی از دسته‌های زیر می‌یابد، دسته‌بندی می‌کند:

- هشدارها برای مقادیر بازگشتی‌ای که به طور کامل نادیده گرفته می‌شوند، یا حالتی که مقدار بازگشتی ذخیره شده باشد اما هرگز استفاده نشده باشد.
- همچنین هشدارها برای مقادیر بازگشتی‌ای که در یک محاسبه استفاده می‌شوند قبل از این که در یک عبارت کنترل جریان مورد آزمایش قرار بگیرند.

هر مقدار بازگشتی‌ای که قبل از این که مورد آزمایش قرار گیرد، به عنوان یک شناسه تابع به یک تابع پاس داده می‌شود، و همچنین مقدار بازگشتی هر اشاره‌گری که بدون آزمایش مورد ارجاع قرار می‌گیرد، پرچم گذاری می‌شود. با این حال انواعی از توابع وجود دارند که باعث می‌شوند روند تجزیه و تحلیل پایا، هشدارهای مثبت کاذب تولید کنند. مشخص کردن این که

مقدار بازگشتی کدام تابع نیاز نیست بررسی شود، بدون دانش قبلی، سخت است. تکنیک‌های کاوش مخزن کد منبع می‌توانند باعث بهبود نتایج تجزیه و تحلیل پایا باشند. به طور خاص، داده‌ی بدست آمده از مخزن کد منبع و از نسخه فعلی نرم‌افزار که ما کاوش می‌کنیم، برای تعیین الگوی استفاده‌ی واقعی برای هر تابع استفاده می‌شود.

به طور کلی، مشاهده شده است که اشکال‌های فهرست شده در پایگاه داده‌های اشکال و آن‌هایی که توسط بررسی تاریخچه‌ی تغییرات کد منبع یافت شده‌اند، در نوع و سطح انتزاع متفاوت هستند. مخازن نرم‌افزار، همه‌ی اشکال‌های رفع شده در هر مرحله از فرآیند تولید را ثبت می‌کنند و در نتیجه آن‌ها اطلاعات بسیار مفیدی را ارائه می‌دهند. بنابراین، ثابت شده است که یک سیستم برای تکنیک‌های تشخیص اشکال، زمانی که داده‌های بدست آمده از مخزن کد منبع را به صورت خودکار کاوش می‌کند، موثرتر است.

### ۳.۲.۳ قوانین انجمنی و کاوش الگوهای تکرار شونده

یک رویکرد در [۴۴] پیشنهاد شده است که از تکنیک‌های استخراج قوانین انجمنی برای تحلیل داده‌های نقص استفاده می‌کند. نقص‌های نرم‌افزاری شامل اشکال‌ها، مشخصات و تغییرات طراحی است. داده‌ی نقص جمع‌آوری شده تحت تجزیه و تحلیل، متغیرهای مقیاس اسمی هستند مانند شرحی از نقص، اولویت برای رفع نقص و وضعیت آن و همچنین متغیر مقیاس فاصله و نسبت، با توجه به تلاش و طول مدت تصحیح نقص. یک روش کاوش قوانین انجمنی گسترش یافته برای استخراج اطلاعات مفید و آشکار کردن قوانین مرتبط با تلاش اصلاح نقص، به کار گرفته شده است.

محققان پیشنهاد استفاده از تکنیک‌های کاوش گراف برای کشف قوانین شرطی ضمنی در یک پایه کد و برای کشف قوانین نقضی که نشان دهنده‌ی شرط نادیده گرفته شده است، را دادند. آن‌ها برنامه‌ها و قوانین برنامه نویسی شرطی را در قالب گراف‌های وابستگی نمایش دادند. سپس آن‌ها از تکنیک‌های کاوش زیرگراف‌های تکرار شونده و کاوش مجموعه فقره‌های تکرار شونده برای کشف قوانین شرطی شامل پیش شرط‌ها و پس شرط‌های فراخوانی‌های تابع و همچنین کشف نقض‌های آن قوانین استفاده می‌کنند.

جدول ۳.۳: روش های کاوش استفاده شده در تست نرم افزار

نتایج تحلیل داده ها	داده ورودی	رویکرد کاوش
شبکه آزمایش کننده تابع	متغیر های I/O سیستم نرم افزار	رده بندی
خوشه های پرو فایل های اجرا	پرو فایل های اجرا	خوشه بندی
رده بندی های رفتار نرم افزار	اجرا های برنامه	خوشه بندی، رده بندی

## ۳.۳ آزمایش کردن

ارزیابی نرم افزار بر اساس آزمایش هایی است که توسط آزمایش کنندگان نرم افزار طراحی شده است. بنابراین ارزیابی خروجی های آزمایش با تلاش های قابل توجه آزمایش کنندگان انسانی مرتبط است، که اغلب دانش ناقصی از مشخصات نیازها دارند.

رویکردهای داده کاوی را می توان برای استخراج اطلاعات مفید از نرم افزار مورد آزمایش، به کار برد که می تواند به تست نرم افزار کمک کند. به طور خاص، مدل های داده کاوی ایجاد شده از نرم افزار آزمایش شده، می تواند برای بازیابی مقادیر گم شده و مشخصات ناقص، طراحی مجموعه ای از آزمون های رگرسیون و ارزیابی درست بودن خروجی های نرم افزار در هنگام منتشر کردن نسخه های جدید سیستم استفاده شود. یک کتابخانه تست رگرسیون باید شامل یک حداقل تعداد تستی باشد که تمام جنبه های ممکن عملکرد سیستم را، پوشش می دهد. برای اطمینان از طراحی موثر تست های رگرسیون جدید، باید نیازهای واقعی یک سیستم موجود را پوشش دهد.

بنابراین، یک آزمایش کننده باید مشخصات سیستم را تجزیه و تحلیل کند، تحلیل ساختاری کد منبع سیستم را انجام دهد و نتایج اجرای سیستم را برای تعریف روابط ورودی خروجی در نرم افزار مورد آزمایش، مشاهده کند. جدول ۳.۳ تکنیک های اصلی داده کاوی را که در زمینه تست نرم افزار مورد استفاده قرار می گیرند خلاصه می کند.



### ۱.۳.۳ خوشه بندی

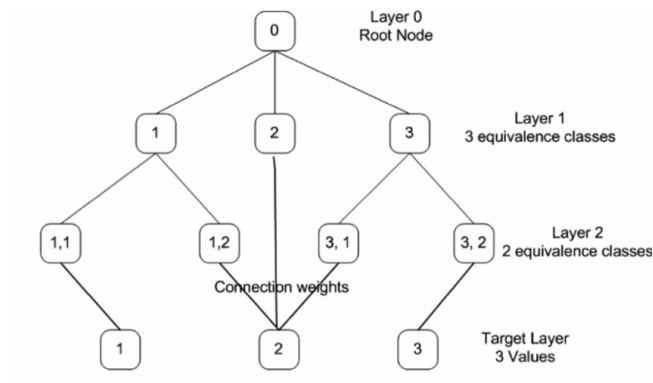
یک رویکرد دیگر، که روش‌های تحلیل خوشه را برای انتخاب مجموعه‌ای از اجراها که برای تبعیت از نیازها مورد ارزیابی قرار می‌گیرند، استخراج می‌کند. رویکرد پیشنهادی از یک مجموعه‌ی پروفایل‌های اجرایی<sup>۵</sup> استفاده می‌کند که اینگونه تعریف می‌شوند، اجرای نسخه‌ی تحت آزمایش نرم‌افزار روی یک مجموعه‌ی داده شده از ورودی‌های برنامه. الگوریتم خوشه‌بندی برای فیلتر کردن پروفایل‌ها بر اساس ویژگی‌های مشابهشان استفاده می‌شود. سپس پروفایل‌های اجرایی از نتایج خوشه‌ها انتخاب می‌شوند.

رویکردی که هدف آن تجزیه و تحلیل مجموعه‌ای از اجراهای برنامه‌ها و تعریف رده‌بندی‌های رفتار نرم‌افزار است، در [۸] پیشنهاد شده است. با توجه به این مطالعه، مدل مارکف برای کدگذاری پروفایل‌های پروژه‌ها مورد استفاده قرار می‌گیرد. سپس مدل‌های مارکف اجراهای برنامه‌های فردی، با استفاده از الگوریتم خوشه‌بندی زنجیره‌ای، خوشه‌بندی شده است. روند خوشه‌بندی قصد دارد اجراهای مشابه برنامه را جمع کند و سپس رده‌بندی‌های موثر رفتار برنامه را تعریف کند. همچنین خودراه‌اندازی<sup>۶</sup> به عنوان یک تکنیک یادگیری فعال استفاده می‌شود، در نتیجه رده‌بندی‌های یادگیری به طور مداوم تحت آموزش قرار می‌گیرند. به طور خاص، این روش با شناسایی نمونه‌های آموزشی جدید به رده‌بندی‌ها کمک می‌کند و سپس رده‌بندی‌ها را می‌توان با استفاده از مجموعه‌ی گسترده‌ای از نمونه‌های آموزشی، بازآموزی کرد.

### ۲.۳.۳ رده بندی

یک رویکرد که هدف آن خودکارسازی تجزیه و تحلیل ورودی خروجی داده‌های اجرا بر اساس یک متدولوژی داده‌کاوی است، که در [۳۵] پیشنهاد شده است. این متدولوژی مبتنی بر شبکه اطلاعات فازی<sup>۷</sup> است که دارای ساختار درخت‌مانند بی نظیری است. اجزای شبکه شامل گره‌ی ریشه، یک تعداد قابل تغییری از لایه‌های پنهان (یک لایه برای هر ورودی انتخاب شده) و

execution profiles<sup>۵</sup>  
Bootstrapping<sup>۶</sup>  
info-fuzzy network (IFN)<sup>۷</sup>



شکل ۱.۳: یک مثال از ساختار شبکه اطلاعات فازی

لایه هدف (خروجی) که نشان‌دهنده مقادیر خروجی ممکن است. یک ویژگی ورودی مشابه در تمام گره‌های یک لایه (سطح) مورد استفاده قرار می‌گیرد در حالی که هر گرهی هدف با یک مقدار (رده) در دامنه ویژگی هدف مرتبط است. اگر مدل شبکه اطلاعات فازی به منظور پیش‌بینی مقادیر یک ویژگی هدف پیوسته استفاده شده باشد، گره‌های هدف، بازه‌های از هم جدا (بدون اشتراک) در محدوده ویژگی را نشان می‌دهند.

یک لایه پنهان  $l$  شامل گره‌هایی است که نشان‌دهنده ترکیبی از مقادیر  $l$  ویژگی ورودی اول هستند که مشابه تعریف یک گرهی داخلی در یک درخت تصمیم استاندارد است. گره‌های نهایی (پایانه) شبکه نشان‌دهنده ارتباطات غیرافزوده از مقادیر ورودی هستند که خروجی‌های متمایزی تولید می‌کنند. با توجه به اینکه شبکه از داده‌های اجرایی یک سیستم نرم‌افزاری استنباط می‌شود، هر اتصال درونی بین گرهی پایانه و هدف، یک خروجی احتمالی یک نمونه آزمایشی را نشان می‌دهد. شکل ۱.۳ ساختار شبکه اطلاعات فازی را نشان می‌دهد.

یک شبکه اطلاعات فازی جداگانه برای نشان دادن هر متغیر خروجی ساخته شده است. مولفه‌های اصلی محیط مبتنی بر شبکه اطلاعات فازی ارائه شده در [۳۵] عبارتند از:

- سیستم میراث<sup>۸</sup>. این مولفه یک برنامه، یک جزء یا یک سیستم مورد

Legacy (LS)system<sup>۸</sup>

آزمایش را در نسخه‌های بعدی نرم‌افزار نشان می‌دهد.

- مشخصات ورودی‌ها و خروجی‌های برنامه<sup>۹</sup>. داده‌های پایه برای هر متغیر ورودی و خروجی در سیستم میراث.
- مولد تست تصادفی<sup>۱۰</sup>. این مولفه، ترکیباتی تصادفی از مقادیر، در محدوده هر متغیر ورودی تولید می‌کند.
- بستر تست<sup>۱۱</sup>. این مولفه موارد آموزشی ایجاد شده توسط مولفه مولد تست تصادفی را به سیستم میراث می‌فرستد.

الگوریتم شبکه اطلاعات فازی بر روی ورودی‌های ارائه شده توسط مولد تست تصادفی و خروجی‌های حاصل از یک سیستم میراث با استفاده از مولفه بستر تست مورد آموزش قرار می‌گیرد. یک مولفه شبکه اطلاعات فازی جداگانه برای هر متغیر خروجی ساخته شده است. الگوریتم شبکه اطلاعات فازی، نمونه‌های آموزشی را که به طور تصادفی از طریق مولفه مولد تست تصادفی تولید می‌شوند، به عنوان ورودی در نظر می‌گیرد و خروجی‌ها توسط سیستم میراث برای هر نمونه آزمایشی تولید می‌شوند. الگوریتم شبکه اطلاعات فازی بطور مکرر اجرا می‌شود تا زیرمجموعه‌ی متغیرهای ورودی مربوط به هر خروجی و مجموعه‌ی متناظر نمونه‌های آزمایشی غیرافزوده را پیدا کند. نمونه‌های آزمایشی واقعی با استفاده از یک سیاست آزمایشی موجود، از رده‌های هم‌ارز به طور خودکار تشخیص داده شده، تولید می‌شوند.

### ۴.۳ اشکال زدایی

برآورد اشکال در نرم‌افزار فعالیتی بسیار اساسی برای برنامه ریزی موفق و صحیح پروژه است. تمام داده‌های مربوط به اشکال نرم‌افزار در مخزن اشکال نگهداری می‌شوند. مخزن اشکال نرم‌افزار در برگیرنده اطلاعات جالبی در زمینه تولید یک پروژه است و به کارگیری داده‌کاوی در این مخزن‌ها می‌تواند الگوهای جالبی را نشان دهد. برای مثال به کارگیری یک روش داده‌کاوی

Specification of Application Inputs and Outputs (SAIO)<sup>۹</sup>  
Random test generator (RTG)<sup>۱۰</sup>  
Test Bed (TB)<sup>۱۱</sup>

جدول ۴.۳: روش های کاوش استفاده شده در اشکال زدایی

نتایج تحلیل داده ها	داده ورودی	رویکرد کاوش
تشخیص خطاهای منطقی	تست های ورودی خروجی	رده بندی احتمالاتی
اشکال های منطقی	اجراهای برنامه	رده بند <i>SVM</i>
درخت تصمیم شکست ها	پروفایل های اجرایی و نتیجه	رده بندی، درخت تصمیم
الگوهای استفاده فراخوانی	کد منبع	قوانین انجمنی

پیش بینی، می تواند برای تشخیص اشکال نرم افزار در مخزن اشکال به کار رود. در ابتدای کار خلاصه و توضیحات اشکالی که نیاز به برآورد دارد با خلاصه و توضیحات اشکال های موجود در مخزن اشکال مقایسه و برای این کار از مدل تشابه وزنی استفاده می شود. گام بعدی محاسبه زمان مورد نیاز برای اصلاح تمام اشکال های مشابه و گرفتن میانگین آن جهت برآورد پیش بینی شده یک اشکال است [۶۰].

خطاهای منطقی برنامه به ندرت موجب تناقض های دسترسی به حافظه می شوند اما خروجی های نادرست تولید می کنند. تعدادی از تکنیک های کاوش برای شناسایی خطاهای منطقی و کمک به اشکال زدایی نرم افزار، مورد استفاده قرار گرفته اند (جدول ۴.۳ را ببینید).

### ۱.۴.۳ رده بندی شکست های نرم افزار

یک استراتژی نیمه خودکار برای رده بندی شکست های نرم افزاری در [۴۶] ارائه شده است. این رویکرد مبتنی بر این ایده است که اگر  $m$  شکست در طی یک دوره زمانی که نرم افزار اجرا می شود، مشاهده شود، احتمال دارد که این خرابی ها به دلیل تعداد نقص های متمایز کمتری باشد. فرض کنید  $F = \{f_1, f_2, \dots, f_m\}$  مجموعه ای از شکست های گزارش شده است و هر شکست تنها به دلیل یک نقص است. سپس  $F$  را می توان به  $m > k$  زیرمجموعه های  $F_1, F_2, \dots, F_k$  تقسیم کرد به طوری که تمام شکست های  $F_i$  توسط یک نقص  $d_i$  مشابه ایجاد شده اند که  $k \geq i \geq 1$ . این بخش بندی، رده بندی شکست واقعی نامیده می شود. در ادامه، ما فازهای اصلی استراتژی را برای تقریب زدن رده بندی شکست واقعی توصیف می کنیم:

۱. نرم افزار برای جمع آوری و انتقال به توسعه ی پروفایل های اجرایی یا

اجراهای گرفته شده پیاده‌سازی شده است، و سپس به کار گرفته می‌شود.

۲. پروفایل‌های اجرایی متناظر با شکست گزارش شده، با یک نمونه تصادفی از پروفایل‌های اجرایی عملیاتی‌ای که برای آن‌ها هیچ شکستی گزارش نشده است، ترکیب شده‌اند. این مجموعه‌ی پروفایل‌ها برای انتخاب یک زیرمجموعه از تمام ویژگی‌های<sup>۱۲</sup> پروفایل برای استفاده در گروه‌بندی شکست‌های مرتبط، مورد تجزیه و تحلیل قرار گرفته است. یک ویژگی از یک پروفایل اجرایی، مربوط به یک خاصیت<sup>۱۳</sup> یا عنصر از آن است. به عنوان مثال، یک پروفایل فراخوانی تابع شامل یک شمارش اجرا برای هر تابع در یک برنامه است و هر شمارش یک ویژگی از پروفایل است. پس استراتژی انتخاب ویژگی به شرح زیر است:

- مجموعه ویژگی‌های کاندید را تولید کنید و از هر یک برای ایجاد و آموزش یک رده‌بند الگو برای تشخیص تمایز بین شکست‌ها از اجراهای موفق، استفاده کنید.
- ویژگی‌های رده‌بندی که بهترین نتایج را به دست می‌دهد، را انتخاب کنید.

۳. پروفایل‌های شکست‌های گزارش شده با استفاده از تجزیه و تحلیل خوشه‌ای، تجزیه و تحلیل می‌شوند تا شکست‌هایی که پروفایل‌هایشان مشابه است با توجه به ویژگی‌های انتخاب شده در فاز ۲، گروه شوند.

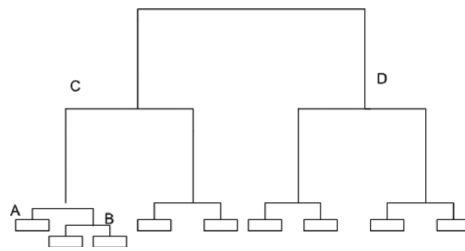
۴. رده‌بندی شکست حاصل در گروه‌ها به منظور تأیید یا اصلاح آن، مورد بررسی قرار می‌گیرد.

استراتژی توصیف شده فوق، رده‌بندی اولیه‌ای از شکست‌های نرم‌افزاری را ارائه می‌دهد. بسته به کاربرد و نیازهای کاربر، این رده‌های اولیه می‌توانند ادغام شوند یا تقسیم شوند تا شکست‌های نرم‌افزاری، به طور مناسبی شناسایی شوند.

در [۱۹]، دو تکنیک مبتنی بر درخت، برای اصلاح رده‌بندی اولیه شکست‌ها پیشنهاد شده است. در زیر ما ایده اصلی این رویکردها را ارائه می‌دهیم.

---

feature<sup>۱۲</sup>  
attribute<sup>۱۳</sup>



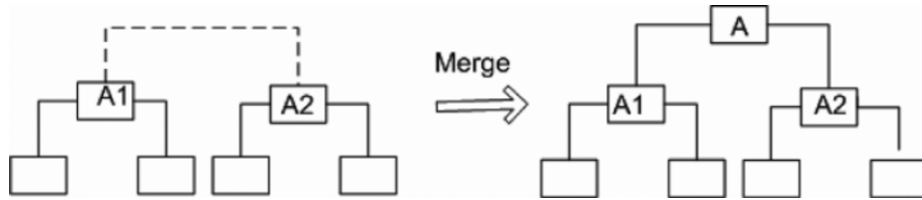
شکل ۲.۳: سلسله مراتب خوشه ها

### ۲.۴.۳ اصلاح رده‌بندی شکست‌ها با استفاده از دندروگرام‌ها

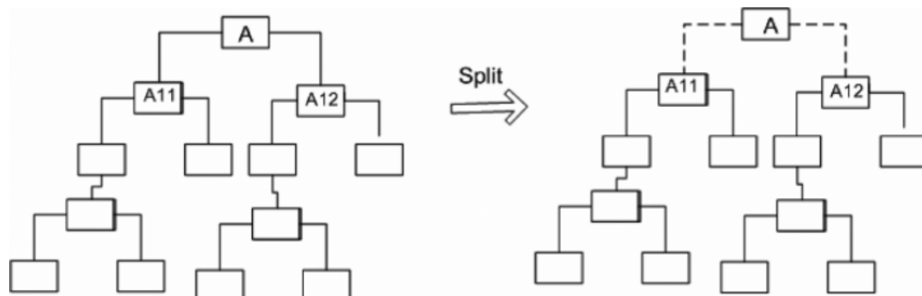
یکی از استراتژی‌هایی که برای اصلاح دسته‌بندی اولیه شکست پیشنهاد شده است، براساس نموداری درخت‌مانند (شناخته شده به عنوان دندروگرام<sup>۱۴</sup>) است. به طور خاص، از آن‌ها استفاده می‌کنند تا تصمیم بگیرند که چگونه خوشه‌های غیر یکپارچه باید به دو یا چند زیر مجموعه تقسیم شوند و تصمیم بگیرند که کدام خوشه‌ها برای ادغام باید در نظر گرفته شوند. یک خوشه در یک دندروگرام به یک زیردرخت نظیر می‌شود که نشان‌دهنده‌ی ارتباط‌های بین زیرخوشه‌های آن است. هر چه جد مشترک دو خوشه، در فاصله‌ی دورتری از ریشه‌ی دندروگرام قرار داشته باشد، آن دو خوشه بهم شبیه‌ترند. به عنوان مثال، بر اساس دندروگرام ارائه شده در شکل ۲.۳ می‌توان دید که خوشه‌های A و B بیشتر از خوشه‌های C و D مشابه هستند. بزرگترین زیردرخت همگن یک خوشه، بزرگترین زیردرختی است که شامل شکست‌هایی با یک علت مشابه است. اگر خوشه‌بندی بیش از حد درشت باشد، ممکن است برخی از خوشه‌ها دارای دو یا چند زیر درخت همگن بزرگ باشند که شامل شکست‌هایی با علل مختلف هستند. چنین خوشه‌ای باید در سطحی که در آن زیردرخت‌های همگن بزرگ آن، متصل هستند، تقسیم شود، به طوری که این زیردرخت‌ها به دو خواهر تبدیل می‌شوند همانطور که شکل ۳.۳ نشان می‌دهد. اگر خیلی خوب باشد، ممکن است دو خواهر خوشه‌هایی با شکست‌هایی با علل مشابه باشند. چنین خواهرانی (خوشه‌ها) باید در سطح والدینشان ادغام شوند (شکل ۴.۳).

بر اساس این تعاریف، استراتژی پیشنهاد شده برای اصلاح رده‌بندی اولیه

<sup>۱۴</sup>dendrogram



شکل ۳.۳: ادغام دو خوشه. خوشه جدید  $A$  حاوی خوشه هایی است که توسط دو زیردرخت همگن  $A1$  و  $A2$  معرفی شده اند.



شکل ۴.۳: تقسیم یک خوشه: دو خوشه جدید (زیر درخت های با ریشه های  $A11$  و  $A12$ ) به زیردرخت های بزرگ همگنی در خوشه قدیمی مربوط می شوند.

شکست‌ها با استفاده از دندروگرام‌ها دارای سه مرحله است:

۱. تعداد خوشه‌هایی که دندروگرام باید به آن تقسیم شود، را انتخاب کنید.
۲. خوشه‌های فردی را برای همگن بودن، با انتخاب دو اجرایی که در خوشه دارای پروفایل‌هایی با بیشترین بی‌شبهی هستند، بررسی کنید. اگر اجراهای انتخاب شده دارای علل مشابه یا مرتبط هستند، احتمال دارد که تمام شکست‌های دیگر در خوشه نیز چنین باشند. اگر اجراهای انتخاب شده علت‌های مشابه یا مرتبط ندارند، خوشه همگن نیست و باید تقسیم شود.
۳. اگر نه خوشه و نه خواهر آن براساس مرحله ۲ تقسیم نشدند و شکست‌های مورد بررسی یک علت داشتند، آنها را ادغام می‌کنیم. خوشه‌هایی که توسط عملیات ادغام یا تقسیم تولید شده‌اند، باید به همان شیوه تجزیه و تحلیل شوند، که امکان تقسیم یا ادغام بازگشتی را بدهند.

### ۳.۴.۳ پالایش با استفاده از درختان رده بندی

تکنیک دوم پیشنهاد شده توسط فرانسیس و همکاران، متکی بر ساخت یک درخت رده‌بندی برای تشخیص اجراهای ناموفق است. یک درخت رده‌بندی، نوعی رده‌بند الگویی است که به شکل یک درخت تصمیم‌گیری دودویی است. هر گره داخلی در درخت با یک عبارت رابطه‌ای برچسب گذاری می‌شود که یک ویژگی عددی از شی‌ای که رده‌بندی شده است، را با یک مقدار تقسیم ثابت، مقایسه می‌کند. از سوی دیگر، هر برگ از درخت با مقدار پیش‌بینی شده‌ای برچسب گذاری می‌شود که رده‌ای است که برگ نشان‌دهنده آن است. اگر درخت رده‌بندی داده شده باشد، ما باید درخت را از ریشه به برگ پیمایش کنیم، تا یک شی را رده‌بندی کنیم. در هر گام از پیمایش قبل از رسیدن به برگ، مقدار عبارت گره فعلی را ارزیابی می‌کنیم. هنگامی که شی به یک برگ می‌رسد، مقدار پیش‌بینی شده‌ی آن برگ به عنوان رده پیش‌بینی شده برای این شی، در نظر گرفته می‌شود. در مورد مساله‌ی رده‌بندی شکست نرم‌افزار، ما دو رده موفقیت و شکست



را در نظر می‌گیریم. الگوریتم‌های رده‌بندی و درخت رگرسیون<sup>۱۵</sup> برای ساخت درخت رده‌بندی مربوط به شکست‌های نرم‌افزار، مورد استفاده قرار گرفت. فرض کنید یک مجموعه آموزشی از پروفایل‌های اجرا به شما داده شده است

$$L = \{(x_1, j_1), \dots, (x_N, j_N)\}$$

که هر  $x_i$  یک پروفایل اجرا را نشان می‌دهد و  $j_i$  نتیجه (موفقیت / شکست) مربوط به آن است. مراحل ساخت درخت رده‌بندی بر اساس  $L$  عبارتند از:

- انحراف گره  $t \in L$  به صورت زیر تعریف می‌شود
- $$d(t) = \frac{1}{N_t} \sum (j_i - \bar{J}(t))^2 \quad (۱.۳)$$

که  $N_t$  اندازه  $t$  است و  $\bar{J}(t)$  مقدار متوسط  $j$  در  $t$  است.

- هر گره  $t$  به دو فرزند  $t_L$  و  $t_R$  تقسیم می‌شود. تقسیم‌بندی‌ای انتخاب می‌شود که دارای بیشترین کاهش انحراف باشد. یعنی، از مجموعه‌ی تقسیم‌های ممکن  $S$ ، تقسیم بهینه توسط فرمول زیر به دست می‌آید:

$$s^* = \operatorname{argmin}_s (d(t) - \frac{N_{t_L}}{N_t} d(t_R) - \frac{N_{t_R}}{N_t} d(t_L)) \quad (۲.۳)$$

- اگر  $\alpha \geq d(t)$  برای یک آستانه‌ی  $\alpha$  باشد، یک گره برگ اعلام می‌شود.
- مقدار پیش‌بینی شده برای یک برگ، میانگین مقدار  $j$  در بین اجراهای آن برگ است.

### ۴.۴.۳ قوانین انجمنی و کاوش الگوهای تکرارشونده

دو رویکرد برای کاوش الگوهای استفاده از فراخوانی از کد منبع وجود دارد. اولین رویکرد، بر مبنای ایده‌ی کاوش مجموعه فقره‌ها است. این روش، زیر مجموعه‌های فقره‌های تکرارشونده را شناسایی می‌کند که حداقل، یک کمترین مقدار پشتیبانی تعریف شده توسط کاربر را برآورده می‌کند. نتایج استفاده از این رویکرد برای کد منبع، الگوهای نامرتب و مرتبط با فراخوانی‌های تابع است. از سوی دیگر، رویکرد استخراج الگوی سری،

<sup>۱۵</sup> Classification And Regression Tree (CART)

مجموعه‌ای از الگوهای مرتب با حداقل میزان پشتیبانی برآورده شده، تولید می‌کند. به طور کلی، این رویکردها می‌توانند به الگوهای کاوش استفاده از فراخوانی و در نتیجه شناسایی بالقوه اشکال‌ها در یک سیستم نرم‌افزاری کمک کنند.

## ۵.۳ نگه داری

یک مشکل که در مهندسی نرم‌افزار با آن مواجه می‌شویم نگه‌داری صحیح نرم‌افزار است. به دست آوردن و شناسایی نقص‌های سیستم قبل از آن که منجر به خطا شوند، برای ما مطلوب است. به نظر می‌رسد خیلی از خطاها به گروه‌هایی تقسیم می‌شوند که هر گروه از نقص خاصی در سیستم ناشی شده‌اند.

تحقیقات اخیر، بر روش‌های داده‌کاوی به منظور آسان کردن دسته‌بندی خطاها بر اساس علتشان تمرکز دارد. به طور خاص، این رویکردها نیاز به سه دسته اطلاعات زیر دارند:

۱. مشخصات اجرایی که نشان دهنده علل خطای به وجود آمده هستند.
  ۲. اطلاعات حسابرسی که وجود خطا را تایید می‌کنند.
  ۳. اطلاعات تشخیصی که در تعیین علت بروز خطا استفاده می‌شوند.
- در ادامه روش‌های داده‌کاوی که نگه‌داری نرم‌افزار را راحت تر می‌کنند، آمده است.

### ۱.۵.۳ خوشه بندی

به طور خاص، روش‌های خوشه‌بندی به مهندسان کمک می‌کنند تا درک بهتری از ساختار کد منبع و ارزیابی قابلیت نگه‌داری آن داشته باشند. رویکرد ارائه شده به زیر مجموعه‌هایی از کد منبع شامل موارد زیر اعمال شده است:

۱. ماهیت‌هایی که متعلق به دامنه‌های ساختاری و یا رفتاری هستند.
۲. ویژگی‌هایی که ماهیت‌های گفته شده را توصیف می‌کنند (مانند نام رده، تابع و ...)

۳. معیارهای استفاده شده به عنوان ویژگی اضافه که به برنامه‌نویسان نرم‌افزار در درک سیستم تحت نگهداری کمک می‌کنند.

بخش دیگر چارچوب فرایند استخراج ویژگی‌های کد منبع است سپس اطلاعات استخراج شده در یک پایگاه داده ذخیره می‌شوند تا روش‌های داده‌کاوی روی آن‌ها اعمال شوند. در این رویکرد خاص، روش‌های خوشه‌بندی برای تحلیل اطلاعات ورودی و فراهم کردن درک نسبی از سیستم نرم‌افزاری استفاده می‌شوند. خوشه‌بندی با ایجاد دسته‌های مستقل از هم با توجه به شباهت‌هایشان، درکی کلی از سیستم نرم‌افزاری ارائه می‌دهد. علاوه بر این، خوشه‌بندی می‌تواند به کشف الگوهای برنامه‌نویسی و داده‌های پرت و موارد غیرمعمول کمک کند.

### ۲.۵.۳ رده بندی

در [۴۹] آن‌ها از داده به دست آمده از یک میلیون پروژه نرم‌افزاری متن باز موجود در پرتال *SourceForge* به منظور ساختن مدل پیش‌بینی‌کننده برای نگهداری نرم‌افزار با استفاده از روش‌های داده‌کاوی و متن‌کاوی استفاده کردند.

با استفاده از *SAS Enterprise Miner* و *SAS Text Miner*، آن‌ها روی جمع‌آوری مقادیر برای متغیرها با توجه به هزینه‌های نگهداری و میانگین زمان بازیابی خطا تمرکز کردند. این کار همچنین حذف پروژه‌های توسعه نیافته و یا پروژه‌های بدون گزارش خطا را به دنبال داشت و تنها پروژه‌های قابل استفاده و شامل مقادیر مورد نیاز باقی ماندند.

به علاوه، آن‌ها پروژه‌های باقی‌مانده را بر اساس توصیفشان خوشه‌بندی کردند تا بتوانند رده‌های مهم در پروژه‌های متن باز موجود در پایگاه داده *SourceForge* را کشف کنند. سپس مقادیر میانگین زمان‌های بازیابی خطا با استفاده از توزیع مقادیرشان تبدیل به شکل دودویی (بالا یا پایین) شدند. و در نهایت از *SAS Enterprise Miner* برای رده‌بندی میانگین زمان‌های بازیابی خطا استفاده کردند.

نتایج گزارش شده، همبستگی جالبی را بین ویژگی‌هایی مانند تعداد دانلود، قدمت پروژه و استفاده از ایمیل با متغیر رده نشان می‌داد. به عنوان مثال پروژه‌های با قدمت بالا، میانگین زمان‌های بازیابی خطای بالاتری نسبت به پروژه‌های جدید داشتند.

رویکردی برای بهره‌برداری از ایده فیلتر کردن اسپم‌ها<sup>۱۶</sup> برای تشخیص نرم‌افزارهای آسیب‌پذیر وجود دارد. این روش بر اساس این حقیقت است که واحدهای نرم‌افزاری معیوب، الگوی کلمات و جمله‌های مشابهی دارند. میزونی و همکارانش پیاده‌سازی ابزاری برای جدا کردن واحدهای معیوب (FP) و سالم (NFP) ارائه دادند. سپس این واحدها در مرحله یادگیری دسته‌بندی استفاده می‌شدند که بتوانند واحدهای معیوب یا سالم جدید را تشخیص بدهند.

همچنین رویکردی دیگر برای رده‌بندی داده‌های بزرگ برای فهم اساس و پایه آن‌ها در [۲۵] ارائه شد. هر چند این داده‌ها معمولاً به عنوان داده پرت در نظر گرفته می‌شوند، ممکن است اطلاعات مفیدی راجع به پروژه و بهبود آن داشته باشند.

هیندل و همکارانش تصمیم به استفاده از روش‌های رده‌بندی برای مشخص کردن رده‌های مختلف داده‌های بزرگ و در نتیجه تشخیص تغییرات مختلف نرم‌افزار گرفتند. این مطالعه نشان داد که در بیشتر موارد، داده‌های بزرگ نشان‌دهنده تغییرات در معماری سیستم هستند.

روشی برای ارزیابی میزان پایداری یک مدل پیش‌بینی‌کننده، ارائه دادند. آن‌ها ۴ پروژه متن‌باز را بررسی کرده و ویژگی‌هایشان را استخراج کردند. سپس یک مدل پیش‌بینی‌نقص با استفاده از درخت تصمیم نرم‌افزار Weka ساخته شد و محققان کیفیت پیش‌بینی را در طول زمان ارزیابی کردند. این مطالعه نشان داد که در طی زمان، تغییرات چشم‌گیری وجود داشته و در نتیجه باید محتاطانه استفاده شود.

### ۳.۵.۳ کاوش الگوهای تکرار شونده و قوانین انجمنی

محققان در مطالعه‌هایشان بهره‌برداری از روش‌های استخراج قوانین انجمنی به منظور تشخیص تغییرات همزمان در سیستم نرم‌افزاری را نشان می‌دهند. به عنوان مثال هدف ما کشف رابطه بین تغییر ماهیت‌های نرم‌افزار است. می‌خواهیم به این سوال پاسخ بدهیم که زمانی که یک ماهیت از کد منبع (مانند تابع A) تغییر می‌کند، چه ماهیت‌های دیگری ممکن است تغییر کنند (مانند توابع B و C)؟

---

<sup>۱۶</sup>Spam-Filtering

به طور خاص، ابزاری برای تجزیه کد منبع ارائه شده که شماره خطوط را به سطوح فیزیکی ماهیت‌ها مرتبط می‌کند. این ماهیت‌ها به صورت سه‌تایی (نام فایل، نوع، شماره) نمایش داده می‌شوند. تغییرات بعدی در مخزن به عنوان تراکنش گروه‌بندی شده‌اند. یک روش کاوش قوانین انجمنی سپس برای تشخیص قوانین به شکل

$B, C \rightarrow A$

استفاده می‌شود.

استفاده از روش‌های خوشه‌بندی و قوانین انجمنی برای بازیابی طراحی معماری سیستم‌های نرم‌افزاری قدیمی، با توجه به اهداف تعیین شده کاربر پیشنهاد شده است. کد منبع سیستم قدیمی تحلیل شده و مجموعه‌ای از بخش‌های پر تکرار آن استخراج می‌شود. با استفاده از روش‌های خوشه‌بندی و تطبیق الگو، الگوریتم ارائه شده بخش‌های مختلف سیستم قدیمی را تعریف می‌کند. بر اساس پرس و جوی کاربر، بهترین بخش مطابق با خواسته او انتخاب می‌شود. همچنین برای هر جواب ممکن می‌توان امتیازی در نظر گرفت که در نهایت لیستی رتبه‌بندی شده را برای ارزیابی به کاربر ارائه می‌دهد.

رویکردی برای شناسایی الگوهای استفاده مجدد کتابخانه‌ها نیز، یک روش پیشنهادی دیگر است. این رویکرد پیشنهادی از قوانین انجمنی برای شناسایی روابط بین رده‌ها در یک کتابخانه بهره می‌برد. با تعمیم قوانین انجمنی به قوانین عمومی، وراثت رده‌ها نیز در نظر گرفته می‌شود. بنابراین یک روش خودکار به منظور کشف الگوهای استفاده مجدد کتابخانه‌ها و شناسایی مشخصه استفاده آن‌ها توسعه می‌یابد.

یک رویکرد برای تحلیل کد نمونه برای پیدا کردن تغییرات استفاده در چارچوب نیز ارائه شده است. فرایند استخراج دو نسخه از کد نمونه را به عنوان ورودی می‌گیرد و هدف از بهره‌برداری از روش‌های کاوش الگوهای تکرار شونده توصیف تغییرات استفاده از چارچوب است. در مرحله اول، اطلاعاتی راجع به این که کد نمونه چگونه از چارچوب استفاده می‌کند ( رده‌های استفاده شده و یا توابع فراخوانی شده) استخراج می‌شود. سپس تراکنش‌ها به وسیله ترکیب اطلاعات نحوه استفاده به دست آمده از دو نسخه کد نمونه ساخته می‌شوند. در نهایت یک الگوریتم استخراج قوانین انجمنی بر تراکنش‌ها اعمال شده و تمام قوانین تغییر ممکن استخراج می‌شود.

## ۶.۳ استفاده مجدد از نرم افزار

یکی از مهم ترین جنبه های افزایش بهره وری و کیفیت نرم افزار، استفاده مجدد از آن است. [۴۰، ۴۵] اگر چه استفاده مجدد از نرم افزار می تواند شکل های مختلفی داشته باشد و موتور های جستجو برای بخش های نرم افزار در زبان های برنامه نویسی گوناگون حد فاصل تحقیقات در زمینه استفاده مجدد از نرم افزار هستند، اخیرا تلاش هایی برای استفاده از روش های داده کاوی برای تشخیص جنبه های تاثیر گذار بر موفقیت استفاده مجدد نرم افزار صورت گرفته است. انگیزه ی این رویکرد ها از این حقیقت ریشه می گیرد که تحقیقات قبلی نشان دهنده احتمال بالای شکست پروژه به دلیل کمبود فرایند استفاده مجدد و همچنین عدم وجود روشی برای تغییر فرایند های غیر قابل استفاده بود [۴۵].

در این راستا، موریسو و همکارانش [۴۵] تلاش کردند تا عوامل مهم موثر بر موفقیت استفاده مجدد از نرم افزار را در پروژه انجام شده پیدا کنند. به طور خاص، از طریق مصاحبه اطلاعاتی از ۲۴ پروژه اروپایی بین سال های ۱۹۹۴ تا ۱۹۹۷ جمع آوری کردند و در تحلیلشان، ۲۷ متغیر که برای توصیف پروژه ها استفاده می شدند تعریف کردند که خلاصه ای از آن ها در [۴۰] آمده است. در حال حاضر مجموعه داده فوق با وجود تعداد کم نمونه ها، بزرگترین مجموعه داده تجربی برای استفاده مجدد نرم افزار را تشکیل می دهد. روی این داده ها الگوریتم های داده کاوی زیادی اعمال شده اند تا الگوهای مربوط به عوامل موثر در موفقیت استفاده مجدد نرم افزار را تشخیص دهند.

### ۱.۶.۳ رده بندی

در همان مطالعه [۴۰]، نویسندگان همچنین تلاش کردند از درخت تصمیم و به طور خاص پیاده سازی درخت *J.48* در نرم افزار *Weka* که پیاده سازی درخت تصمیم *C4.5* است برای تحلیل همان داده ها استفاده کنند. کاربرد درخت تصمیم *C4.5* به این صورت بود که مهم ترین ویژگی ها را از بین ۲۷ ویژگی با استفاده از آزمایش حذف ویژگی ها تشخیص دهد. به طور خاص آن ها ریشه درخت تصمیم را در هر مورد در زمانی که مهم ترین ویژگی حذف شود (مثلا ریشه درخت)، به دست آوردند و درخت مجددا بدون آن ویژگی ساخته می

شد. این روش به آن‌ها اجازه می‌داد که ویژگی‌های ضعیف، که در طی مراحل حذف از درخت در ریشه ظاهر نمی‌شوند و یا ویژگی‌هایی که حذف آن‌ها تغییری در دقت رده بندی ایجاد نمی‌کند را شناسایی کنند. علاوه بر تجزیه و تحلیل فوق، آن‌ها از یک الگوریتم یادگیری به اسم *treatment learning* استفاده کرده و به طور خاص، الگوریتم *TAR2* را اعمال کردند. ایده ساده *treatment learning* این است که این الگوریتم، زیر مجموعه ای مانند *A* را از مجموعه داده‌های آموزشی *D* انتخاب می‌کند که بیشتر شامل رده‌های مطلوب بوده و تعداد کمی از رده‌های نامطلوب را داشته باشد. معیاری که بر اساس آن زیر مجموعه انتخاب می‌شود بر اساس  $R_x$  انتخاب شده است. *TAR2* از کاربر می‌خواهد که عددی را به عنوان امتیاز به هر رده اختصاص دهد تا رده‌ها بر اساس مطلوبیت مشخص شوند. در تشخیص قابلیت استفاده مجدد، پروژه‌های موفق از این نظر باید امتیاز بالا تری نسبت به پروژه‌هایی با قابلیت استفاده مجدد پایین بگیرند. پس از تحلیل‌های انجام شده، آن‌ها می‌توانند ویژگی‌ها و مقادیر مربوط به آن‌ها که در پروژه‌های موفق استفاده شده‌اند، مشخص کنند.

بخش جالب کاربرد این الگوریتم داده کاوی این حقیقت است که این الگوریتم می‌تواند ویژگی‌ها و مقادیری را پیدا کند که مطالعه تجربی در [۴۵] نتوانسته بود آن‌ها را بیابد. این موضوع نشان دهنده اهمیت داده کاوی در مهندسی نرم‌افزار است.

از سوی دیگر، البته با استفاده از همان مجموعه داده برای استفاده مجدد نرم‌افزار، Jiang و همکارانش [۳۰] یک رویکرد یادگیری گروهی (به عنوان مثال روشی که تصمیم روش‌های رده بندی مختلف را ترکیب کرده و بهترین را ارائه بدهد) اعمال کردند. آن‌ها از الگوریتم گروهی *RF2TREE* که جنگل تصادفی یا مجموعه‌ای از درخت‌های تصمیم را به یک درخت تصمیم تبدیل می‌کند استفاده کردند. تضمین شده است که این الگوریتم به شرط برقراری موارد زیر به خوبی کار می‌کند:

۱. مجموعه داده آموزشی اصلی همانند مجموعه موجود در [۴۵] بسیار کوچک باشد.

۲. در صورتی که جنگل تصادفی و یک درخت تصمیم هر دو توسط مجموعه داده آموزشی اصلی آموزش دیده باشند، دقت درخت تصادفی

از درخت تصمیم بالا تر باشد.

الگوریتم ابتدا از روی مجموعه داده اصلی جنگل تصادفی ایجاد می کند و سپس این جنگل تصادفی برای ساخت نمونه های مجازی که برای آموزش درخت های تصمیم استفاده می شوند، به کار گرفته می شود.

بر اساس آزمایش های انجام شده، آن ها دریافتند که مهم ترین ویژگی هایی که بر موفقیت استفاده مجدد نرم افزار اثر می گذارند عبارت اند از: عوامل انسانی، فرایند استفاده مجدد معرفی شده، نوع نرم افزار تولید شده، تعهد مدیر ارشد، و تغییرات فرایند های غیر قابل استفاده مجدد که با تحلیل های تجربی و تحلیل های داده کاوی تفاوت دارند.

### ۲.۶.۳ کاوش الگوهای تکرار شونده و قوانین انجمنی

Stefano Di و Menzies [۴۰] بر روی داده ذکر شده کار کردند تا بتوانند نتایج بعدی عوامل موثر روی استفاده مجدد نرم افزار را بررسی کرده و الگوهای حاصل از اعمال روش های خودکار داده کاوی را با روش های تجربی گفته شده در [۴۵] مقایسه کنند. از بین روش های مختلف، کاوش قوانین انجمنی به کار گرفته شد تا ارتباط معنا دار بین ۲۷ ویژگی را استخراج کند. استخراج قوانین انجمنی با استفاده از الگوریتم *Apriori* [۱] پیاده سازی شده توسط نرم افزار *Weka* انجام شد. ۱۰ قانون برتر با اطمینان %۹۰ در [۴۰] آمده است. با استفاده از قوانین انجمنی به دست آمده، نتایج جالبی به دست آمد. به عنوان مثال این حقیقت که زمانی که یک محصول نرم افزاری در محصولی استفاده شده بود، استفاده از سیاست پاداش برای استفاده مجدد از نرم افزار فعال نشده بود.

SoftwareandProduct = product => RewardsPolicy = no

یکی دیگر از جنبه های مهم داده کاوی در مهندسی نرم افزار استخراج الگوهای طراحی از نرم افزار به منظور استفاده مجدد آن ها است. به طور خاص، فرایند داده کاوی با هدف استخراج تصمیم های طراحی که معمولاً در کد منبع پنهان شده اند انجام می شود. به طور معمول در طی فرایند طراحی سیستم نرم افزاری، قسمت های مختلف سیستم با طراحی اعمال شده برچسب گذاری نمی شوند و در نتیجه تصمیم های طراحی با سیستم موجود سازگار نیستند. این مشکل باعث عدم درک جزئیات سیستم می شود.



در این راستا، روش ها و ابزار هایی با هدف استخراج الگوهای طراحی از یک سیستم نرم افزاری ارائه شده اند.

## فصل ۴

# نتیجه گیری

با توجه به تولید داده‌های زیاد در مهندسی نرم‌افزار، کاربرد صحیح این داده‌ها برای حل مشکلات چرخه تولید نرم‌افزار اهمیت بسزایی دارد. برخی مشکلات جدی همانند ایجاد اشکال، افزایش هزینه نگهداری برنامه یا نیازمندی‌های نامعلوم می‌تواند خلاقیت و کیفیت نرم‌افزار را کاهش دهند. این نوشتار به برخی روش‌های داده‌کاوی که توانایی به کارگیری همراه داده‌های مهندسی نرم‌افزار را در چالش‌های تولید، مدیریت، اشکال‌زدایی و نگهداری دارند، اشاره کرده است. ارزش نتایج داده‌کاوی است که آن را توجیه‌پذیر می‌کند و به همین دلیل کیفیت و کمیت داده‌های در دسترس و هزینه عملیات محاسباتی آن‌ها موفقیت این روش را در فرایند تولید نرم‌افزاری تعیین می‌کند. مهندسان و داده‌کاوان باید روش‌های داده‌کاوی را طوری انتخاب کنند که هزینه ابزار کاهش یابد و از داده‌های در دسترس بهترین استفاده به عمل آید. پیش‌بینی می‌شود در آینده تحقیقات این رشته در زمینه افزایش اتوماسیون و ساده‌تر شدن فرایندها متمرکز شود.

# کتاب نامه

- [۱] M. Halkidia, D. Spinellis, G. Tsatsaronis and M. Vazirgiannis, Data mining in software engineering
- [۲] J. Dong, Y. Zhao and T. Peng, A review of design pattern mining techniques, International Journal of Software Engineering and Knowledge Engineering (۶) ۱۹, (۲۰۰۹). ۸۵۵-۸۲۳
- [۳] P. Francis, D. Leon, M. Minch and A. Podguraki, Treebased method for classifying software failures. In Proceedings of the ۱۵th International Symposium on Software Reliability Engineering, .۲۰۰۴ software modifications. In Proceedings of ۲۰th IEEE International Conference on Software Maintenance (ICSM'۰۴), .۲۰۰۴
- [۴] A. Hindle, D. German and R. Holt, What do large commits tell us? a taxonomical study of large commits. In Proceedings of the IEEE Working Conference on Mining Software Repositories, .۲۰۰۸
- [۵] Y. Jiang, M. Li and Z.H. Zhou, Mining extremely small data sets with application to software reuse, Softw. Pract Exper (۴) ۳۹, ۲۰۰۹. ۴۴۰-۴۲۳
- [۶] M. Last, M. Friedman and A. Kandel, The Data Mining Approach to Automated Software Testing, In Proceeding of the SIGKDD Conference, .۲۰۰۵
- [۷] T. Menzies and J.S. Di Stefano, More success and failure factors in software reuse, IEEE Trans Software Eng (۵) ۲۹, (۲۰۰۳). ۴۷۷-۴۷۴
- [۸] S. Morisaki, A. Monden and T. Matsumura, Defect data analysis based on extended association rule mining. In Proceedings of International Workshop on Mining Software

Repositories (MSR), . 2007

[9] M. Morisio, M. Ezran and C. Tully, Success and failure factors in software reuse, IEEE Trans Software Eng (4) 28, (2002) . 357-340

[10] A. Podgurski, W. Masri, Y. McCleese, M. Minch, J. Sun, B. Wang and W. Masri, Automated support for classifying software failure reports. In Proceedings of the 25th International Conference on Software Engineering, . 2003

[11] U. Raza and M.J. Tretter, Predicting software outcomes using data mining and text mining. In SAS Global Forum, . 2007

[12] C.C. Williams and J.K. Hollingsworth, Automating mining of source code repositories to improve bug finding techniques, IEEE Transactions on Software Engineering (6) 31, (2005) . 480-466

[13] Lovedeep, Varinder Kaur Atri, Applications of Data Mining Techniques in Software Engineering, ISSN (Online): , 2347-2820, Volume , 2- Issue-5, 6, 2014

## **Abstract**

The increased availability of data created as part of the software development process allows us to apply novel analysis techniques on the data and use the results to guide the process's optimization. In this paper we describe various data sources and discuss the principles and techniques of data mining as applied on software engineering data. Data that can be mined is generated by most parts of the development process: requirements elicitation, development analysis, testing, debugging, and maintenance. Based on this classification we survey the mining approaches that have been used and categorize them according to the corresponding parts of the development process and the task they assist. Thus the survey provides researchers with a concise overview of data mining techniques applied to software engineering data, and aids practitioners on the selection of appropriate data mining techniques for their work.



Faculty of Science  
School of mathematics, statistics and computer science

# Application of Data mining in Software engineering

By

**Zahra Omid**

Supervisor

**Ebrahim Naghibzade Mashayekh**

Project for receiving bachelor degree  
Computer Science

July 2018