

دریافت مقاله: ۱۳۹۳/۱۱/۲۴

پذیرش مقاله: ۱۳۹۵/۳/۳۱

## مروری بر روش‌های تخمین هزینه نرم‌افزار مبتنی بر یادگیری ماشین

صبا بیرانوند<sup>۱\*</sup>، محمدعلی زارع چاهوکی<sup>۲</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، دانشکده برق و کامپیوتر، دانشگاه یزد، یزد، ایران

saba.beiranvand@stu.yazd.ac.ir

<sup>۲</sup> استادیار، دانشکده برق و کامپیوتر، دانشگاه یزد، یزد، ایران

chahooki@yazd.ac.ir

### چکیده

مدیریت پروژه نرم‌افزاری از مهم‌ترین فعالیت‌ها در توسعه محصول نرم‌افزاری است؛ زیرا تمامی فرایندها توسعه نرم‌افزار، از ابتدا تا انتها را شامل می‌شود. تخمین هزینه نرم‌افزار، یک فعالیت چالشی در مدیریت پروژه نرم‌افزاری است. مفهوم تخمین هزینه نرم‌افزار، همزمان با شروع صنعت کامپیوتر در سال ۱۹۴۰ مورد توجه قرار گرفته و همچنان پژوهش در این حوزه ادامه دارد. با اینکه تلاش، فقط دربرگیرنده بخشی از هزینه‌های توسعه یک پروژه نرم‌افزاری است، عامل اساسی برای تعیین هزینه محسوب می‌شود. از این رو، در پژوهش‌های این حوزه، دو اصطلاح تخمین تلاش و تخمین هزینه به صورت معادل به کار می‌روند. مدل تخمین هزینه نرم‌افزاری در صورتی که قبل از عقد قرارداد، دقت و اطمینان زیادی برای پیش‌بینی هزینه پروژه‌های نرم‌افزاری فراهم کند، مناسب است. به علت ذات غیرقطعی تخمین و برای افزایش دقت، به مرور توجه پژوهشگران به استفاده از روش‌های یادگیری ماشین در این حوزه معطوف شده است. در پژوهش حاضر، به بررسی مطالعات صورت گرفته در تخمین هزینه نرم‌افزار با روش‌های یادگیری ماشین پرداخته و روش‌های تخمین ارائه شده، معیارهای ارزیابی دقت این روش‌ها و دادگان مورد استفاده در کارهای پژوهشی و همچنین پژوهش‌های آتی در این حوزه را معرفی کرده‌ایم.

واژه‌های کلیدی: تخمین هزینه نرم‌افزار، تخمین تلاش نرم‌افزار، یادگیری ماشین.

## ۱. مقدمه

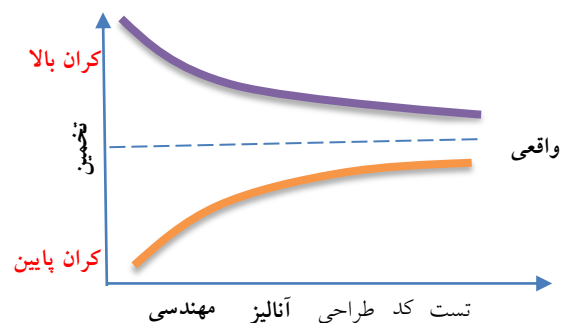
تلاش، منابع و هزینه هر فرایند توسعه استفاده می‌شود. بنابراین، عمل تخمین هزینه نرم‌افزار، برای مدیریت و نظارت بر توسعه نرم‌افزار مورد نیاز است. تا قبل از ارائه روش‌های تخمین، فرایند تخمین هزینه نرم‌افزار به برخی الگوریتم‌های ساده وابسته بود [۸].

در واقع، تخمین هزینه نرم‌افزار از سال ۱۹۶۰ شروع شد؛ زمانی که فرانک فریمن<sup>۳</sup> مفاهیم مدل‌های تخمین پارامتریک را توسعه داد و به توسعه مدل قیمت<sup>۴</sup> برای سخت‌افزار منجر شد. در سال ۱۹۷۰، محققان در این زمینه برای خیلی از پروژه‌ها، از روش‌های آماری، برای شناخت فاکتورهای مؤثر روی هزینه توسعه نرم‌افزار، با استفاده از روش‌های رگرسیون و وابستگی استفاده کردند. با اتمام این دهه، مدل COCOMO توسط باری بوهم<sup>۵</sup> و آبتس<sup>۶</sup> فرموله شد. مدل پارامتریک قیمت هم با اتمام این دوره ارائه شد. سپس تجزیه و تحلیل تابع نقطه‌ای (FPA<sup>۷</sup>)، توسط چند تن از کارکنان شرکت IBM برای تخمین هزینه و تلاش نرم‌افزار ارائه شد [۸].

به‌طور کلی، می‌توان گفت خیلی از تحولات روی مدل‌ها و روش‌های ارائه‌شده، در سال ۱۹۸۰ صورت گرفت. در همین سال بود که یکی از محققان پیشرو در تخمین هزینه نرم‌افزار، باری بوهم و همکارانش، روی مدل COCOMO که قبلاً خودشان ارائه کرده بودند، تغییراتی اعمال نمودند که حاصل آن مدل جدیدی به نام COCOMOII بود. از دهه ۹۰ به بعد، تحقیقات به سمت بهبودنمایی در صنعت نرم‌افزار و فناوری اطلاعات کشیده شد و بیشتر محققان در بهبود دقت الگوریتم‌های مختلف ارائه‌شده در این زمینه سعی کرده و روی این بحث متمرکز شده‌اند [۸]. عبدالله زید و همکارانش در مقاله پژوهشی خود [۸] روند کلی تحولات تاریخی فرایند تخمین هزینه نرم‌افزار را در شکل (۲) به تصویر کشیده‌اند.

مسئله مهم برای مدیران پروژه، تخمین‌های دقیق و مطمئن تلاش توسعه نرم‌افزار، به‌خصوص در مراحل اولیه چرخه توسعه نرم‌افزار است [۱ و ۲]. دلیل آن، ذات غیرقطعی تخمین و همچنین تغییرات سریع در روش‌های توسعه نرم‌افزار است [۲]. شکست پروژه، مسئله‌ای است که امروزه بسیار بدان توجه شده است و تخمین نادرست هزینه نرم‌افزار یکی از عوامل مؤثر بر شکست پروژه محسوب می‌شود. دو اصطلاح تخمین هزینه و تخمین تلاش مورد نیاز برای توسعه پروژه نرم‌افزاری (SDEE<sup>۱</sup>) در پژوهش‌های مهندسی نرم‌افزار و مدیریت پروژه، معادل هم، مورد استفاده قرار می‌گیرند. در واقع تلاش، عامل اصلی تعیین‌کننده هزینه است و این واقعیت وجود دارد که تخمین هزینه نرم‌افزار، در نهایت یک مسئله تخمین تلاش نرم‌افزار است [۳-۶].

هزینه و تلاش به شدت به هم وابسته‌اند. انتخاب بهترین مدل مبتنی بر میزان داده‌های در دسترس، روش‌های توسعه، روش‌های استفاده شده و عوامل دیگری از این قبیل است. شکل (۱) بیانگر رابطه تخمین دقیق، با تخمین انجام‌شده بر اساس مراحل مختلف توسعه پروژه است.



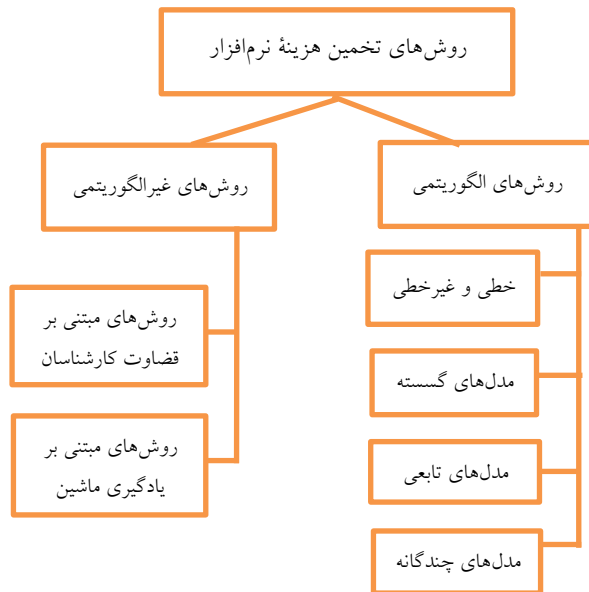
شکل (۱): رابطه تخمین دقیق با مراحل مختلف توسعه [۷]

## ۱.۱. تاریخچه

تخمین هزینه نرم‌افزار<sup>۲</sup> فرایندی قدیمی است که با شروع صنعت کامپیوتر در سال ۱۹۴۰ به وجود آمد و چندین بار توسعه یافت. از این فرایند در صنعت توسعه نرم‌افزار برای پیش‌بینی

3. Frank Freiman  
4. PRICE  
5. Barry W. Boehm  
6. C. Abts  
7. Function point analysis

1. Software Development Effort Estimation  
2. Software Cost Estimation



شکل (۳): دسته بندی روش های مختلف تخمین هزینه نرم افزار

در جدول (۱)، مقایسه ای بین دو دسته تعیین شده انجام شده است.

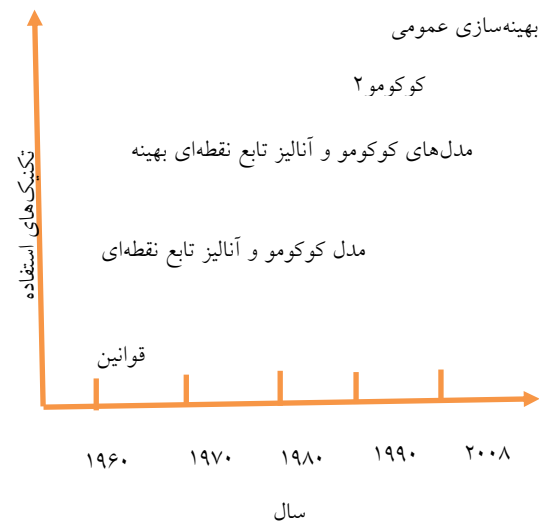
جدول (۱) مقایسه روش های الگوریتمی و غیر الگوریتمی [۷]

روش های الگوریتمی	روش های غیر الگوریتمی
انعطاف پذیری کم	انعطاف پذیری بالا
مبتنی بر مدل	مبتنی بر یادگیری
به اطلاعات زیادی نیاز دارد.	برخی از ویژگی های پروژه کافی اند.
روش های آماری	روش های مختلف
فرایند تخمین ساده	ممکن است فرایند پیچیده باشد.
نیازمند به روزرسانی است.	سازگار با تغییرات جدید است.
تخمین سریع	اغلب، تخمین زمان بر است.
هیچ انسانی نمی تواند در مدل مداخله داشته باشد.	کارشناسان می توانند روش را تنظیم کنند.
برآورد غلط در مراحل اولیه پروژه	تخمین دقیق در مراحل اولیه
قابل استفاده با پارامترهای مخصوص	قابل استفاده با پارامترهای مختلف

به طور کلی، تخمین تلاش مبتنی بر مدل های الگوریتمی، براساس رابطه (۱) طراحی می شود [۹].

$$EFFORT = F(x_1, x_2, \dots) \quad (1)$$

در این رابطه، منظور از  $x_1, x_2, \dots$  ویژگی هایی است که برای هر پروژه بیان می شود. تخمین تلاش مبتنی بر مدل های الگوریتمی مختلف، معمولاً به طور قابل توجهی، متفاوت از سایر مدل ها است. در واقع، در این روش ها مدلی براساس یک الگوریتم



شکل (۲): سیر تاریخی تحولات روش های تخمین هزینه نرم افزار [۸]

### ۲.۱. روش های مختلف تخمین هزینه نرم افزار

روش های مختلفی از آغاز تاکنون در زمینه تخمین هزینه و تلاش نرم افزار ارائه شده و مقالات مختلف، دسته بندی های متفاوتی برای این روش ها ارائه داده اند. برایان<sup>۱</sup> و همکارانش در [۴] روش های ارائه شده را در شش دسته مبتنی بر مدل، مبتنی بر خبره، مبتنی بر یادگیری، پویایی، رگرسیون و بیزین مرکب دسته بندی کرده اند. عبدالله زید و همکارانش در [۸]، تقسیم بندی روش های تخمین هزینه و تلاش در دو دسته کمی و کیفی را توصیه کرده اند که براساس آن، روش های کیفی بیشتر به شباهت و تفاوت وابسته اند و روش های کمی مثل روش های پارامتریک و تحلیلی، به جزئیات وابسته بوده و می توانند در فاز پیشرفته ای اجرا استفاده شوند.

عمده مقالات مانند [۶ و ۹] روش های تخمین هزینه نرم افزار را به سه دسته مدل های الگوریتمی، روش های مبتنی بر قضاوت کارشناسان<sup>۲</sup> و روش های یادگیری ماشین<sup>۳</sup> (ML) تقسیم کرده اند. البته دسته ای از مقالات، از جمله مقالات ارائه شده در [۷ و ۱۰] مدل مبتنی بر قضاوت کارشناسان را هم از دسته مدل های الگوریتمی می دانند و در اصل، روش های تخمین هزینه نرم افزار را در دو دسته کلی الگوریتمی و غیر الگوریتمی معرفی کرده اند. شکل (۳) نمایش تصویری این دسته بندی است.

1. Briand
2. Expert Judgment
3. Machine learning

مشابهی باشند. کار روش‌های یادگیری ماشین برای یادگیری الگوهای ذاتی ویژگی‌های پروژه‌ها، ارتباطات بین آن‌ها و تلاش‌های پروژه است که می‌تواند برای پیش‌بینی تلاش پروژه جدید استفاده شود [۶]. مهم‌ترین دلایل اهمیت تخمین هزینه نرم‌افزار به این شرح است که به طبقه‌بندی پروژه‌های بهبودیافته، با توجه به طرح جامع اقتصادی کمک می‌کند، در استفاده بهینه از منابع تأثیرگذار است، در ارزیابی تغییرات مورد استفاده است و مشتریان انتظار دارند هزینه در راستای هزینه واقعی باشد [۹].

با توجه به وجود عدم قطعیت در تخمین نرم‌افزاری، استفاده از روش‌های غیرقطعی و انعطاف‌پذیر یادگیری ماشین، می‌تواند نقش بسزایی در افزایش دقت تخمین داشته باشد. مزیت‌های این نوع از روش‌های تخمین، شامل توانایی مدل‌های هوش محاسباتی برای مدل کردن مجموعه‌ای پیچیده از ارتباطات بین تلاش و افزایشندگان آن و همین‌طور توانایی آن‌ها در یادگیری از داده‌های پروژه‌های قدیمی است [۱۱]. بدین منظور از پژوهش‌هایی که به مرور محاسبات نرم [۱۲] و یادگیری ماشین [۱۳] در این زمینه پرداخته‌اند، استفاده می‌شود.

### ۱.۳. روش تحقیق

مقالات فراوانی تاکنون در راستای کاربرد یادگیری ماشین در این زمینه شناسایی شده‌اند که در طول دوران پژوهش، به مطالعه این مقالات پرداخته شده است. دسته‌ای از این مقالات به مرور کارهای مختلف در این زمینه پرداخته‌اند. از میان مقالات انگلیسی، مقاله‌های [۴ و ۱۳] مقالاتی هستند که از نظر حوزه مورد بررسی، مشابه مقاله حاضرند. به همین دلیل، در مرور انجام‌شده، از آن‌ها نیز استفاده شده است. این مقالات، کارهای انگلیسی در محدوده سال‌های ۱۹۹۱ تا ۲۰۱۰ را دربرمی‌گیرند. در میان مقالات فارسی هم، مقالات مروری در زمینه تخمین هزینه نرم‌افزار ارائه شده است، اما این مقالات روی روش‌های یادگیری ماشین در این زمینه تمرکز نکرده‌اند. تفاوت مقاله حاضر با سایر مقالات مروری ارائه‌شده در این است که اولین مقاله مروری فارسی، در زمینه تخمین هزینه نرم‌افزار با روش‌های یادگیری ماشین است؛ ضمن اینکه مقالات سال‌های ۲۰۱۰ به بعد را هم در بردارد. در مرور

مشخص فرموله شده و از آن فرمول، برای تخمین هزینه نرم‌افزار استفاده می‌شود. تاکنون مدل‌های الگوریتمی مختلفی برای تخمین هزینه نرم‌افزار ارائه شده است (نک: شکل (۳)). از جمله معایب مدل‌های الگوریتمی، یکی این است که کار را برای مدیران سخت کرده‌اند؛ چراکه تصمیم‌گیری درباره میزان منابع واقعی، برای آن‌ها سخت شده و دیگر اینکه مبتنی بر داده‌های قدیمی‌اند و نمی‌توانند پیشرفت‌های فعلی در زبان‌های برنامه‌نویسی، سخت‌افزار و مهندسی نرم‌افزار را بازتاب دهند [۶].

مدل‌های غیرالگوریتمی بر استفاده از آزمایش‌ها و روش‌های خاص خود مبتنی هستند [۹]. معمولاً داده‌های قدیمی برای قضاوت خبرگان لازم نیستند. قضاوت خبرگان اغلب مبتنی بر استفاده مجدد از موارد تعیین‌کننده پروژه‌های قبلی است که حال ممکن است مستند هم نشده باشند. نتایج به دست آمده بیانگر آن است که ۶۲ درصد از تخمین هزینه پروژه‌های نرم‌افزاری در سازمان‌ها، از این روش استفاده می‌کنند. مزیتی این نوع از روش‌ها این است که برای فرهنگ سازمانی خاصی، سفارشی شده‌اند. پس در مقایسه با رویکرد الگوریتمی، دقیق‌تر خواهند بود و در خیلی از موارد ثابت شده است که دقت بهتری نسبت به سایر مدل‌ها ارائه داده‌اند. ولی این تخمین، ذهنی و براساس منطق خبره است و در واقع، مزیت آن را می‌توان به عنوان عیب آن هم در نظر گرفت؛ زیرا هر خبره فقط براساس تجربیات و فرهنگی که در آن سازمان خاص وجود دارد، هزینه را تخمین می‌زند و همان خبره در سازمان دیگر ممکن است هزینه نرم‌افزار را با دقتی کمتر تخمین بزند. پس لزوماً در محیط‌های مختلف، همان دقت را نخواهند داشت [۶].

روش‌های یادگیری ماشین، الگوهایی را از داده‌های پروژه‌های قدیمی یاد می‌گیرند و از این الگوها برای پیش‌بینی تلاش استفاده می‌کنند. ایده اساسی استفاده از روش‌های یادگیری ماشین برای تخمین و پیش‌بینی تلاش، این است که دادگان<sup>۱</sup> قدیمی شامل پروژه‌های قدیمی زیادی هستند که با خصوصیات بارزششان برای توصیف آن پروژه تشریح شده‌اند. پروژه‌های با خصوصیتی دارای مقادیر مشابه، ممکن است شامل تلاش‌های پروژه تقریباً

صورت گرفته مقالات فارسی هم بررسی شده‌اند.

در مرور صورت گرفته در این مقاله، به دنبال یافتن پاسخی برای پرسش‌های پژوهشی مطرح در زمینه تخمین هزینه نرم‌افزار با استفاده از رویکردهای یادگیری ماشین هستیم؛ پرسش‌هایی نظیر:

- کدام یک از روش‌های ML در SDEE استفاده شده‌اند؟
- دقت کلی تخمین مدل‌های ML چقدر است؟ (دادگان استفاده شده، متدهای اعتبارسنجی، متریک‌های استفاده شده برای اندازه‌گیری دقت در تخمین و میزان دقت)
- آیا مدل‌های ML بهتر از مدل‌های غیر ML هستند؟
- کدام مدل ML بهتر از سایر مدل‌های ML است؟ (مدل‌های ML در مقایسه با هم، درجه بهبود)
- هر روش ML در چه وضعیتی بهتر است؟ (نقاط قوت و ضعف مدل‌های ML از نظر SDEE)

این پرسش‌های پژوهشی همچنین در [۴] که مقاله مروری خوبی در این حوزه است، استفاده شده است. در [۴] مقالات بین سال‌های ۱۹۹۱ تا ۲۰۱۰ براساس پرسش‌های پژوهشی انتخاب و مرور شده‌اند. در این مقاله هم براساس همان پرسش‌های پژوهشی، مقالات ارائه شده در سال‌های بعد از ۲۰۱۰ انتخاب و مرور شده‌اند. البته برای تکمیل مطالب، از برخی مقالات مربوط به سال‌های قبل هم استفاده شده است.

در بخش ۲، به بررسی الگوریتم‌های یادگیری ماشین مختلفی که در زمینه تخمین هزینه نرم‌افزار استفاده شده‌اند، خواهیم پرداخت. دادگان مختلفی که در این زمینه استفاده شده‌اند، در بخش ۳ بررسی می‌شوند. در بخش ۴ معیارهای ارزیابی پرکاربرد در تخمین هزینه نرم‌افزار معرفی شده و در بخش ۵ به جمع‌بندی و ارائه پیشنهادهایی برای پژوهش‌های آتی در زمینه مورد بحث پرداخته می‌شود.

## ۲. یادگیری ماشین در تخمین هزینه نرم‌افزار

تخمین دقیق تلاش مورد نیاز در مراحل اولیه توسعه نرم‌افزار، نقش بسیار مهمی در مدیریت پروژه بر عهده دارد [۴]. از سال ۱۹۸۰ براساس پژوهش [۸] یعنی برهه‌ای از زمان که تحولات اساسی در

زمینه تخمین تلاش نرم‌افزار رخ داده است، بسیاری از روش‌های تخمین در زمینه تخمین هزینه نرم‌افزار ارائه شده‌اند که از میان آن‌ها از روش مبتنی بر قضاوت متخصصان و روش مقایسه‌ای<sup>۱</sup> با بیشتری استفاده شده است. در سال‌های اخیر، به روش‌های مبتنی بر یادگیری ماشین در این زمینه، بیشتر توجه شده است [۴].

ون<sup>۲</sup> و همکارانش در [۴] مروری سامان‌مند روی مقالاتی که از سال ۱۹۹۱ تا ۲۰۱۰ در زمینه تخمین هزینه نرم‌افزار با استفاده از روش‌های یادگیری ماشین ارائه شده‌اند، انجام داده‌اند. در این مقاله مروری، اشاره شده است که برخی از محققان حتی نشان داده‌اند که روش مبتنی بر یادگیری ماشین، یکی از سه دسته اصلی روش‌های تخمین (دو دسته دیگر شامل مدل‌های الگوریتمی و قضاوت هستند) است.

مقالات مختلفی تاکنون در این زمینه ارائه شده‌اند که دقت‌های مختلفی را با الگوریتم‌های یادگیری ماشین گزارش داده‌اند. اما معمولاً دقت‌های گزارش شده در کارهای مختلف، با هم متفاوت بوده و همین اختلاف در مطالعات تجربی موجود که علت آن هم هنوز قابل درک نیست، ممکن است باعث ممانعت از استفاده این الگوریتم‌ها توسط متخصصان در این زمینه شود. در مقایسه با سایر حوزه‌هایی که الگوریتم‌های یادگیری ماشین کاربرد دارند و با موفقیت پیاده‌سازی شده‌اند، در حوزه تخمین هزینه نرم‌افزار، چالش‌های زیادی از جمله دادگان کوچک، عدم قطعیت اطلاعات، معیارهای کیفی و عوامل انسانی وجود دارند. علاوه بر این، باید گفت که این نوع استفاده از روش‌های ML بسیار پیچیده‌تر از روش‌های تخمین معمولی است [۴].

با اینکه در سال‌های اخیر با هدف افزایش دقت، استفاده از مدل‌های یادگیری ماشین در تخمین تلاش توسعه نرم‌افزار، بیشتر مورد توجه قرار گرفته است، هیچ‌یک از مدل‌های موجود در تمام شرایط مناسب نیستند و از دادگانی به دادگان دیگر، دقت متفاوتی را ارائه می‌دهند. بنابراین، نیاز به ساختن مدلی تخمینی وجود دارد که مطمئن باشد و دقت بیشتری را فراهم کند. با همین هدف، ترکیبات مختلفی از این الگوریتم‌ها چه در ترکیب با الگوریتم‌های

یادگیری ماشین به کار گرفته شده‌اند.

GA با توجه به مروری که انجام شده است، فقط در فرم ترکیبی استفاده می‌شود. در این ترکیبات، روش GA تنها به منظور وزن دهی یا انتخاب ویژگی<sup>۴</sup>، استفاده شده است. منطق فازی هم برای مقابله با عدم اطمینان و عدم صحت اطلاعات استفاده می‌شود. در عمل، برای بهبود عملکرد مدل با پیش پردازش ورودی‌های مدل، از روش منطق فازی استفاده شده است [۴].

البته شایان ذکر است که برخی مقالات مثل [۳] به جای رگرسیون بردار پشتیبان، روش‌های رگرسیون را نام برده‌اند. در مطالعات مرور شده، روش‌های ANN و CBR به عنوان پرکاربردترین روش‌ها نسبت به سایر روش‌ها معرفی شده‌اند. در ادامه کار، با انواع دسته‌های الگوریتم‌های یادگیری ماشین آشنا خواهیم شد. در ضمیمه الف، شرح مختصری از مقالاتی که از روش‌های یادگیری ماشین در تخمین هزینه نرم افزار استفاده و در این پژوهش بررسی شده‌اند، به صورت جدول آمده است.

#### ۲.۱.۱. استدلال مبتنی بر مورد

روش استدلال مبتنی بر مورد، یکی از روش‌های یادگیری ماشین در تخمین هزینه نرم افزار است که امکان تطبیق با موارد مشابهی را که در گذشته رخ داده، برای متخصصان این زمینه فراهم می‌کند. CBR منبع استدلال قیاسی، حافظه پویا و نقش موقعیت‌های قبلی است که توجه بسیاری را در سال‌های اخیر به خود جلب کرده است. برخی از پژوهشگران، CBR را به صورت دوره‌ای و متشکل از چهار مرحله توصیف کرده‌اند: بازیابی مورد یا موارد خیلی مشابه مثل پروژه‌های توسعه یافته قبلی؛ استفاده مجدد از موارد بازیابی شده برای یافتن یک راه حل برای مسئله؛ دریافت راه حل یا تجدیدنظر در راه حل پیشنهادی؛ نگهداری راه حل به عنوان یک راه حل جدید.

در سال‌های اخیر، برخی تحقیقات ملموس در کاربرد CBR برای تخمین تلاش، نشان داده‌اند که CBR می‌توان امکان ویژه‌ای برای مدیریت توسعه نرم افزار فراهم کند [۳]. در این روش، از الگوریتم‌های مختلف برای محاسبه شباهت از معیارهایی مثل

یادگیری ماشین دیگر و چه در ترکیب با الگوریتم‌های غیر یادگیری ماشین، تاکنون ارائه شده است. محمود الیش<sup>۱</sup> و همکارانش در [۱۱] ترکیبات همگن<sup>۲</sup> و ناهمگن<sup>۳</sup> متفاوتی از الگوریتم‌های مختلف یادگیری ماشین ارائه داده و آزمایش‌های گوناگونی را روی آن‌ها انجام داده‌اند. ترکیبات ناهمگن شامل اعضایی با داشتن الگوریتم‌های یادگیری پایه و ترکیبات همگن شامل اعضایی با الگوریتم یادگیری تکنوع هستند. در ادامه، مروری بر کارهایی از این دست که می‌توان در زمینه تخمین تلاش و روش‌های یادگیری ماشین انجام داد، خواهیم داشت.

#### ۲.۱. معرفی روش‌ها

بر اساس مروری که در این پژوهش بر مطالعات صورت گرفته در زمینه تخمین هزینه نرم افزار از جمله [۴] انجام شده است، از سال ۱۹۹۱ تاکنون ۸ دسته الگوریتم یادگیری ماشین در حوزه تخمین هزینه نرم افزار استفاده شده‌اند. در جدول (۲) این دسته بندی‌ها ارائه شده و در ادامه، شرح مختصری از آن‌ها را خواهیم داشت. روش‌های پر کاربرد با \* نشان داده شده‌اند.

جدول (۲) روش‌های یادگیری ماشین مورد استفاده در تخمین هزینه نرم افزار

روش	عنوان	مخفف
CBR	Case-Based Reasoning	استدلال مبتنی بر مورد*
ANN	Artificial Neural Networks	شبکه‌های عصبی مصنوعی*
DT	Decision Trees	درخت تصمیم*
BN	Bayesian Networks	شبکه‌های بی‌زین
SVR	Support Vector Regression	رگرسیون بردار پشتیبان
GA	Genetic Algorithms	الگوریتم ژنتیک
GP	Genetic Programming	برنامه‌نویسی ژنتیک
AR	Association Rules	قوانین وابستگی

البته نکته مهم این است که علاوه بر دسته بندی ارائه شده در جدول (۲)، روش‌های یادگیری ماشین دیگری هم طی سال‌های اخیر ارائه شده و به تبع، در این حوزه به کار رفته‌اند. نکته گفتنی اینکه در برخی مطالعات، الگوریتم‌های یادگیری ماشین در ترکیب با الگوریتم‌های دیگر مورد استفاده قرار گرفته‌اند؛ مثلاً روش‌هایی مثل منطق فازی و GA اغلب برای ترکیب با سایر روش‌های

1. Mahmoud O. Elish  
2. Homogeneous  
3. Heterogeneous

تخمین هزینه نرم‌افزار استفاده کرده و دقت آن‌ها را مورد ارزیابی قرار داده‌اند.

### ۲.۱.۳. درخت تصمیم

دسته‌بندی داده‌ها، یک فرایند دومارحله‌ای است متشکل از: ۱. مرحله آموزش که در آن، مدل ساخته می‌شود؛ ۲. مرحله دسته‌بندی که در آن از مدل ساخته‌شده برای دسته‌بندی و برچسب‌زنی روی داده‌های موردنظر استفاده می‌شود.

درخت‌های تصمیم با چندتایی‌های آموزش و دسته‌های آن‌ها در معیار هدف ایجاد می‌شوند. این درخت از دو نوع نود تشکیل شده است که شامل نود تصمیم یا نود داخلی<sup>۱۰</sup> و نود برگ<sup>۱۱</sup> است. این نوع از روش‌های یادگیری ماشین، یک نمودار جریان<sup>۱۲</sup> با ساختاری مانند درخت است که هر گره داخلی یعنی نودهایی که برگ نیستند، یک آزمون را روی یکی از خصوصیات چندتایی انجام می‌دهند. هر شاخه<sup>۱۳</sup> بیانگر یک خروجی از آزمون و هر برگ نشان‌دهنده برچسب یا مقدار تخمین‌زده‌شده برای آن نمونه است. قوی‌ترین گره یک درخت، ریشه آن است.

به‌منظور انتخاب خصوصیتی به‌عنوان خصوصیت جداساز، خصوصیت موردنظر باید کمترین عدم خلوص را ارائه دهد. در صورتی یک نود خالص است که تمام نمونه‌هایی که در یک شاخه از درخت واقع شده‌اند، در یک دسته باشند. معیارهای انتخاب خصوصیت، روش‌های ابتکاری هستند که برای انتخاب بهترین جداسازی به‌منظور تقسیم پایگاه داده به‌کار می‌روند. اگر داده‌های D به بخش‌های کوچک‌تر تقسیم شود، در حالت ایدئال باید این قسمت‌ها خالص باشند. معیارهای جداسازی به‌عنوان قوانین جداساز<sup>۱۴</sup> هم شناخته می‌شوند؛ زیرا چگونگی جداسازی چندتایی‌ها را نیز تعیین می‌کنند. گره N چندتایی‌های D را نگاه می‌دارد. خصوصیتی که بیشترین بهره اطلاعاتی را دارد، به‌عنوان معیار جداسازی برای گره N انتخاب می‌شود. این خصوصیت، اطلاعات موردنیاز برای شناسایی چندتایی‌های درون گره‌های

فاصله مَنهت<sup>۱</sup>، اقلیدسی<sup>۲</sup> و معیارهایی از این دست استفاده می‌شود. در تألیفات گوناگون، این دسته از الگوریتم‌های یادگیری ماشین را با عنوان الگوریتم‌های غیر پارامتریک یاد می‌کنند [۱۴]. در مطالعه پژوهشی-آزمایشی صورت‌گرفته در [۱۵]، از الگوریتم تعیین فاصله اقلیدوسی، برای محاسبه فاصله یا میزان شباهت نمونه‌ها، برای استفاده از کاربرد الگوریتم CBR در تخمین هزینه نرم‌افزار استفاده شده است.

### ۲.۱.۲. شبکه‌های عصبی مصنوعی

ANNها تعداد زیادی از عناصر پردازشی به‌شدت بهم‌پیوسته به نام نورون‌ها<sup>۳</sup> دارند که معمولاً در عمل به‌صورت موازی هستند و تحت معماری منظمی، پیکربندی شده‌اند. هر نورون به‌وسیله یک لینک ارتباطی<sup>۴</sup>، با سایر نورون‌ها در ارتباط است. هر لینک ارتباطی وزنی دارد که بیانگر اطلاعاتی درباره سیگنال ورودی است. نورون، یک مجموع از وزن‌های ورودی را محاسبه می‌کند و اگر مجموع وزن‌ها بالاتر از یک حد آستانه<sup>۵</sup> باشد، یک خروجی را تولید می‌کند. این فرایند تا زمانی که یک یا چند خروجی تولید شوند، ادامه می‌یابد. این مدل‌های تخمین می‌توانند با استفاده از داده‌های آموزشی قدیمی، برای تولید نتایج، با تنظیم دقیق مقادیر پارامترهای الگوریتم، آموزش ببینند تا اختلاف بین مقدار واقعی شناخته‌شده با مدل پیش‌بینی را کاهش دهند [۳].

همچنین پژوهشگرانی در این زمینه مثل [۱۶] شبکه‌های عصبی مورد استفاده در تخمین نرم‌افزار را در انواع شبکه عصبی روبه‌جلو<sup>۶</sup>، شبکه عصبی مکرر<sup>۷</sup>، شبکه عصبی تابع پایه‌ای شعاعی (RBF)<sup>۸</sup>، شبکه عصبی‌فازی<sup>۹</sup> خلاصه کرده‌اند. پژوهشگران در برخی مقالات از جمله [۱۵]، از دو روش شبکه عصبی چندلایه (MLP) و شبکه عصبی تابع پایه‌ای شعاعی را به‌عنوان نمایندگانی از الگوریتم‌های شبکه عصبی در مطالعه آزمایشی خود، در حوزه

1. Manhattan
2. Eughlo
3. Neuron
4. Communication link
5. Threshold
6. Feed-forward neural network
7. Recurrent neural networks
8. Radial basis function network
9. Neuro-fuzzy networks

10. Internal node

11. leaf

12. Flowchart

13. Branch

14. Splitting Rules

حوزه، می توانید به این مقاله مراجعه کنید. در [۲۰] که پژوهشگران به بررسی روش های مختلفی روی دادگان محلی در این حوزه پرداخته اند، از این الگوریتم هم در مقایسات خود استفاده کرده اند.

#### ۲.۱.۵. الگوریتم ژنتیک

فرض کنید مجموعه خصوصیات انسان توسط کروموزوم های او به نسل بعدی منتقل می شوند؛ برای مثال، ژن ۱ می تواند رنگ چشم باشد، ژن ۲ طول قد، ژن ۳ رنگ مو و... حال اگر این کروموزوم به تمامی، به نسل بعد انتقال یابد، تمامی خصوصیات نسل بعدی شبیه به خصوصیات نسل قبل خواهد بود. بدیهی است که در عمل چنین اتفاقی رخ نمی دهد. در واقع، به صورت همزمان دو اتفاق برای کروموزوم ها می افتد. اتفاق اول جهش<sup>۴</sup> است. جهش به این صورت است که بعضی ژن ها به صورت کاملاً تصادفی تغییر می کنند. البته تعداد این گونه ژن ها بسیار کم است، اما در هر حال، این تغییر تصادفی همان گونه که پیش تر دیدیم، بسیار مهم است؛ مثلاً ژن رنگ چشم می تواند به صورت تصادفی باعث شود تا در نسل بعدی، یک نفر دارای چشمان سبز باشد، در حالی که تمامی نسل قبل دارای چشم قهوه ای بوده اند. علاوه بر جهش، اتفاق دیگری که می افتد و البته این اتفاق به تعداد بسیار بیشتری نسبت به جهش رخ می دهد، چسبیدن دو کروموزوم از طول به یکدیگر و تبادل برخی قطعات بین دو کروموزوم است. این مسئله با نام معبر<sup>۵</sup> شناخته می شود. این همان چیزی است که باعث می شود تا فرزندان، ترکیب ژن های متفاوتی را (نسبت به والدین خود) به فرزندان خود انتقال دهند [۲۱].

معمولاً الگوریتم های ژنتیک یک عدد احتمال اتصال دارند که احتمال به وجود آمدن فرزند را نشان می دهد. ارگانسیم ها با این احتمال دوباره با هم ترکیب می شوند. اتصال دو کروموزوم، فرزند ایجاد می کند که به نسل بعدی اضافه می شوند. این کارها انجام می شوند تا اینکه کاندیداهای مناسبی برای جواب، در نسل بعدی پیدا شوند. مرحله بعدی تغییر دادن فرزندان جدید است. الگوریتم های ژنتیک یک احتمال تغییر کوچک و ثابت دارند که معمولاً درجه ای در حدود ۰/۰۱ یا کمتر دارد. براساس این

جدید را حداقل می کند. به این ترتیب حداقل عدم خلوص و تصادفی بودن چندتایی های درون قسمت ها ایجاد می شود. این روش، تعداد آزمون های مورد نیاز برای دسته بندی چندتایی های مدنظر را به حداقل می رساند و موجب ایجاد درخت های تصمیم ساده تری می شود. از رابطه هایی که برای تعیین خصوصیت جداساز استفاده می شود، رابطه ایتروپی<sup>۱</sup> است که در زیر مشاهده می کنید [۱۷].

$$ENTROPY = - \sum_{i=1}^k p_m^i \log_2 p_m^i \quad (2)$$

در پژوهش ارائه شده در [۱۵]، از دو الگوریتم درخت های رگرسیون و کلاس بندی (CART<sup>۲</sup>) و روش درخت مدل<sup>۳</sup> از جمله الگوریتم های درخت تصمیم استفاده شده در حوزه تخمین هزینه نرم افزار معرفی و در مطالعه ای آزمایشی در همین حوزه استفاده شده است.

#### ۲.۱.۴. شبکه های بیزین

یک شبکه بیزی یا شبکه پاور یا شبکه پاوربیزی، یک گراف جهت دار غیر مدور است که مجموعه ای از متغیرهای تصادفی و نحوه ارتباط مستقل آن ها را نشان می دهد؛ برای نمونه، یک شبکه بیزی می تواند نشان دهنده ارتباط بین علت بیماری ها با خود آن ها باشد. پس با داشتن عوامل، می توان احتمال یک بیماری خاص را در یک مریض تشخیص داد. شبکه بیزین یک ابزار نسبتاً جدید برای شناسایی (هویت) روابط احتمالی، به منظور پیشگویی یا ارزیابی کلاس عضویت است. به طور خلاصه، می توان گفت شبکه بیزین، نمایش با معنی روابط نامشخص مابین پارامترها در یک حوزه است. شبکه بیزین، گراف جهت دار غیر حلقوی از نودها، برای نمایش متغیرهای تصادفی و کمان ها برای نمایش روابط احتمالی مابین متغیرها به شمار می رود [۱۸].

در [۱۹] به بررسی ۲۳ کار پژوهشی صورت گرفته طی سال های ۲۰۰۲ تا ۲۰۱۰ که براساس شبکه های بیزین در حوزه تخمین هزینه نرم افزار کار کرده اند، پرداخته شده است. برای درک بیشتر این نوع از الگوریتم های یادگیری ماشین استفاده شده در این

4. Mutation  
5. Crossover

1. Entropy  
2. Classification and Regression Trees  
3. Model tree technique



از برنامه‌نویسی ژنتیک برای ارائه تخمین هزینه توسعه نرم‌افزار، آن‌هم به صورت دقیق، استفاده کرده و آن را با مدل رگرسیون آماری و الگوریتم شبکه عصبی مقایسه نموده‌اند.

## ۲.۱.۷. قوانین وابستگی

ما به دنبال یافتن قوانین وابستگی<sup>۱</sup> از روی مجموعه داده‌های موجود در دادگان مورد استفاده برای ارزیابی هستیم. یک قانون وابستگی را می‌توان به صورت  $A \rightarrow B$  نشان داد که در آن  $A$  و  $B$  زیرمجموعه‌های غیر تهی از  $I$  (بیانگر مجموعه اقلام) هستند و هیچ اشتراکی با همدیگر ندارند. به طوری که اگر  $A$  و احتمال وقوع آن منجر شود که احتمال وقوع  $B$  هم بیش از یک حدی باشد که مورد علاقه ماست، آنگاه  $A$  و  $B$  هر دو، گزاره‌هایی از این نقش خواهند بود. به زبان ساده‌تر، یعنی اینکه ما به دنبال یافتن رابطه رخ دادن دو گزاره باهم هستیم. اطمینان<sup>۲</sup> و پشتیبان<sup>۳</sup> از معیارهای سنجش مورد علاقه بودن قوانین محسوب می‌شوند و در واقع، قطعیت و کاربرد قوانین شناسایی شده را می‌سنجند. قانون  $A \rightarrow B$  دارای پشتیبان  $S$  در  $D$  است که  $S$  درصدی از تراکنش‌های  $D$  است که حاوی  $A \cup B$  باشد و می‌توان آن را با احتمال  $p(A \cup B)$  نیز نشان داد. اطمینان  $A \rightarrow B$  نیز با  $C$  در مجموعه تراکنش  $D$  نشان داده می‌شود که نشان‌دهنده درصدی از تراکنش‌های  $D$  است که شامل  $A$  و همچنین  $B$  است. اطمینان را نیز می‌توان با احتمال شرطی  $p(A|B)$  نشان داد [۱۷]. به این ترتیب، رابطه‌ها به شکل رابطه‌های (۳) و (۴) هستند.

$$(3) \quad p(A \rightarrow B) = p(A \cup B) \text{ پشتیبان}$$

$$(4) \quad p(A \rightarrow B) = p(A|B) \text{ اطمینان}$$

محققانی در [۲۴] به ارائه مدلی از ترکیب قوانین وابستگی و درخت تصمیم CART پرداخته‌اند. آن‌ها با آزمایش‌های خود روی دادگان مختلف، نشان داده‌اند که استفاده از این الگوریتم‌ها سبب افزایش دقت در این حوزه از تخمین در مهندسی نرم‌افزار می‌شود.

احتمال، کروموزوم‌های فرزند به‌طور تصادفی تغییر می‌کنند یا جهش می‌یابند؛ به‌خصوص با جهش بیت‌ها در کروموزوم ساختمان داده انسان. این فرایند باعث به‌وجود آمدن نسل جدیدی از کروموزوم‌هایی می‌شود که با نسل قبلی متفاوت است. کل فرایند برای نسل بعدی هم تکرار می‌شود. جفت‌ها برای ترکیب انتخاب می‌شوند، جمعیت نسل سوم به‌وجود می‌آیند و این عمل به همین شکل ادامه می‌یابد. این فرایند تکرار می‌شود تا به آخرین مرحله برسیم [۲۱]. شرایط خاتمه الگوریتم‌های ژنتیک عبارت‌اند از اینکه به تعداد ثابتی از نسل‌ها برسیم، بودجه اختصاص داده شده تمام شود (زمان محاسبه/پول)، یک فرد (فرزند تولید شده) پیدا شود که مینیمم (کمترین) ملاک را برآورده کند، بیشترین درجه برازش فرزندان حاصل شود یا دیگر نتایج بهتری حاصل نشود، بازرسی دستی و یا اینکه از ترکیب‌های بالا استفاده شود.

این الگوریتم از نوع الگوریتم‌های یادگیری نظارتی است که در آن، از قانون دلتا استفاده می‌شود. در این الگوریتم، به یک آموزگار نیاز داریم که خروجی مناسب برای یک ورودی را بداند یا اینکه بتواند آن را حساب کند [۲۱].

محققان دیگری در [۲۲] به ارائه الگوریتمی ترکیبی بر پایه دو الگوریتم یادگیری ماشین ژنتیک و شبکه عصبی برای تخمین هزینه نرم‌افزار پرداخته‌اند. الگوریتم ژنتیک برای مسائل بهینه‌سازی کاربرد دارد. در این مقاله، هدف آن است که از طریق الگوریتم ژنتیک، فقط ورودی‌هایی که عامل مهمی برای هزینه هستند یا پیچیدگی محاسباتی فراوانی ندارند، انتخاب شوند و به‌عنوان ورودی به شبکه مصنوعی داده شوند.

## ۲.۱.۶. برنامه‌نویسی ژنتیک

برنامه‌نویسی ژنتیک یکی از روش‌های در حال ظهور در تخمین هزینه نرم‌افزار است. ایده اصلی استفاده شده مبتنی بر تئوری تکاملی داروین است که می‌گوید عملیات ژنتیک بین کروموزوم‌ها در نهایت، به افرادی مشخص که احتمال بیشتری برای زنده ماندن دارند، منجر می‌شود. GP در واقع توسعه یافته GA است که محدودیت کروموزوم نشان‌دهنده فردی را که دارای یک رشته دودویی ثابت است، حذف کرده است [۳]. در [۲۳] پژوهشگران

1. Association Rules  
2. Confidant  
3. Support

## ۲.۱.۸ رگرسیون بردار پشتیبان

دسته‌ای<sup>۱۲</sup> و داده‌های گم شده<sup>۱۳</sup> هستند. خصوصیات روش‌های ML را طبق آنچه در جدول (۳) وجود دارد، می‌توان شرح داد. علامت‌های استفاده شده در این جدول بدین شرح هستند که \* به معنای مطلوب، + به معنای نامطلوب و - به معنای ذکر نشده است.

جدول (۳): خصوصیات روش‌های ML [۴]

داده‌های گم شده	ویژگی‌های دسته‌ای	داده پرت	دادگان کوچک	
+	+	*	*	CBR
-	+ <sup>۱۲</sup>	*	+	ANN
-	*	*	-	DT
* <sup>۱۵</sup>	* <sup>۱۶</sup>	* <sup>۱۷</sup>	*	BN

از جدول (۳) می‌توان به این نتایج رسید که روش DT مستعد بیش برآزش<sup>۱۸</sup> روی دادگان آموزشی کوچک است. دو روش ANN و CBR تا زمانی که ویژگی‌های دسته‌ای به ویژگی‌های کمی تبدیل نشوند، نمی‌توانند با آن‌ها به فرم استانداردشان کار کنند. مثل دو روش ANN و CBR، روش DT هم نمی‌تواند تا زمانی که به داده‌های گم شده مقدار نسبت داده نشده، نمی‌تواند با آن‌ها کار کند. علاوه بر نتایج حاصل از جدول بالا می‌توان گفت دو روش CBR و DT بصری و به آسانی قابل درک هستند. روش ANN توانایی یادگیری توابع پیچیده‌ای را دارد، اما مسائلی که باید در نظر گرفت، بدین شرح است که به مقادیر زیادی از داده‌های آموزشی نیازمند است، اغلب از مسئله پیش‌برآزش رنج می‌برد و همچنین توانایی توضیحی آن ضعیف است. روش BN قادر به یادگیری روابط علت و معلولی است. برای سه روش GP، SVR و AR اطلاعات چندانی در دست نیست [۴].

## ۲.۳.۲ مروری مقایسه‌ای بر مدل‌های مختلف ارائه شده

همان‌طور که در [۴] آمده، اوج انتشارات در سال ۲۰۰۸ است و به نظر می‌رسد که از دو روش CBR و ANN در مقایسه با سایر روش‌های یادگیری ماشین، در سال‌های اخیر بیشتر استفاده شده

رگرسیون بردار پشتیبان<sup>۱</sup> یا SVR یکی دیگر از روش‌های یادگیری ماشین پرکاربرد در زمینه تخمین هزینه و تلاش مورد نیاز برای توسعه محصول نرم‌افزاری محسوب می‌شود. در این نوع از روش‌های یادگیری ماشین، فرضیه‌ای<sup>۲</sup> برای جداسازی نمونه‌های یک دسته از دسته‌های دیگر را مشخص می‌کنیم. در رویکرد نرم این نوع روش، فرضیه خود را با میزان مشخصی تحطی از حاشیه<sup>۳</sup> تعیین شده برای آن، در نظر می‌گیریم. در واقع ایده اصلی رویکرد نرم، این است که خطاهای کمتر از یک حد آستانه تعیین شده و همین‌طور خطاهای نزدیک به صفر را نادیده بگیریم [۱۴]. در [۱۵] از ماشین‌های بردار پشتیبان حداقل مربعات (LS-SVM)<sup>۴</sup> به نمایندگی از روش‌های رگرسیون بردار پشتیبان، استفاده شده و این روش با سایر روش‌های استفاده شده در این پژوهش از جمله ols<sup>۵</sup>، ترکیب ols و تغییر شکل لگاریتمی<sup>۶</sup>، رگرسیون قوی<sup>۷</sup>، رگرسیون ols به همراه تغییر شکل BC<sup>۸</sup>، رگرسیون میانگین مربعات<sup>۹</sup> و دو الگوریتم MARS<sup>۱۰</sup> و رگرسیون RIDGE مورد مقایسه قرار گرفته است.

## ۲.۲.۲ بستر مطلوب هریک از روش‌های یادگیری ماشین

به منظور انتخاب یک مدل از میان مدل‌های ممکن، برای تخمین تلاش توسعه نرم‌افزار باید به نقاط ضعف و نقاط قوت مربوط به آن‌ها توجه کرد و در انتخاب روش، این ویژگی‌ها را مدنظر قرار داد. در بررسی‌هایی که در [۴] روی مطالعات انتخابی انجام شده است، خصوصیات روش‌های یادگیری ماشین، عمدتاً در چهار بستر<sup>۱۱</sup> از زمینه‌های تخمین، قابل دسترسی و مرتبط هستند. این بسترها شامل دادگان کوچک، داده‌های پرت، خصوصیات

1. Support Vector Regression
2. Hypothesis
3. Margin
4. Least squares support vector machines
5. Ordinary least squares regression
6. OLS regression with log transformation
7. Robust regression
8. Box Cox transformation
9. Least median squares regression
10. Multivariate Adaptive Regression Splines
11. Contexts

12. Categorical

13. Miss value

۱۴. تنها توسط مقاله ۶۸ از مطالعات انتخابی پشتیبانی می‌شود.

۱۵. تنها توسط مقاله ۶۷ از مطالعات انتخابی پشتیبانی می‌شود.

۱۶. تنها توسط مقاله ۷۸ از مطالعات انتخابی پشتیبانی می‌شود.

۱۷. تنها توسط مقاله ۷۸ از مطالعات انتخابی پشتیبانی می‌شود.

18. Over-fitting

توصیه شده که هر دو نوع مدل را به صورت موازی در مراحل اولیه کار استفاده کنید. تنها زمانی که مدل ML به طور قابل توجه و پیوسته‌ای بهتر از سایر مدل‌های موجود بود، می‌توان به متخصصان پیشنهاد کرد که مدل‌های موجود را با تمرکز بر مدل ML و تغییر به سیستم تخمین جدید تشویق کنند. نکته مهمی که باید همیشه مدنظر داشت، این است که قبل از انتخاب مدل ML، متخصصان نه تنها به آگاهی از بسترهای تخمین نیاز دارند، بلکه به فهم دقیق ویژگی‌های مدل‌های ML کاندید هم نیازمندند. معمولاً میزان تطبیق بین بسترهای تخمین و ویژگی‌های مدل ML انتخابی، تأثیر مستقیم و قابل توجهی روی عملکرد مدل دارند [۴].

#### ۴.۲. پیش‌پردازش‌های استفاده شده

در برخی از کارهای پژوهشی قبل از اینکه از دادگان استفاده کنند، به منظورهای خاصی از جمله افزایش دقت تخمین، از روش‌هایی استفاده کرده‌اند. همان‌طور که قبلاً هم بیان شد، در این حوزه از مهندسی نرم‌افزار از الگوریتم ژنتیک برای وزندهی به ویژگی‌های دادگان مورد استفاده در تخمین هزینه نرم‌افزار استفاده شده است.

دولت‌آبادی‌نژاد به همراه همکارانش در [۲]، به بررسی تأثیر کاهش انحراف توزیع تلاش بر دقت تخمین تلاش توسعه نرم‌افزاری پرداخته‌اند. کاهش انحراف توزیع تلاش، از ایجاد سربارهای مالی و زمانی در توسعه محصول نرم‌افزاری جلوگیری می‌کند. از جمله مواردی که گروه مدیریت بایستی پس از تخمین تلاش به انجام برسانند، توزیع تلاش در بین فازهای توسعه محصول نرم‌افزاری است. انحراف توزیع تلاش نیز مانند خود تخمین تلاش و هزینه توسعه نرم‌افزار، یکی از چالش‌های پیش روی گروه مدیریت توسعه نرم‌افزاری محسوب می‌شود. جهت توزیع تلاش نیز مدل‌های گوناگونی تاکنون ارائه گردیده است. آن‌ها پس از انجام آزمایش‌ها به این نتیجه رسیده‌اند که کاهش انحراف توزیع تلاش، در بالا بردن دقت، تأثیر مثبتی خواهد داشت. در پایان این پژوهش، به بیان دو راهکار مدیریتی برای توزیع و تخمین تلاش در مدیریت توسعه نرم‌افزار پرداخته شده است: اول اینکه اگر گروه مدیریت پروژه، با در نظر گرفتن موارد ساینز، دامنه و نوع توسعه، اقدام به توزیع تلاش در فازهای

است. البته به این نکته هم باید توجه داشت که در برخی از مطالعات، بیش از یک روش یادگیری ماشین، مورد مطالعه قرار گرفته است.

#### ۳.۲.۱. مدل‌های ML در مقایسه با هم

در مقاله [۴] مقایساتی بین الگوریتم‌های یادگیری ماشین استفاده شده در تخمین هزینه نرم‌افزار صورت پذیرفته است. از این مقایسات، سه نتیجه به دست می‌آید:

- روش‌های CBR و ANN دقیق‌تر از روش DT هستند و این نتیجه توسط بیشتر آزمایش‌ها در طول مقایسه اثبات شده است.
- در مواردی از مقایسات بین دو روش CBR و ANN، به نظر می‌رسد که روش CBR از روش ANN دقیق‌تر است. با این حال، این نتایج با آنچه در مقایسه دقت بین روش‌های ML گفته شده، که روش ANN از روش CBR بهتر است، تناقض دارد. این تناقض ممکن است به سه دلیل رخ داده باشد: اول اینکه تمام آزمایش‌های انتخابی که حاکی از برتری CBR نسبت به ANN است، از مطالعاتی که روی روش CBR کار کرده‌اند انتخاب شده‌اند و همین امر ممکن است باعث بایاس<sup>۱</sup> شدن مقایسات شده باشد؛ دوم اینکه مطالعاتی که روی روش ANN کار کرده‌اند، به ندرت آزمایشی برای مقایسه آن با روش CBR انجام داده‌اند، به طوری که مطالعاتی که دقت بیشتری را برای این روش گزارش داده‌اند، کمکی به این مقایسه نکرده‌اند؛ دلیل سوم هم این است که تعداد آزمایش‌های مربوط به CBR دو برابر آزمایش‌های مربوط به ANN برای این مقایسه است.
- برای مقایسه بین مدل‌های ML، تعداد آزمایش‌ها نسبتاً کم است و برخی از مقایسات حتی در تناقض با هم هستند. بنابراین در مقایسه آن دسته از روش‌های ML که به ندرت با یکدیگر مقایسه شده‌اند، تعیین دقیق‌ترین آن‌ها کار نسبتاً سختی خواهد بود. با توجه به نتایج متناقض مقایسه بین مدل‌های ML و مدل‌های غیر ML که در این بررسی آمده است، در این مقاله

خوشه بندی فازی و تخمین بر پایه مقایسه، که اساس کار این روش هاست، به صورت مختصر شرح داده شده است.

خوشه بندی داده ها روشی برای تقسیم همه داده ها در چندین خوشه است به قسمی که داده های مشابه در یک خوشه باشند و خوشه ها تا حد ممکن از یکدیگر مجزا باشند. ایده اصلی خوشه بندی فازی این است که «یک داده می تواند در چندین خوشه باشد». از این رو، یک سطح عضویت برای تعیین چگونگی وجود یک داده در یک خوشه تعریف می شود. مقدار سطح عضویت از صفر تا یک متغیر است. خوشه بندی C-MEANS فازی<sup>۲</sup> مهم ترین روش خوشه بندی فازی است که در سال ۱۹۸۱ مطرح شد. هدف این الگوریتم به حداقل رساندن تابع تعلق است و از آن برای حل مسائل تشخیص و پیش بینی استفاده می کنند. الگوریتم خوشه بندی C-Means فازی در ادامه این مقاله شرح داده شده است.

خوشه بندی پروژه های نرم افزاری، بخش کلیدی روش تخمین مطرح شده در این مقاله است. در اولین مرحله از این روش، از روش خوشه بندی فازی که بر اساس درجه شباهت هر پروژه به خوشه ها عمل می کند، استفاده شده است. پیشنهاد پژوهشگران این مقاله برای غلبه بر مشکل داده های پرت، آموزش شبکه عصبی با استفاده از خوشه های به دست آمده است. پس از مقداردهی دادگان با خوشه ها، ممکن است خوشه هایی شامل چند پروژه که مناسب و قابل استفاده نیستند، در مجموعه آموزشی شبکه عصبی ظاهر شوند. این مشکل ممکن است روش پیشنهادی را مجبور به نادیده گرفتن این خوشه ها کند. برای اجتناب از این محدودیت، ABE<sup>۳</sup> به روش پیشنهادی اضافه شده است؛ زیرا این روش می تواند در هر دادگانی، با هر تعداد پروژه ای برای تخمین استفاده شود. بنابراین با ترکیب شبکه های عصبی و ABE همه پروژه های موجود در دادگان می توانند در روند تخمین مشارکت کنند. به منظور شرح روند ترکیب سه بخش مطرح شده، روش ارائه شده در این کار پژوهشی، در دو فاز اصلی به نام های آموزش و آزمایش، سازمان دهی شده است. در مرحله

توسعه محصول نرم افزاری کنند، می تواند باعث کاهش انحراف توزیع تلاش شده و از به وجود آمدن سربارهای مالی و زمانی جلوگیری کند. دوم اینکه استفاده از مدل های توزیع تلاش، با در نظر گرفتن ویژگی های مؤثر در انحراف توزیع تلاش، در نهایت می تواند باعث افزایش دقت تخمین تلاش در توسعه محصول نرم افزاری گردد و گام مهمی در مدیریت بهینه پروژه نرم افزاری محسوب شود [۲].

از بررسی نتایج حاصل از [۹] که اخوان بیطرف و همکارانش روی روش های یادگیری ماشین برای پیمایش دادگان در تخمین هزینه نرم افزار کار کرده بودند، می توان به این نتیجه کلی رسید که پیمایش با توجه به کاهش کمیته میانگین خطا، اقدام درستی بوده و باعث شده که تخمین ها به مقدار واقعی نزدیک تر شوند. افزایش مقدار بیشینه میانگین خطا هم، بدین شکل توجیه شده که با توجه به عدم حساسیت به داده ها ممکن است تعداد معدودی از داده ها، دارای اختلاف زیاد خروجی با مقدار واقعی باشند.

خطیبی بردسیری به همراه همکارانش در [۲۵] برای افزایش دقت تخمین تلاش توسعه نرم افزار، دست به استفاده از خوشه بندی پروژه ها زده اند. در ابتدای کار، به این نکته پرداخته شده است که اساساً ویژگی های پروژه های نرم افزاری، خیلی پیچیده، غیر خطی و غیر قابل فهم هستند و این به دلیل نامشخص و بی ثبات بودن طبیعت پروژه های نرم افزاری است. روش های غیر الگوریتمی موجود از یک دادگان، شامل پروژه های نرم افزاری متنوع به منظور تخمین استفاده می کنند. در این چنین دادگانی پروژه های نامرتب و متناقض زیادی وجود دارد. اصطلاحاً در بحث داده کاوی به این نوع داده ها، داده های پرت می گوئیم. روش های یادگیری ماشین، برای خود یک مدل آموزشی با استفاده از داده های موجود در دادگان می سازند و بر اساس این مدل آموزشی، عملیات تخمین را انجام می دهند. واضح است که وجود داده پرت، کیفیت آموزش این الگوریتم ها را تحت تأثیر قرار داده و سبب تخمین های غیر دقیق و بی ثبات خواهد شد. از این رو، این مقاله با هدف رفع این مشکل نوشته شده و این مشکل را از طریق خوشه بندی<sup>۱</sup> پروژه ها حل کرده است. در ادامه، روش های

2. Fuzzy C-Means Clustering

3. Analogy Based Estimation

1. Clustering

آموزش، ساختار روش ترکیبی مطرح شده پیکربندی شده است. در مرحله آزمون هم روند تخمین تلاش انجام می شود.

برخی دیگر از محققان در مقالات خود [۲۶ و ۲۷] از روش نرمال سازی MIN-MAX به عنوان پیش پردازش دادگان خود استفاده کرده اند. برخی هم مانند [۱۳] از روش کاهش بعد PCA به عنوان پیش پردازشی روی دادگان خود استفاده کرده و نشان داده اند که این عمل باعث افزایش دقت و کاهش میانگین خطای مقدار تخمین زده شده، با مقدار واقعی هزینه مورد نیاز برای توسعه محصول نرم افزاری می شود.

در کاری که در [۹] انجام شده، از الگوریتم MP5 به عنوان پیش پردازشی برای تعیین وزن خصوصیات در شبکه عصبی استفاده شده است. در این کار پژوهشی، چون درخت های تصمیم به دنبال بهترین خصوصیات برای ارائه جداسازی با کمترین عدم خلوص هستند، تأثیر خصوصیات مشخص می شود. بر اساس آزمایش های نشان داده شده در این مقاله، اثبات شده که این روش برای تعیین وزن هم مناسب است.

دسته ای دیگر از مقالات [۲۸] از انتخاب ویژگی به منظور پیش پردازشی روی دادگان مورد استفاده، برای افزایش دقت تخمین هزینه نرم افزار و ارائه نزدیک ترین مقدار تخمین زده شده به میزان تلاش واقعی نمونه های موجود در دادگان استفاده کرده اند. البته این دسته از الگوریتم های یادگیری ماشین، انواع مختلفی دارند. نوعی از آن ها از یک ویژگی شروع و هر بار ویژگی ها اضافه می شوند تا زمانی که یا تأثیری در دقت نداشته باشند یا اینکه دقت را کاهش دهند. به این نوع از انتخاب ویژگی روبه جلو می گویند. دسته دیگر، از همه ویژگی ها کار را شروع می کند و برعکس حالت قبل، هر بار یکی از ویژگی ها را کم می کند. شرط اتمام این دسته هم مانند دسته اول است. به این دسته از الگوریتم های انتخاب ویژگی روبه عقب می گویند. این دو دسته الگوریتم های معرفی شده از دسته متدهای بسته بندی هستند. علاوه بر این متدها، متدهای فیلتر و یکپارچه سازی نیز از انواع متدهای انتخاب ویژگی محسوب می شوند.

## ۵.۲. سایر کارهای انجام شده

درک و محاسبه مدل های مبتنی بر داده های قدیمی، با توجه به

روابط پیچیده ای که بین خصوصیات وجود دارد، مشکل است. ویژگی ها و ارتباطات استفاده شده برای تخمین تلاش توسعه نرم افزار، می تواند در طول زمان و برای محیط های توسعه متفاوت، تغییر کند. به منظور رسیدگی و غلبه بر این مشکلات، یک مدل جدید با تخمین دقیق، مورد نیاز خواهد بود. در همین راستا در مقاله پژوهشی [۲۹] نسبت به ارائه یک مدل تخمین هزینه نرم افزار جدید اقدام نموده اند. پژوهشگرانی در [۳] به معرفی مدل های تخمینی که تاکنون به صورت موفقیت آمیزی مورد استفاده قرار گرفته اند، پرداخته اند. پژوهشگران در این مقاله آورده اند که چندین مدل تاکنون برای تخمین هزینه نرم افزار ارائه شده اند، اما تصمیم گیری درباره اینکه چه مدلی را انتخاب کنیم، سخت است. بنابراین برای حل این مشکل باید شناخت کاملی درباره روش های مختلف داشته باشیم. نتایج تحقیقات نشان می دهند که بین ۳۰ تا ۴۰ درصد پروژه های نرم افزاری به شکلی کامل، تکمیل شده و سایر پروژه ها شکست می خورند. هدف اصلی در انتخاب مدل، به حداکثر رسانی دقت است. کمار و همکارانش روش های تخمین در حال ظهور مختلفی را که به عنوان روش های تکی، آن هم برای همه شرایط بهترین هستند، آنالیز کرده اند.

پژوهشگران دیگری در مقاله [۲۷] روی جنبه های کمی استفاده از روش های یادگیری ماشین، به منظور ساخت مدل هایی برای تخمین تلاش نرم افزار، بررسی هایی انجام و نتایج را تشریح کرده اند. روش های انتخابی در این مقاله، رگرسیون خطی چندگانه، رگرسیون LOGISTIC و CART هستند. رگرسیون خطی چندگانه به پیداکردن ارتباط بین متغیرهای وابسته و متغیرهای مستقل کمک می کند. این روش، نوعی از رگرسیون خطی است که به متغیرهای مستقل چندگانه می پردازد. در پایان این مقاله نتیجه گرفته شده است که هر روشی روی یک دادگان خوب عمل می کند.

همان پژوهشگران ایرانی که تأثیر کاهش انحراف توزیع تلاش را بررسی کرده اند، در مقاله ای دیگر [۳۰] روش منطق فازی در تخمین هزینه و تلاش نرم افزار را بررسی نموده اند. آن ها برای رسیدن به این هدف، به بررسی جایگاه استفاده از این روش در

هزینه نرم‌افزار پرداخته‌اند. در این مقاله بیان شده است که باینکه تاکنون متدهای تخمین زیادی ارائه شده است، هیچ‌کدام از آن‌ها نتیجه یکسانی روی همه دادگان ندارند. از سال ۱۹۸۰ تاکنون، متدهای زیادی در این زمینه ارائه شده و مانند مقالات قبلی این مقاله هم متدهای ارائه‌شده را در سه دسته مدل‌های الگوریتمی، مبتنی بر قضاوت کارشناسان و هوش محاسباتی که مقالات قبلی از آن با عنوان یادگیری ماشین یاد کرده بودند، دسته‌بندی کرده است. این مقاله گروه‌های همگن و ناهمگنی از ترکیب مدل‌های یادگیری ماشین برای تخمین تلاش توسعه نرم‌افزار ارائه داده است.

### ۳. دادگان مورد استفاده

تخمین هزینه نرم‌افزار در ابتدا و قبل از ارائه مدل‌های تخمین هزینه نرم‌افزار، با استفاده از یک سری قوانین سرانگشتی و برخی الگوریتم‌های ساده انجام می‌شده است. تا اینکه مدل تابع نقطه‌ای توسط آلبرت<sup>۵</sup> معرفی شد. این مدل در حال حاضر هم از درجه اطمینان فوق‌العاده‌ای برخوردار است. آنالیز تابع نقطه‌ای متدی است مرتبط با کمیت سایز و پیچیدگی یک سیستم نرم‌افزاری یا به اصطلاح، توابع مورد استفاده برای نوشتن کد مربوط به هر نمونه پروژه است که به کاربر تحویل داده می‌شود. این توابع به زبان برنامه‌نویسی مورد استفاده و یا ابزارهای استفاده‌شده برای توسعه پروژه‌های نرم‌افزاری وابسته نیستند. از آن زمان تاکنون، مطالعات زیادی برای تنظیم پارامترها و وزن‌های آن با هدف منعکس کردن پیشرفت‌های جاری در زمینه فناوری و صنعت توسعه نرم‌افزار صورت گرفته است. مدل‌هایی از این دست در دسته مدل‌های الگوریتمی قرار می‌گیرند.

پژوهشگران در [۳۲] به مرور مدل‌های الگوریتمی ارائه‌شده در این زمینه پرداخته‌اند. همان‌طور که در شکل (۳) می‌بینید، این دسته از متدها، خود به چهار دسته تقسیم‌بندی شده‌اند. مدل ارائه‌شده توسط بیلی و بیزلی<sup>۶</sup> را می‌توان نمونه‌ای از مدل‌های غیرخطی نام برد که در سال ۱۹۸۱ ارائه شده است [۳۳]. دسته دیگری از مدل‌های الگوریتمی دربرگیرنده مدل‌های گسسته

تخمین نرم‌افزاری با مطالعه ۶۴ مقاله از مقالات منتشره در سال‌های ۱۹۹۵ تا ۲۰۱۲ پرداخته و نتایج مقایسه‌ای خود را در [۳۰] ارائه داده‌اند.

پژوهش دیگری که در این حوزه صورت گرفته و منجر به مقاله پژوهشی [۳۱] شده است، مطالعه‌ای روی پیش‌بینی تلاش نرم‌افزار با استفاده از روش‌های یادگیری ماشین را نشان می‌دهد. پژوهشگران در مطالعات خود در این مقاله، دو دسته روش یادگیری نظارتی و غیرنظارتی را برای پیش‌بینی تلاش نرم‌افزار با استفاده از دادگان قدیمی معرفی کرده‌اند. پژوهشگران در این مقاله معتقدند که اگرچه مدل‌های یادگیری ماشین ممکن است بهترین راه‌حل را ارائه ندهند، می‌توان حداقل از آن‌ها برای تکمیل سایر مدل‌ها توسط مدیران استفاده کرد. تخمین دقیق تلاش توسعه نرم‌افزار، مخصوصاً در بازار به شدت رقابتی نرم‌افزار، قطعاً روی موفقیت پروژه نرم‌افزاری تأثیرگذار است.

محققان در مهندسی نرم‌افزار، بر این فرضیه تأکید دارند که پروژه‌های با خصوصیات مشابه، تلاش تقریباً معادلی خواهند داشت. در یادگیری غیرنظارتی، از نظر یادگیری ماشین، خوشه‌بندی پروژه‌های نرم‌افزاری براساس یک زیرمجموعه تصادفی می‌تواند اطلاعات مورد نیاز برای خصوصیات مشخص‌شده را فراهم کند. اگر ما تلاش را به‌عنوان یک خصوصیت در دادگان بشناسیم، می‌توان آن را با خوشه‌بندی پروژه‌ها براساس دیگر خصوصیات استنباط کرد. در این کار پژوهشی، از روش خوشه‌بندی K-MEDOID با سه معیار شباهت فاصله<sup>۱</sup>، جاکارد<sup>۲</sup> و ضرایب کالزینسکی<sup>۳</sup>، برای اندازه‌گیری شباهت بردارهای بولین که ارائه‌دهنده پروژه‌های نرم‌افزاری است، به‌عنوان نماینده‌ای برای یادگیری غیرنظارتی استفاده کرده‌اند. سه روش درخت تصمیم J48<sup>۴</sup>، بیز ساده و BPNN<sup>۴</sup> به‌عنوان نمایندگانی برای یادگیری نظارتی انتخاب و پیاده‌سازی شده‌اند [۳۱].

پژوهشگران در مقاله [۱۱]، به مطالعه آزمایشی روی ترکیبات همگن و ناهمگن روش‌های یادگیری ماشین در زمینه تخمین

1. Distance
2. Jaccard
3. Kulzinsky coefficients
4. Back Propagation neural network

5. albercht

6. John W. Bailey and Victor R. Basili

تعیین شده تعداد توابع کاربر (UFC<sup>۴</sup>) را تعیین می‌کند و به‌عنوان ورودی به مرحله بعد ارائه می‌شود.

ب. تنظیم پیچیدگی پردازش محیطی

تابع نقطه‌ای نهایی با ضرب UFC در یک عامل تنظیم که با در نظر گرفتن ۱۴ ویژگی سیستمی عمومی تعیین می‌گردد، محاسبه می‌شود. این ۱۴ ویژگی به‌همراه شرح آن‌ها در جدول (۴) آمده است. عامل تنظیم ارزش (VAF<sup>۵</sup>) با در نظر گرفتن این ۱۴ ویژگی که نرخ عمومی عملکرد برنامه‌ها را تعیین می‌کنند و با ارزشی که با مقدار C که تعیین‌کننده درجه نفوذ هر ویژگی است، مقارنه می‌شوند. درجه نفوذ با عددی در مقیاس صفر تا ۵ که تعیین‌کننده میزان تأثیر و نفوذ هر ویژگی از حالت بدون تأثیر تا تأثیر قوی است، مقارنه می‌شود. جدول (۴) جدولی است که برای معرفی ویژگی‌های عمومی سیستم توسط گروه کاربران عمومی تابع نقطه‌ای ارائه شده است. پس از پاسخگویی به این ۱۴ ویژگی، این مقادیر برای تعیین VAF مورد استفاده قرار می‌گیرند. عامل تنظیم ارزش با استفاده از رابطه (۵) برای هر ویژگی محاسبه می‌شود. در این رابطه C<sub>i</sub> درجه نفوذ هر ویژگی است و این عامل برای تعداد i ویژگی که معادل ۱۴ است، محاسبه می‌شود.

$$VAF = 0.65 + \left[ \frac{C_i}{100} \right] \quad (5)$$

در ادامه، عامل VAF نهایی از مجموع VAF تمام ویژگی‌ها حاصل می‌شود. در پایان کار، تابع نقطه‌ای از حاصل ضرب دو عامل تعیین شده در این دو مرحله، طبق رابطه (۶) محاسبه می‌شود.

$$FP = UAF * VAF \quad (6)$$

همان‌طور که گفته شد، تابع نقطه‌ای تعیین‌کننده کمیت سبزی و پیچیدگی نرم‌افزار است. این کمیت عامل مهمی در تعیین هزینه و تلاش مورد نیاز توسعه نرم‌افزار در مدل‌هایی است که در ادامه ارائه شده است. از جمله مدل‌هایی که از این مدل برای تعیین سبزی استفاده کرده‌اند، دو مدل ESTIMATS و SPQR/20 هستند.

تخمین هزینه نرم‌افزار هستند. مدل ولورتن<sup>۱</sup> نمونه‌ای از این دسته مدل‌های الگوریتمی است که ولورتن در ۱۹۷۴ ارائه کرده است [۳۴]. دسته سوم از مدل‌های الگوریتمی، مدل‌های تابعی نامیده می‌شوند. مدل پوتنم<sup>۲</sup> ارائه شده توسط پوتنم و همکارش در سال ۱۹۷۸ از دسته مدل‌های تابعی است [۳۵]. دسته آخر از مدل‌های الگوریتمی با عنوان مدل‌های چندگانه دسته‌بندی شده‌اند. مدل والستون فلیکس<sup>۳</sup> که والستون و فلیکس در سال ۱۹۹۷ ارائه کرده‌اند، از جمله مدل‌های الگوریتمی چندگانه محسوب می‌شود [۳۶].

ادامه ارائه مدل‌های تخمین هزینه نرم‌افزار، در دهه ۱۹۶۰ فرانک فریمن مفاهیم مدل‌های تخمین پارامتریک را توسعه داد که منجر به توسعه مدل قیمت برای سخت‌افزار شد. در ادامه و در دهه ۱۹۷۰، با استفاده از روش‌های آماری، مبادرت به شناخت ویژگی‌های مؤثر بر هزینه توسعه نرم‌افزار با استفاده از رگرسیون و وابستگی کردند.

مدل آنالیز تابع نقطه‌ای یکی از اولین مدل‌های ارائه شده در زمینه تخمین هزینه نرم‌افزار است که برای تعیین کمیت سبزی و پیچیدگی نرم‌افزار مورد استفاده قرار گرفته است. این مدل در ادامه پیشرفت‌های تخمین هزینه نرم‌افزار به‌عنوان یکی از عوامل هزینه هم مدنظر قرار گرفته است [۱۰].

تعداد توابع نقطه‌ای به انواع متمایزی از موارد از نظر فرمت و منطق وابسته است. تعیین تعداد توابع نقطه‌ای شامل دو مرحله است:

الف. شمارش توابع کاربر

تعداد توابع خام با در نظر گرفتن پنج جزء عمومی نرم‌افزار مشخص می‌شود. این اجزای نرم‌افزاری شامل ورودی‌های خارجی، خروجی‌های خارجی، استعلامات خارجی، فایل‌های داخلی منطقی و رابط خارجی اند که با سه سطح پیچیدگی ساده، متوسط و پیچیده، مقارنه می‌شوند. مقدار پیچیدگی تعیین شده برای هر جزء، تعیین‌کننده وزن هر جزء است. مجموع وزن

4. User function count  
5. value adjustment factor

1. Wolverton (TRW) Model  
2. Putnam  
3. Walston & Felix Model

با اتمام دهه هفتاد، مدل هزینه سودمند COCOMO را بوهم و همکارانش ارائه دادند. در دهه هشتاد، محققان نسبت به بهبود روش های ارائه شده در این زمینه اقدام کردند. با اتمام این دهه مدل پارامتریک قیمت توسعه داده شد و محققان شرکت IBM مدل تابع نقطه ای را برای تخمین هزینه و تلاش نرم افزار توسعه دادند. اندازه گیری سایز نرم افزار مبتنی بر تعداد عملگرها و عملوندها نیز در پایان این دهه توسط هالاستیت<sup>۱۵</sup> توسعه پیدا کرد. در انتهای این دهه، بوهم و همکارانش مدل COCOMO را که قبلاً ارائه داده بودند، به نام ADA COCOMO برای نرم افزارهای نوشته شده با زبان برنامه نویسی ADA ارائه کردند. کاپرس<sup>۱۶</sup> قابلیت آنالیز تابع نقطه ای را با افزودن اثر الگوریتم های پیچیده ای در محاسبه توابع نقطه ای بهبود بخشید. چارلز سیمون نیز در ادامه توسعه مدل های ارائه شده، نسخه ای از آنالیز تابع نقطه ای را برای حل مسائل فردیت در آنالیز تابع نقطه ای بهبود بخشید. در سال ۱۹۹۰، بوهم و همکارانش مدل خود را دوباره بهبود داده و مدل COCOMO II را ارائه کردند. این بهینه سازی همچنان ادامه دارد [۸].

به مرور زمان و با هدف افزایش دقت تخمین هزینه نرم افزار، امروزه روش های یادگیری ماشین به دلیل دقت خوبی که ارائه کردند، مورد توجه واقع شده اند. اما روال کاری این روش ها بر یادگیری الگو از نمونه های قدیمی متکی است و قطعاً نیازمند دادگانی هستند که شامل نمونه هایی تشریح شده با یک سری ویژگی است. مدل های الگوریتمی ارائه شده تا پیش از این، باعث ایجاد دادگانی مانند COCOMO شدند. دادگان ارائه شده برای مدل های الگوریتمی به عنوان منابعی برای یادگیری روش های یادگیری ماشین مورد استفاده قرار می گیرند. این روش ها از بردار ویژگی تشریح گر نمونه پروژه های دادگان، الگویی استخراج می کنند و با استفاده از این الگو، و مقادیر بردار ویژگی نمونه جدید نسبت به تخمین هزینه برای آن نمونه مبادرت می کنند.

مدل COCOMO یک مدل تخمین نرم افزار مبتنی بر رگرسیون است و در دسته مدل های الگوریتمی تابعی قرار می گیرد. این مدل در سال ۱۹۸۱ توسط بوهم و همکارانش توسعه یافت و مورد استانداردترین، معروف ترین و پذیرفته ترین مدل

جدول (۴): ویژگی های عمومی سیستم برای تعیین تابع نقطه ای [۳۷]

شماره ویژگی	عنوان ویژگی	شرح ویژگی
۱	ارتباطات داده ای <sup>۱</sup>	میزان امکانات ارتباطی برای کمک به انتقال یا تبادل اطلاعات توسط نرم افزار
۲	پردازش داده توزیع شده <sup>۲</sup>	میزان به کارگیری داده و توابع توزیع شده
۳	عملکرد <sup>۳</sup>	زمان پاسخگویی و توان مورد نیاز برای هر کاربر
۴	میزان سنگینی پیکربندی استفاده شده <sup>۴</sup>	میزان سختی پلت فرم سخت افزار استفاده شده برای اجرای نرم افزار
۵	نرخ تراکنش <sup>۵</sup>	میزان اجرای تراکنش ها (به صورت روزانه، هفتگی یا ماهیانه)
۶	ورود داده آنلاین <sup>۶</sup>	درصد اطلاعات وارد شده به صورت آنلاین
۷	بهره وری کاربر نهایی <sup>۷</sup>	میزان مؤثر بودن نرم افزار برای کاربر نهایی
۸	به روز رسانی آنلاین <sup>۸</sup>	میزان به روز رسانی فایل های منطقی خارجی با تراکنش های آنلاین
۹	پردازش های پیچیده <sup>۹</sup>	میزان وجود پردازش های پیچیده منطقی یا ریاضی
۱۰	قابلیت استفاده مجدد <sup>۱۰</sup>	میزان قابلیت توسعه نرم افزار برای یک یا چند کاربر
۱۱	سهولت نصب و راه اندازی <sup>۱۱</sup>	میزان سختی تبدیل یا نصب نرم افزار
۱۲	سهولت عملکردی <sup>۱۲</sup>	میزان مؤثر بودن و یا انجام راه اندازی خودکار، پشتیبان گیری و بازیابی فرایندها
۱۳	سایت های مختلف <sup>۱۳</sup>	میزان امکان طراحی، توسعه و پشتیبانی خاص نرم افزار برای سایت ها و سازمان های مختلف
۱۴	تسهیل در تغییرات <sup>۱۴</sup>	امکان طراحی، توسعه و پشتیبانی خاص برای تغییرات آسان نرم افزار

1. Data communication
2. Distributed data processing
3. Performance
4. Heavily used configuration
5. Transaction rate
6. On-Line data entry
7. End-user efficiency
8. On-Line update
9. Complex processing
10. Reusability
11. Installation ease
12. Operational ease
13. Multiple sites
14. Facilitate change

15. halasted

16. Capres



این دادگان ارائه شده است. مقادیر اختصاص یافته به هر ویژگی در محدوده very low تا extra high است [۲۹].

جدول (۵): شرح ویژگی‌های مدل COCOMO			
شماره	نام ویژگی	نام کامل انگلیسی	معادل فارسی
۱	RELY	required software reliability	قابلیت اعتماد مورد نیاز برای نرم افزار
۲	DATA	data base size	اندازه پایگاه داده‌ها (در آزمون نرم افزار)
۳	CPLX	process complexity	پیچیدگی محصول
۴	TIME	time constraint for cpu	محدودیت زمان اجرایی
۵	STOR	main memory constraint	محدودیت حافظه اصلی
۶	VIRT	machine volatility	تغییرپذیری ماشین
۷	TURN	turnaround time	مدت زمان بین شروع کار نرم افزار تا دریافت خروجی
۸	ACAP	analysts capability	توانایی تحلیلگرها
۹	AEXP	application experience	تجربه کاربردی
۱۰	PCAP	programmers capability	توانایی برنامه‌نویس‌ها
۱۱	VEXP	virtual machine experience	تجربه ماشین مجازی
۱۲	LEXP	language experience	تجربه زبان
۱۳	MODP	modern programing practices	تجربه‌های برنامه‌نویسی نوین
۱۴	TOOL	use of software tools	استفاده از ابزارهای نرم‌افزاری
۱۵	SCED	schedule constraint	زمان بندی و محدودیت‌های زمانی
۱۶	LOC	Line of code	تعداد خطوط کد
۱۷	ACTUAL	Effort	میزان تلاش مورد نیاز برای توسعه نرم افزار

در ادامه، کاربرد بردار ویژگی تعیین شده توسط بوهوم و همکارانش در مدل COCOMO، سازمان ناسا در سال ۲۰۰۴ دادگان COCONASA شامل ۶۰ پروژه از نرم‌افزارهای تولیدی توسط مراکز مختلف ناسا ارائه داد که توسط همین بردار ویژگی شرح داده شده‌اند. سپس در سال ۲۰۰۶، این سازمان یک نسخه جدید از دادگان خود را با نام COCONASA\_2 با ۹۳ نمونه پروژه ارائه داد. تفاوت این دادگان با نسخه قبلی در بردار ویژگی تعیین شده برای هر نمونه پروژه است. در این دادگان، علاوه بر ۱۷ ویژگی تعیین شده برای نمونه پروژه‌ها، ۷ ویژگی جدید به بردار ویژگی اضافه شده است. ویژگی‌های اضافه شده در جدول

الگوریتمی پیش‌بینی هزینه نرم‌افزار است. مدل COCOMO می‌تواند برای محاسبه میزان تلاش و زمان مورد نیاز برای پروژه‌های نرم‌افزاری استفاده شود. یکی از انواع مدل‌های COCOMO، مدل COCOMO81 است که زمان در آن یک مدل باثبات است. از مشکلات این مدل عدم تطابق آن با پیشرفت‌های محیطی توسعه نرم‌افزاری است. این مشکلی است که همه مدل‌های الگوریتمی با آن درگیرند. بنابراین در اواخر دهه ۹۰ مدل COCOMOII ارائه شد. این مدل دارای سه زیرمدل متفاوت از مدل COCOMO81 است: ۱. مدل ترکیبی کاربردی که برای انجام تخمین هزینه پروژه‌های ساخته شده با ابزارهای GUI مناسب است. ۲. مدل طراحی اولیه که برای ارائه تخمینی واقعی از هزینه یک پروژه به شکلی که قبل از تخمین معماری کلی نرم‌افزار تعیین می‌شود، مورد استفاده قرار می‌گیرد. این مدل از مجموعه کوچکی از تعیین‌کنندگان هزینه و یک سری معادلات جدید برای تخمین تشکیل شده است. ۳. مدل معماری-پسین که بعد از تعیین معماری کلی پروژه، از آن استفاده می‌شود. این مدل از یکی از دو روش تابع نقطه‌ای و تعداد خطوط (LOC<sup>۱</sup>) برای تخمین سایز که از عوامل اصلی تعیین‌کننده هزینه است، استفاده می‌کند و دربرگیرنده تلاش توسعه و نگهداری است.

این مدل شامل ۱۷ تعیین‌کننده هزینه است که روی مقیاسی از very low تا extra high به همان روش COCOMO81 رتبه‌بندی شده‌اند. مدل معماری-پسین با رابطه (۷) فرموله شده:

$$EFFORT = A * [size]^B * \prod_{i=1}^{17} EFFORT \text{ multiplier } i$$

where

$$B = 1.01 + 0.01 * \sum_{j=1}^5 scale \text{ factor } j$$

در این فرمول، پارامتر A ضریب ثابت است. size حجم پروژه نرم‌افزاری اندازه‌گیری شده با معیار ksloc (هزاران خط کد، تابع نقطه‌ای یا اشیا نقطه‌ای<sup>۲</sup>) است. انتخاب معیار مقیاس یا بزرگی براساس ارتباطاتی است که یک منبع قابل توجه از تغییرات نهایی در یک تلاش نرم‌افزار یا بهره‌وری پروژه است.

از آنجایی که این دادگان با وجود سیر بهبودی که در صنعت توسعه نرم‌افزار و مدل‌های تخمین باز هم شامل همان بردار ویژگی دادگان COCOMO81 است، در جدول (۵) بردار ویژگی

1. Line Of Code  
2. object point

پردازش داده واقعی، علوم پایه، شبیه سازی و کاربردی است. حالت توسعه پروژه ها به یکی از سه شکل تعبیه شده، ارگانیک و یا نیمه مجزاست. سیستم هوایی با حرف f و سیستم زمینی با حرف g در ستون FORG مشخص است.

(۶) شرح داده شده اند. در جدول (۶)، ویژگی کاربرد بیانگر نوع کاربرد پروژه در زمینه های ارتباطات هوایی، کاربردهای زمینی، نظارت هوایی، پردازش داده های دسته ای، ارتباطات، ضبط داده ها، پردازش پرتاب، برنامه نویسی پرواز، کنترل مانیتور، سیستم عامل،

جدول (۶): شرح ویژگی های اضافه بر مدل COCOMO در دادگان 2COCONASA

شماره	نام ویژگی	معادل فارسی	مقادیر
۱	UNIQUE ID	شماره واقعی نمونه	عدد از ۱ تا ۹۳
۲	PROJECTNAME	نام پروژه	de,erb,gal,X,hst,slp,spl,Y
۳	CAT2	کاربرد	Avionics, application_ground, avionicsmonitoring, batch data processing, communications, data capture, launchprocessing, missionplanning, monitor_control, operatingsystem, realdataprocessing, science, simulation, utility
۴	FORG	سیستم هوایی یا زمینی؟	f,g
۵	CENTER	کدام مرکز ناسا؟	۱,۲,۳,۴,۵,۶
۶	YEAR REAL	سال توسعه نرم افزار؟	مقدار عددی سال
۷	MODE	حالت توسعه؟	embedded,organic,semidetached

مجموعه دادگان حذف می کنند. در واقع این دادگان شامل ۷۷ پروژه نرم افزاری کامل و ۴ نمونه ناقص است. ویژگی های مستقل شامل تجربه گروه، تجربه مدیر، طول پروژه، معاملات، موجودیت ها، گروه تابع نقطه ای، محیط توسعه و زبان برنامه نویسی است. متغیر وابسته هم تلاش توسعه نرم افزار، اندازه گیری شده با ساعت در فرد است [۱۱]. ویژگی های این دادگان در جدول (۷) تشریح شده اند. این ویژگی ها کاملاً از ویژگی های دادگان COCOMO متفاوت اند.

دادگان پر کاربرد دیگری که در زمینه تخمین هزینه نرم افزار استفاده می شوند، دادگان DESHARNAIS است. این دادگان دربرگیرنده پروژه های تولیدی طی سال های ۱۹۸۱ تا ۱۹۸۸ یک خانه نرم افزاری است. نسخه اصلی این دادگان شامل ۸۱ نمونه پروژه توصیف شده توسط ۱۲ ویژگی از یک خانه نرم افزار کاندید است. اما ۴ نمونه از نمونه های این دادگان دارای مقدار گم شده در ۴ ویژگی هستند. محققان برخوردهای مختلفی با این دادگان برای استفاده دارند. دسته ای ویژگی های حاوی مقدار گم شده را کنار می گذارند و دسته ای نیز نمونه های حاوی مقدار گم شده را از

جدول (۷): شرح ویژگی های DESHARNAIS

شماره	نام ویژگی	معادل فارسی	مقادیر
۱	Project	شماره پروژه	عدد از ۱ تا ۸۱
۲	Exp Team	تجربه گروه	اندازه گیری شده با معیار سال
۳	Manager Exp	تجربه مدیر	اندازه گیری شده با معیار سال
۴	End Year	سال اتمام پروژه	سال
۵	Length	تعداد ماه مورد نیاز برای تولید پروژه	مقدار عددی
۶	Effort	تلاش مورد نیاز	شخص در ساعت
۷	Transactions	تعداد تراکنش های عمومی منطقی در سیستم	مقدار عددی
۸	Entities	تعداد موجودیت ها در سامانه های مدل داده	مقدار عددی
۹	Adjust Points	اندازه پروژه براساس توابع نقطه ای تعدیل نشده	این ویژگی به عنوان تراکنش ها به علاوه موجودیت محاسبه می شود
۱۰	Enver gure	تابع نقطه عامل تنظیم پیچیدگی	براساس ویژگی های عمومی سیستم (GSC) تعیین می شود. GSC شامل ۱۴ ویژگی است که هر کدام براساس یک مقیاس ترتیبی شش نقطه ای امتیازدهی شده اند.
۱۱	Points Non Ajust	اندازه پروژه	اندازه گیری شده با نقاط تابعی
۱۲	Language	زبان	پروژه بیان شده با مقادیر ۱، ۲ یا ۳. مقدار ۱ مربوط به «کوبول عمومی»، که در آن مقدار ۲ مربوط به «پیشرفته کوبول» و ارزش ۳ مربوط به زبان GL۴

خاص در این دسته قرار می‌گیرند.

#### • متغیرهای مربوط به توسعه

شامل اطلاعاتی درباره جنبه‌های مدیریتی و یا جنبه‌های روش توسعه پروژه‌های نرم‌افزاری، از جمله زبان برنامه‌نویسی یا نوع سیستم پایگاه داده‌ای که در طول توسعه مورد استفاده قرار گرفته است.

در ادامه، به تشریح اجمالی تعدادی از دادگان معروف که در اغلب کارهای پژوهشی در زمینه تخمین هزینه نرم‌افزار از آن‌ها استفاده می‌شود، می‌پردازیم [۱۱، ۱۳، ۱۵، ۳۱ و ۲۸].

#### ۳.۱. دادگان ALBRECHT

این دادگان شامل ۲۴ پروژه نرم‌افزاری توسعه‌یافته توسط سازمان سرویس‌های IBM DP است که دارای ۶ متغیر یا ورودی مستقل تعداد ورودی، تعداد خروجی، تعداد پرس‌وجو، تعداد فایل، توابع نقطه‌ای و خطوط کد است. ورودی وابسته آن‌هم تلاش اندازه‌گیری شده با تعداد ساعت کاری مورد نیاز برای طراحی، توسعه و آزمون هر نرم‌افزار است [۱۱].

#### ۳.۲. دادگان MIYAZAKI

این دادگان شامل ۴۸ پروژه نرم‌افزاری گردآوری شده توسط گروه کاربران سامانه‌های بزرگ فوجیتسو<sup>۲</sup> است که دارای ۷ خصوصیت مستقل شامل تعداد صفحه‌ها، تعداد فرم‌ها، تعداد فایل‌ها، تعداد اجزای داده‌ای در صفحه، تعداد اجزای داده‌ای در فرم‌ها، تعداد اجزای داده‌ای در فایل‌ها و خطوط کد است. همچنین متغیر وابسته تلاش توسعه نرم‌افزار اندازه‌گیری شده با تعداد فرد در ماه برای طراحی سیستم‌ها در آزمون سیستم‌ها، شامل تلاش مستقیم مثل مدیریت پروژه است. هر شخص در ماه معادل ۱۶۰ ساعت کار کردن است [۱۱].

#### ۳.۳. دادگان MAXWELL

شامل ۶۲ پروژه نرم‌افزاری از بزرگ‌ترین بانک‌های جهانی در فنلاند است که ۲۵ متغیر مستقل دارد که با ویژگی‌های نرم‌افزاری مختلف، مثل نوع کاربرد و سایز تعیین می‌شوند. متغیر وابسته آن هم تلاش توسعه نرم‌افزار تعیین شده با تعداد ساعت کار انجام شده

بنابراین روش‌های یادگیری ماشین، برای انجام تخمین هزینه نرم‌افزار هیچ‌گونه وابستگی به مدل‌های تخمین هزینه نرم‌افزار ارائه شده ندارند و تنها از دادگان آن‌ها برای یادگیری الگو استفاده می‌کنند. در جدول (۹)، دقت روش‌های یادگیری ماشینی روی دادگان مورد استفاده برای آموزش آن‌ها در این حوزه به دقت بالاتری نسبت به سایرین رسیده است. دادگان فراوانی تاکنون در این زمینه استفاده شده‌اند. برخی مطالعات ۳۱ دادگان را برای این زمینه معرفی کرده‌اند. اما نکته مهم این است که اکثر دادگان مورد استفاده، نسبتاً کوچک‌اند؛ مثلاً از ۳۱ دادگان معرفی شده در این مطالعات، فقط سه دادگان مربوط به بیش از ۵۰ پروژه‌اند. دادگان، عموماً شامل مجموعه‌ای از خصوصیات هستند که می‌توان آن‌ها را در دسته‌های زیر دسته‌بندی کرد [۱۵].

#### • خصوصیات مربوط به سایز

خصوصیاتی هستند که شامل اطلاعات مربوط به اندازه پروژه نرم‌افزاری هستند. این اطلاعات می‌توانند به شکل خطوط کد، تابع نقطه‌ای و یا خیلی از معیارهای دیگر فراهم شوند. متغیرهای بیانگر سایز، اغلب به عنوان ویژگی‌های مهمی برای تخمین تلاش در نظر گرفته می‌شوند. جدول (۸) هم دو روش تعداد خطوط و تابع نقطه‌ای را با هم مقایسه می‌کند.

جدول (۸): مقایسه دو روش تعیین سایز [۷]

تعداد خطوط	تابع نقطه‌ای
در اواخر چرخه پروژه	در هر گام از پروژه
تعداد خطوط چه چیزی است	روش‌های ساخت یافته
وابسته به فناوری	مستقل از فناوری
تعیین اندازه به سبک برنامه‌نویسی وابسته است	تعیین اندازه به سبک برنامه‌نویسی وابسته نیست

#### • اطلاعات محیطی

شامل اطلاعات پیش‌زمینه‌ای درباره گروه توسعه، شرکت و خود پروژه، مثل تعداد توسعه‌دهندگان و تجربه آن‌هاست.

#### • اطلاعات پروژه

شامل خصوصیات است که برای هدف مشخصی از پروژه و نوع پروژه بیان می‌شوند. همچنین ویژگی‌های مورد نیاز پروژه‌های

پژوهشگران در این کار به نتیجه واضحی نرسیدند و به همین دلیل، آن‌ها در مقاله‌ای که اخیراً در سال ۲۰۱۴ ارائه کرده‌اند [۳۹]، کار خود را ادامه داده‌اند.

به‌طور کلی، نتیجه‌ای که در این بررسی‌ها حاصل شد، این است که نیمی از شواهد آنالیز شده نشان داده‌اند که مدل‌های ساخته شده از داده‌های برون‌شرکتی به‌طور چشمگیری، از مدل‌های تخمین ساخته شده از داده‌های درون‌شرکتی بدتر نیستند. علاوه بر این، شواهدی وجود دارد که بیانگر این امر هستند که ساین بیشتر نمونه‌ها به معنای دقت بیشتر در تخمین نیست. نکته قابل ذکر دیگر این است که نمونه‌ها برای ساختن یک مدل تخمین، باید به دقت و براساس کنترل جنبه‌های کیفی و شباهت پروژه‌ها انتخاب شوند [۳۹].

پژوهشگران در [۲۵]، چالش اصلی در تخمین تلاش را، ارائه راهکاری برای پیشنهاد یک روش تخمین که روی همه دادگان دقت مشابهی داشته باشد، معرفی می‌کنند. روش ارائه شده در این مقاله، شامل دو گام اصلی است:

گام اول: یک متخصص، متغیرهای هر دادگان را بررسی کند و فقط آن‌هایی را نگه دارد که پتانسیل پیش‌بینی تلاش را دارند. بنابراین متغیرهایی مثل مدت زمان و تعداد نقص‌ها به دلیل اینکه در ادامه ناشناخته‌اند، حذف شده‌اند.

گام دوم: تغییر شکل داده به مقیاس جدید، مثل گسسته‌سازی را پوشش می‌دهد. این مرحله بدان سبب لازم است که بیشتر روش‌هایی که برای این آزمایش انتخاب شده‌اند، فقط با داده‌های طبقه‌بندی شده کار می‌کنند. روش‌های مختلفی مثل به دست آوردن بازه‌هایی با فرکانس یا عرض مساوی برای انجام این وظیفه وجود دارد. در این آزمایش، از یک متخصص برای گسسته‌سازی همه متغیرهای عددی با هدف تعیین بازه‌های با مقادیر معنی دار و فرکانس‌های مشابه استفاده شده است.

#### ۴. معیارهای ارزیابی

هرگونه کار پژوهشی در نهایت برای تعیین میزان دقت<sup>۴</sup> کار انجام شده، بایستی با یک سری از معیارها، ارزیابی شود.

توسط عرضه کننده نرم افزار، از مشخصات فنی تا زمان تحویل است [۱۱].

#### ۳.۴. دادگان CSBSG

شامل ۱۱۰۳ پروژه نرم‌افزاری است که در سال ۲۰۰۶ با مأموریت ترویج استانداردهای تولیدات نرم‌افزاری چینی ایجاد شده‌اند. پروژه‌های این دادگان از ۱۴۰ سازمان و ۱۵ منطقه چینی مرتبط با صنعت نرم‌افزار جمع‌آوری شده‌اند. هر پروژه نرم‌افزاری در این دادگان با ۱۷۹ خصوصیت تشریح شده‌اند [۳۱].

#### ۳.۵. دادگان ISBSG

شامل ۱۲۳۸ پروژه از سازمان‌های مختلفی همچون بیمه و سازمان‌های دولتی دیگر گردآوری شده از ۲۰ کشور مختلف است. در این دادگان، هر پروژه نرم‌افزاری با استفاده از ۷۰ خصوصیت تشریح شده است. در برخی از کارهای پژوهشی، این دادگان مورد پیش‌پردازش قرار می‌گیرند [۳۱].

علاوه بر دادگان تشریح شده در فوق، دادگان فراوان دیگری هم تاکنون در مطالعات مختلف برای آموزش الگوریتم‌های یادگیری ماشین استفاده شده‌اند. برخی از پژوهشگران در شرکت‌های مختلف، از دادگان درون‌شرکتی<sup>۱</sup> و برخی دیگر از دادگان برون‌شرکتی<sup>۲</sup> برای ساخت مدل‌هایی جهت تخمین هزینه توسعه نرم‌افزار توسط شرکت استفاده می‌کنند. پژوهشگران در بررسی‌هایی که در سال ۲۰۰۷ روی مطالعات در این زمینه انجام داده‌اند [۳۸]، آورده‌اند که از میان ۱۰ مقاله انتخاب شده در این زمینه، ۷ مقاله نتایج مستقلی ارائه داده و از بین این ۷ مقاله هم ۳ مقاله اثبات کرده‌اند که مدل‌های برون‌شرکتی به‌طور قابل توجهی از مدل‌های درون‌شرکتی بدترند؛ البته به دلیل اینکه روش‌های متفاوتی در مطالعات مختلف استفاده می‌شوند، آنالیز کامل مطالعات غیرممکن است. در نهایت، آن‌ها به این نتیجه رسیده‌اند که مطالعات انجام شده روی دادگان درون‌شرکتی با کمتر از ۲۰ نمونه پروژه که از روش اعتبارسنجی کنار گذاشتن یکی در هر بار<sup>۳</sup> استفاده کرده‌اند، همگی از مدل‌های برون‌شرکتی بهتر بوده‌اند.

1. Within-company data
2. Cross-company data
3. leave-one-out

واقعی هزینه نمونه پروژه‌های نرم‌افزاری انتخابی، MRE، به‌عنوان معیاری برای سنجش دقت الگوریتم ارائه‌شده استفاده می‌کنند [۴۰]. این میزان خطا با استفاده از رابطه (۵) محاسبه می‌شود.

$$MRE = \frac{|Actual\ Effort - Estimated\ Effort|}{Actual\ Effort} \quad (5)$$

Actual Effort میزان تلاش واقعی نمونه پروژه موردنظر در دادگان را دربرمی‌گیرد. Estimated Effort هم میزان تلاش تخمین‌زده‌شده توسط الگوریتم موردسنجش را نشان می‌دهد.

#### ۲.۴. میانگین خطای تخمین از میزان واقعی

در برخی دیگر از کارهای پژوهشی در این حوزه، آنچه به‌عنوان میزان دقت الگوریتم ارزیابی می‌شود، میانگین اختلاف هزینه تخمین‌زده‌شده توسط الگوریتم موردنظر با هزینه واقعی برای تمام نمونه‌های مورد نظر MMRE است [۱۶]. رابطه تعیین‌کننده این مورد به شکل رابطه (۶) است.

$$MMRE = \frac{1}{n} \sum_{i=1}^n \left( \frac{|Estimated(i) - Actual(i)|}{Actual(i)} \right) \quad (6)$$

n بیانگر تعداد پروژه‌های موردارزیابی، Estimated میزان هزینه تعیین شده توسط الگوریتم موردارزیابی و Actual بیانگر میزان هزینه یا تلاش واقعی است. بنابراین هرچه میزان MMRE کمتر باشد، یعنی میزان خطای تخمین الگوریتم کمتر بوده و همین نکته بیانگر دقت بهتر است.

#### ۳.۴. میانه خطای تخمین از میزان واقعی

در کارهای پژوهشی در حوزه تخمین تلاش یا هزینه نرم‌افزار، معیار دیگری که برای سنجش دقت الگوریتم‌های موردآزمایش، استفاده می‌شود، میانه خطای مقدار تعیین شده توسط الگوریتم موردنظر از مقدار واقعی تلاش نمونه‌های موردنظر یا همان MDMRE است [۲۲]. رابطه تعیین‌کننده MDMRE به شکل رابطه (۷) است.

$$MDMRE = \text{Median (MRE)} \quad (7)$$

#### ۴.۴. احتمال میزان مشخصی از خطا

در اغلب کارهای تحقیقاتی در زمینه موردبحث، معیار دیگری که برای ارزیابی میزان دقت الگوریتم‌های موردآزمایش استفاده

در این باره، بی‌شک در زمینه تخمین هزینه یا تلاش موردنیاز برای توسعه نرم‌افزار نیز، معیارهایی برای سنجش دقت الگوریتم‌های ارائه‌شده وجود دارد. در ادامه، به بیان معیارهای موجود برای این سنجش خواهیم پرداخت.

در این زمینه، معیارهای سنجش دقت زیادی تاکنون معرفی و استفاده شده‌اند. سه معیار سنجش دقت پرکاربرد براساس میزان خطای الگوریتم‌ها در اغلب مطالعات مورد استفاده قرار می‌گیرند: ۱. MRE که بیانگر میزان اختلاف هزینه تخمین‌زده‌شده توسط الگوریتم‌ها با هزینه واقعی است؛ ۲. MRE که به‌وسیله علامت انحصاری MMRE نشان داده می‌شود و بیانگر میانگین خطای تخمین برای کل نمونه‌های موردنظر (نمونه‌های آموزشی و یا نمونه‌های آزمون) است. ۳. PRED (x) که نشان‌دهنده درصد نمونه‌هایی است که میزان خطای تخمین آن‌ها کمتر یا مساوی با مقدار x است. البته در برخی از کارها میانه خطای تخمین هم با علامت MDMRE، به‌عنوان یکی دیگر از معیارهای سنجش تخمین انجام‌شده توسط الگوریتم موردنظر و میزان دقت آن مورد استفاده قرار می‌گیرد. علاوه بر این روش‌ها، روش‌های دیگری هم توسط برخی دیگر از مطالعات استفاده شده‌اند که در مرور این دسته از کارها، از آن‌ها در جدولی که در ضمیمه الف آمده، نام برده شده است. بیان جامع هر یک از این معیارهای سنجش دقت به‌صورت فرمول‌وار در ادامه این فصل آمده است. روش‌های یادگیری ماشینی، در مرحله آموزش، الگو را از ویژگی‌های نمونه‌های آموزشی استخراج می‌کنند. در مرحله بعدی، تمام ویژگی‌های نمونه‌های آزمون، به‌جز ویژگی تلاش را دریافت و براساس الگوی استخراج‌شده در مرحله آموزش، تلاش موردنیاز برای نمونه‌های آزمون را پیش‌بینی می‌کنند. با استفاده از این معیارها میزان خطای هر الگوریتم را تعیین می‌کنند.

#### ۴.۱. خطای تخمین از میزان واقعی

در برخی از کارهای انجام‌شده در زمینه موضوع مورد بحث، از تعیین میزان اختلاف بین میزان هزینه تخمین زده شده با میزان

1. Magnitude of Relative Error
2. Mean Magnitude of Relative Error
3. Prediction accuracy
4. Median Magnitude of Relative Error

محققان در پژوهش [۴۲]، تخمین هزینه نرم افزار را به عنوان یک مسئله چند هدفه در نظر گرفته و روشی نو برای ارائه یک معیار جدید از ترکیب معیارهای ارزیابی ارائه داده اند با این هدف که مدل یادگیری ای تولید شود تا همه معیارهای ارزیابی را بهینه سازد.

مطالعات صورت گرفته در [۴] و تحلیلی که روی نتایج حاصل از پیاده سازی مدل های یادگیری ماشین روی دادگان تخمین هزینه نرم افزار آمده است، به جز BN بقیه مدل های یادگیری MMRE در محدوده ۰/۳۵ تا ۰/۵۵ و PRED (۰/۲۵) در محدوده ۰/۴۵ تا ۰/۷۵ را به ارمغان می آورند. این نتایج نشان دهنده این است که دقت تخمین با استفاده از مدل های یادگیری ماشین در سطح قابل قبولی قرار دارد.

تحلیل ارائه شده در [۴] روی مقالات منتشر شده طی سال های ۱۹۹۱ تا ۲۰۱۰ صورت گرفته است. در جدول (۹)، دقت روش های یادگیری ماشین روی چهار دادگان معروف و پر کاربرد در این حوزه که توسط محققان در مقالاتی که طی سال های ۲۰۱۰ به بعد منتشر شده اند، قابل مشاهده است. همان طور که در این جدول می توان دید، روی دادگان COCOMO81 روش ANN از نظر دو معیار MMRE و PRED و روش DT از نظر معیار MDMRE بهترین نتایج را داشته اند. برای دادگان MAXWELL هم، روش های DT و SVR بهترین نتایج را داشته اند. روش CBR روی دادگان ISBSG بهترین دقت را داشته است. روی دادگان DESHARNAIS هم روش CBR از نظر MMRE و PRED بهترین دقت را ارائه داده، ضمن اینکه روش SVR هم از نظر معیار ارزیابی MDMRE روی این دادگان بهتر از سایرین عمل کرده است. همان طور که از جدول برمی آید، میانگین خطای تمام روش های استفاده شده در این مقالات، روی دادگان DESHARNAIS بهتر از سایر داده هاست.

می شود،  $PRED(x)$  یا احتمال این است که میزان خطای مقدار تخمین زده شده برای تمام نمونه های مورد ارزیابی کمتر یا مساوی مقدار  $X$  باشد [۱۲]. رابطه ای که برای محاسبه این احتمال به کار می رود، در رابطه (۸) آمده است.

$$PRED(x) = \frac{k}{n} \quad (۸)$$

$X$  بیانگر میزان اختلاف مورد نظر است که در اغلب کارهای پژوهشی با میزان ۰/۲۵ مقدار است.  $K$  بیانگر تعداد نمونه هایی است که اختلاف مقدار هزینه تخمین زده شده توسط الگوریتم مورد ارزیابی برای آن ها، با مقدار واقعی هزینه آن ها کوچک تر یا مساوی مقدار  $X$  باشد.  $n$  هم بیانگر تعداد کل نمونه های مورد ارزیابی است. بنابراین هرچه میزان  $PRED(0/25)$  بیشتر باشد، میزان خطای الگوریتم مورد ارزیابی هم کمتر است و هزینه تخمین زده شده برای تعداد بیشتری از نمونه های دادگان مورد ارزیابی، میزان خطای کمتر یا مساوی ۰/۲۵ دارند.

#### ۴.۵. معیار ترکیبی احتمال و میانگین خطا

معیار دیگری که در سال ۲۰۰۹ به عنوان معیاری برای ارزیابی دقت معرفی شده و در برخی کارهای پژوهشی مثل [۱۱] مورد استفاده قرار گرفته، تابع ارزیابی<sup>۱</sup> یا  $EF$  است که از ترکیب دو معیار  $PRED(0/25)$  و میانگین خطا با هم حاصل شده و از رابطه (۹) برای محاسبه آن استفاده می شود.

$$EF = \frac{PRED(0.25)}{1 + MMRE} \quad (۹)$$

نکته قابل ذکر این است که در بیشتر کارهای پژوهشی، از چندین روش برای ارزیابی کار خود استفاده می کنند و این نکته کاملاً مشخص است که الگوریتمی که دارای میانگین خطای کمتر و  $PRED$  بیشتر باشد، الگوریتم بهتری خواهد بود؛ زیرا دقت بیشتر و خطای کمتری را دارد.

در [۴۱] با در نظر گرفتن هر دو معیار  $PRED$  و  $MMRE$  دو روش پر کاربرد از روش های انتخاب ویژگی بسته بندی (انتخاب ویژگی روبه جلو و روبه عقب) روی دادگان زمینه تخمین هزینه نرم افزار پیاده سازی شد و ثابت شد که این ایده تأثیر مثبتی در راستای شناسایی ویژگی های مؤثر در این حوزه خواهد گذاشت.

جدول (۹): مرور دقت روش های یادگیری ماشین مورد استفاده در تخمین هزینه نرم افزار

دادگان	مرجع	سال	روش	MMIRE	MdMIRE	Pred(25)	دادگان	مرجع	سال	روش	MMIRE	MdMIRE	Pred(25)
COCOMO81	[۱۱]	۲۰۱۳	ANN	۲۵۰,۷۷	-	۲۳,۱	ISB SG	[۱۵]	۲۰۱۲	CBR	-	۵۸	۱۸,۸
	[۱۱]	۲۰۱۳	SVR	۸۰۳,۹۶	-	۷,۱		[۱۵]	۲۰۱۲	CBR	-	۵۲,۷	۲۲
	[۱۵]	۲۰۱۲	DT	-	۷۶,۴	۱۱,۱		[۱۵]	۲۰۱۲	CBR	-	۵۳,۶	۲۵,۸
	[۱۵]	۲۰۱۲	SVR	-	۱۰,۵	۱۲,۷		[۱۵]	۲۰۱۲	CBR	-	۵۷,۳	۲۵,۵
	[۱۵]	۲۰۱۲	ANN	-	۹۹	۱۵,۹		[۱۵]	۲۰۱۲	ANN	-	۶۱,۶	۱۶,۳
	[۱۵]	۲۰۱۲	CBR	-	۷۵,۴	۱۵,۹		میانگین			۶۴,۳	۲۱,۲۱۲	
Maxwell	[۱۵]	۲۰۱۲	CBR	-	۸۶,۱	۱۹	DESHARNAIS	[۱۱]	۲۰۱۳	ANN	۷۸,۲۴	-	۲۰,۰
	[۱۵]	۲۰۱۲	CBR	-	۸۵,۶	۱۱,۱		[۱۱]	۲۰۱۳	SVR	۷۳,۸۳	-	۳۳,۳
	[۱۵]	۲۰۱۲	CBR	-	۹۱,۳	۱۱,۱		[۱۵]	۲۰۱۲	SVR	-	۳۵,۶	۳۵,۸
	[۱۵]	۲۰۱۲	ANN	-	۸۹,۲	۹,۵۲		[۱۵]	۲۰۱۲	ANN	-	۳۶,۱	۳۵,۸
	میانگین			۵۲۷,۳۶	۸۸,۵	۱۳,۶۵۲		[۱۵]	۲۰۱۲	CBR	-	۴۷,۲	۲۳,۵
	[۱۵]	۲۰۱۲	DT	-	۴۸,۴	۳۴,۰۶		[۱۵]	۲۰۱۲	CBR	-	۴۵,۸	۲۷,۲
	[۱۵]	۲۰۱۲	SVR	-	۴۵,۷	۲۵,۸		[۱۵]	۲۰۱۲	CBR	-	۴۵,۷	۳۳,۳
	[۱۵]	۲۰۱۲	ANN	-	۶۵,۵	۱۹,۴		[۱۵]	۲۰۱۲	CBR	-	۴۰	۳۵,۸
	[۱۵]	۲۰۱۲	CBR	-	۵۳,۴	۲۲,۶		[۱۵]	۲۰۱۲	ANN	-	۱۰۰	۲,۴۷
	[۱۵]	۲۰۱۲	CBR	-	۶۱,۸	۱۷,۷۷		[۲۵]	۲۰۱۲	ANN	۵۱	-	۳۳,۶۴
ISBSG	[۱۵]	۲۰۱۲	CBR	-	۵۵,۴	۲۲,۶	DESHARNAIS	[۲۵]	۲۰۱۲	CBR	۷۰	-	۲۷,۵۳
	[۱۵]	۲۰۱۲	CBR	-	۵۶	۳۰,۶		[۲۵]	۲۰۱۲	DT	۶۳	-	۳۱,۲۸
	[۱۵]	۲۰۱۲	ANN	-	۱۰۰	۱۴,۵		[۲۵]	۲۰۱۲	CBR	۶۷	-	۳۴
	میانگین			-	۶۰,۷۷۵	۲۳,۴۱۶		[۲۵]	۲۰۱۲	CBR	۶۳	-	۳۳,۹۴
	[۱۵]	۲۰۱۲	DT	-	۵۵,۷	۲۲,۷		[۲۵]	۲۰۱۲	CBR	۶۳	-	۳۶,۵۱
	[۱۵]	۲۰۱۲	SVR	-	۶۰,۵	۲۴,۷		[۲۵]	۲۰۱۲	CBR	۳۱	-	۶۷,۵۸
میانگین معیارهای دقت روی کل دادگان			MMRE=146.8	۱۱۵	۱۳,۹	میانگین			۶۲,۲۳	۵۰,۰۵	۳۱,۹۷۸		
میانگین معیارهای دقت روی کل دادگان			MdMRE=66.41935	۱۱۵	۱۳,۹	میانگین معیارهای دقت روی کل دادگان			Pred=23.93333	۱۱۵	۱۳,۹	۳۱,۹۷۸	

#### ۴.۶. اعتبارسنجی الگوریتم ها

حاصل شده، به کار می رود) تقسیم می شود. نتیجه این ارزیابی فقط برای انتخاب بهترین بخش آموزشی مورد استفاده قرار می گیرد. آنچه مسلم است این است که نتیجه حاصل از ارزیابی دادگانی که در آموزش استفاده شده اند، به عنوان نتیجه ارزیابی کلی مدنظر قرار نمی گیرد بلکه آنچه به عنوان نتیجه ارزیابی دقت الگوریتم مدنظر است، دقت آن الگوریتم روی پیش بینی هزینه نمونه هایی است که در دسته دادگان آزمون قرار دارند.

روش دسته بندی به نام اعتبارسنجی n بخشی<sup>۵</sup>، بدین شرح است که کل داده ها را به صورت تصادفی در n دسته تقسیم کرده و هر بار یکی از این دسته ها را به عنوان داده آزمون در نظر گرفته

بر اساس بررسی های انجام شده، در کارهایی که در حوزه تخمین هزینه نرم افزار با الگوریتم های یادگیری ماشین انجام شده، دادگان مورد ارزیابی را به واسطه الگوریتم های اعتبارسنجی<sup>۱</sup> به دو دسته دادگان آموزشی<sup>۲</sup> و دادگان آزمون<sup>۳</sup> تقسیم کرده و از دادگان آموزشی، برای آموزش الگوریتم ها استفاده می کنند. البته شایان ذکر است که این داده آموزشی هم، خود به دو دسته آموزش و اعتبارسنج<sup>۴</sup> (که برای ارزیابی میزان آموزشی که از دادگان آموزشی

1. Cross-validation
2. Training set
3. Test set
4. Validate

5. n-fold

با به کار بردن این روش‌ها، چه به صورت تکی و چه به صورت ترکیبی، اعم از ترکیب با روش‌های یادگیری ماشین دیگر و یا ترکیب با سایر روش‌ها، به دقت‌های مورد توجهی دست یافته‌اند. با اینکه هنوز اتفاق نظری صورت نگرفته، شایان ذکر است که بر اساس مطالعات مقایسه‌ای صورت گرفته، دقت روش‌های یادگیری ماشین نسبت به سایر روش‌ها، برتری ملموسی داشته است.

در مطالعاتی که به مرور کارهای پژوهشی از ابتدا تا به امروز در این زمینه پرداخته‌اند، روش‌های استدلال مبتنی بر مورد، شبکه‌های بیزین، قوانین وابستگی، درخت تصمیم، الگوریتم ژنتیک، برنامه‌نویسی ژنتیک، رگرسیون بردار پشتیبان و شبکه‌های عصبی مصنوعی، ۸ روش یادگیری ماشینی هستند که تا به امروز در تخمین هزینه نرم‌افزار به کار برده شده‌اند و هر کدام دقت‌های مختلفی را ارائه داده‌اند. البته هشدار این مسئله ضروری است که هریک از روش‌های فوق، نیازمند برقراری شرایطی مناسب برای ارائه دقت دلخواه هستند؛ برای مثال، روش درخت تصمیم، علاوه بر اینکه نمی‌تواند با داده‌های گم‌شده کار کند، مستعد پیش‌پرازش روی دادگان آموزشی کوچک هم است. روش‌های شبکه عصبی مصنوعی و استدلال مبتنی بر مورد، هم نمی‌توانند روی دادگانی با داده‌های گم‌شده کار کنند هم اینکه برای کاربرد آن‌ها روی دادگان با داده‌های دسته‌ای باید ابتدا این دادگان را به دادگانی کمی تبدیل نمود، سپس از این دو روش برای تخمین هزینه نرم‌افزار استفاده کرد. بر اساس مشاهدات و بررسی‌های صورت گرفته، دو روش CBR و ANN به نسبت سایرین، پرتکرارترین روش‌ها محسوب می‌شوند.

به منظور آموزش الگوریتم یادگیری ماشین، نیازمند استفاده از دادگانی هستیم که در فصل سوم پژوهش حاضر، به معرفی و تشریح برخی از این دادگان پرداخته شده است. دادگان در ابتدای امر به دو دسته دادگان درون‌شرکتی و برون‌شرکتی تقسیم می‌شوند. پژوهشگرانی که به مقایسه تأثیر استفاده از این دو نوع دادگان پرداخته‌اند، پیشنهاد می‌دهند که تا حد امکان سعی شود از دادگان درون‌شرکتی برای آموزش الگوریتم‌ها استفاده گردد. در صورتی که این نوع از دادگان در دسترس نبودند، بهتر است

و n-1 دسته دیگر را به عنوان داده آموزشی در نظر می‌گیرند؛ این کار را به تعداد n بار تکرار کرده تا از بین این دسته‌بندی‌ها، دسته‌ای که خطای کمتری دارد برگزینند. در مقاله مروری ون و همکارانش [۴] روش‌های اعتبارسنجی HOLDOUT و LOOCV<sup>۱</sup> به عنوان دو روش دیگر برای اعتبارسنجی معرفی شده‌اند. تعیین اینکه چه نمونه‌هایی در دسته دادگان آموزشی قرار بگیرند و چه نمونه‌هایی در دسته دادگان آزمون واقع شوند، به انتخاب خود پژوهشگران برمی‌گردد.

برخی از پژوهشگران [۳۴] نمونه‌ها را فقط به دو بخش آموزشی و آزمونی تقسیم می‌کنند. برخی دیگر [۲۷] هم درصد مشخصی را دیتای آموزشی و بقیه را دیتای آزمون در نظر گرفته و درصد مشخصی هم از دادگان آموزشی را به عنوان اعتبارسنج مورد استفاده قرار می‌دهند؛ البته داده‌ای را به عنوان داده آموزشی در نظر می‌گیرند که میزان خطای آموزشی آن کمتر از حد مشخصی باشد.

## ۵. نتیجه‌گیری

در این کار پژوهشی، به بررسی اندکی از کارهای انجام‌شده در زمینه تخمین هزینه نرم‌افزار پرداختیم. در زمینه تخمین هزینه نرم‌افزار از آغاز راه تاکنون، کارهای فراوانی انجام شده که نیازمند پژوهش بیش‌ازپیش در این زمینه است. در پژوهش حاضر، تنها موفق به مرور قطره‌ای از دریای گسترده مطالعات در این زمینه شدیم. اما آنچه مسلم است مطالعات کاربردی در این زمینه همچنان ادامه دارد. بر اساس مطالعات صورت گرفته از سال ۱۹۴۰، مفهوم تخمین هزینه نرم‌افزار معرفی شده و به مرور، روش‌های مختلفی توسط پژوهشگران این حوزه ارائه گردیده است. اکثر مقالات مثل [۷] روش‌های تخمین تلاش را در دو دسته الگوریتمی و غیرالگوریتمی تقسیم‌بندی کرده‌اند. روش‌های یادگیری ماشین هم در دسته مدل‌های غیرالگوریتمی قرار دارند. به دلیل شرایط ذاتی عمل تخمین، یعنی غیرقطعی بودن آن و همچنین هدف جامعه پژوهشی که افزایش دقت در تخمین است، به مرور توجه پژوهشگران به کاربرد روش‌های یادگیری ماشین در تخمین هزینه و تلاش نرم‌افزار جلب شده است. آن‌ها تا به امروز

1. Leave-one-out



هزینه نرم‌افزار پرداخته شده است. در [۴۵] مروری بر چهار روش بهینه‌سازی، منطق فازی، شبکه‌های هوشمند عصبی و الگوریتم‌های دسته‌بندی که برای افزایش دقت تخمین مدل COCOMO بر روی آن اعمال شده، آمده است.

### ۵. ۱. پژوهش‌های پیشنهادی

بر اساس مطالعات صورت گرفته می‌توان به ایده‌هایی در جهت ارائه روش‌هایی نو برای تخمین هزینه نرم‌افزار، آن‌هم به صورت دقیق تر دست یافت. برای پیشنهاد کارهای پیشرو می‌توان به موارد زیر اشاره کرد:

روش‌های یادگیری ماشین نتایج امیدوارکننده‌ای در تخمین هزینه نرم‌افزار ارائه داده‌اند. با این حال استفاده از روش‌های یادگیری ماشین در صنعت، هنوز گسترش چندانی نداشته است. بنابراین باید تلاش بیشتری در جهت تشویق صنعت و صنعتگران، برای استفاده از این روش‌ها در حوزه تخمین هزینه نرم‌افزار و تسهیل استفاده از آن‌ها صورت پذیرد. برای رسیدن به این هدف باید مطالعات آزمایشگاهی، مروری و جزئی‌تری روی این حوزه صورت پذیرد تا اطلاعات کامل‌تری از این روش‌ها در اختیار محققان تخمین هزینه در حوزه صنعت قرار گیرند.

استفاده از الگوریتم‌های انتخاب ویژگی‌هایی مثل انتخاب ویژگی روبه‌جلو و انتخاب ویژگی روبه‌عقب در ترکیب با روش‌های یادگیری ماشین استفاده شده، به ندرت صورت گرفته است. اما به نظر می‌رسد استفاده از این نوع ترکیبات در افزایش دقت تخمین مؤثر واقع گردد.

در پژوهش‌های صورت گرفته در زمینه تخمین هزینه نرم‌افزار از ابتدا تاکنون بیشتر از ۸ روش یادگیری ماشین استفاده شده است. استفاده از الگوریتم‌هایی که این دسته از روش‌ها را با استفاده از تغییراتی جزئی به سمت دقت بیشتر برده‌اند، مثل نسخه‌های جدید الگوریتم SVR، احتمالاً دقت آن‌ها را در این زمینه هم بالاتر ببرند.

استفاده از سایر الگوریتم‌های یادگیری ماشین، که تاکنون در این حوزه مورد استفاده قرار نگرفته‌اند هم می‌تواند به ارائه تخمینی دقیق منجر شود. مطالعات مورد بررسی در این پژوهش،

دادگان برون‌شرکتی که بیشترین شباهت را به پروژه‌ای که قرار است هزینه مورد نیاز جهت توسعه آن تخمین زده شود، انتخاب گردد. دادگانی مثل COCOMO، ISBSG و TUKUTUKU در دسته دادگان برون‌شرکتی محسوب می‌شوند و دادگانی چون NASSA و ALBRECHT، DESHARNAIS از جمله دادگان درون‌شرکتی هستند [۳۹]. شایان توجه است که در انتخاب دادگان، برای آموزش الگوریتم مورد استفاده، علاوه بر رعایت بستر مطلوب آن الگوریتم، باید جنبه‌های کیفی و شباهت پروژه‌ها را نیز در نظر بگیریم.

برای ارزیابی عملکرد روش‌های به کار گرفته شده در تخمین هزینه نرم‌افزار، باید میزان دقت تخمین ارائه شده سنجیده شود. برای نیل به این هدف، معیارهای ارزیابی گوناگونی تا به امروز در این حوزه استفاده شده‌اند. به منظور سنجش دقت الگوریتم ارائه شده در مطالعات این حوزه، ابتدا دادگان مورد استفاده توسط الگوریتم‌های اعتبارسنجی مثل اعتبارسنجی متقابل<sup>۱</sup>، به دو دسته دادگان آموزشی و آزمون تقسیم شده و هر بار با بررسی میزان دقت آموزش بر اساس یک یا چند نمونه از معیارهای ارزیابی دقت، بهترین دادگان آموزشی برگزیده شده است. سپس مقدار هزینه مورد نیاز برای نمونه‌هایی که در دسته دادگان آزمون واقع شده‌اند، تخمین زده شده و برای ارزیابی دقت الگوریتم، آنچه مدنظر قرار می‌گیرد، دقت سنجیده شده الگوریتم مورد نظر روی دادگان آزمون است. معیارهایی مثل MRE، MMRE، MDMRE، PRED(X) از جمله معیارهای مورد استفاده برای سنجش دقت الگوریتم‌های یادگیری ماشین در تخمین هزینه نرم‌افزارند.

در این زمینه، تاکنون تحقیقات زیادی انجام شده است. کتاب [۴۳] به بررسی زمینه تخمین تلاش یا هزینه نرم‌افزار پرداخته است و دید جامعی در این زمینه برای مخاطبان فراهم می‌کند. این کتاب انواع مدل‌های تخمین را معرفی کرده و روی دسته خاصی از این مدل‌ها تأکید ندارد. البته برخی مقالات هم سعی در ترکیب روش‌های یادگیری ماشین با روش‌های دیگر برای افزایش دقت تخمین هزینه نرم‌افزار بهره جسته‌اند. در [۴۴] به بررسی دقت تنها تعدادی از روش‌های یادگیری ماشین در حوزه تخمین

دارند تا با ترکیب روش خوشه بندی با دیگر روش های هوشمند، دقت تخمین را افزایش دهند.

عبدالله زید و همکارانش در [۸] که به بیان موضوعات داغ در زمینه تخمین تلاش نرم افزار پرداخته اند، عناوینی مثل تبدیل معیار، ساین بندی فاز طراحی و نیازمندی ها، آنالیز پیچیدگی سیستم های نرم افزاری، آنالیز ROI، بررسی OSS<sup>V</sup> و بهاگذاری مجدد در تابع نقطه ای را عناوینی می دانند که به تمرکز بیشتر محققان در این زمینه نیاز دارند. آن ها در پیشنهادی که برای کارهای پیش رو دارند، روی تحقیقات بیشتر درباره تأثیر موضوعات داغ مفید در این زمینه برای گشودن درهایی به روی محققان و تلاش های متمرکز محققان برای بهینه کردن تخمین تلاش نرم افزار تأکید دارند.

در انتهای این بررسی، ذکر توصیه ها و دستورالعمل هایی برای متخصصان، به منظور اجرای مدل های ML در صنعت خالی از لطف نیست [۴].

- درک منطقی از مدل های ML و همکاری با محققان باتجربه در اجرای مدل های ML.
- انتخاب آن دسته از مدل های ML که نقاط قوتی مطابق با زمینه ها و اهداف تخمین دارند.
- اولویت بندی دادگان درون شرکتی در زمان ساخت مدل های ML، در صورتی که تعداد پروژه ها در دادگان کافی است.
- استفاده از مدل های ML به صورت موازی با مدل های معمولی موجود در مراحل اولیه استفاده از مدل های ML
- به اشتراک گذاری دادگان پروژه های اختصاصی با جامعه پژوهشی.

پیشنهاد های مختلفی برای مطالعات و پژوهش های آتی در این حوزه ارائه داده اند که در ادامه به بررسی برخی از آن ها خواهیم پرداخت.

بنالا و همکارانش در [۲۲] برای انجام کارهای آینده در این زمینه، پیشنهاد داده اند از آنجایی که استفاده از الگوریتم های ترکیبی به تازگی ارائه شده و کار چندانی بر روی آن صورت نگرفته است و کارایی الگوریتم ارائه شده در پژوهش آن ها نیز نتایج بهینه ای را داشته است، به نظر می رسد می توان از سایر الگوریتم های یادگیری ماشین نیز برای الگوریتم های ترکیبی استفاده کرد. الگوریتم های پیشنهادی در این مقاله شامل بر کلونی زنبور عسل مصنوعی<sup>۱</sup>، تکاملی افتراقی<sup>۲</sup>، سیستم ایمنی مصنوعی<sup>۳</sup>، بهینه سازی جست و جوی غذا با کتریایی<sup>۴</sup>، فازی عصبی هستند.

راندینسکی<sup>۵</sup> و همکارش در [۲۵] با مطالعه ای که روی روش های یادگیری ماشین با داده های محلی انجام داده اند، به نتایج خوبی رسیده اند. این مطالعه فقط روی داده های محلی بود؛ یعنی داده های یک کمپانی را مورد استفاده قرار داده اند. بنابراین دادگان آن ها بزرگ نبودند؛ چرا که یک شرکت پروژه های زیادی را ایجاد نمی کند. آن ها در آینده قصد دارند مطالعه خود را روی داده های برون شرکتی انجام دهند.

خطیبی بردسیری و همکارانش در [۲۵] که با ارائه یک الگوریتم ترکیبی با ادغام خوشه بندی، شبکه عصبی و ABE در زمینه تخمین تلاش نرم افزار، به دقت خوبی رسیدند، در بیان کارهای پیش رو آورده اند که به هر حال با اینکه مدل پیشنهادی تخمین متقاعدکننده ای بر روی دو دیتاست ارائه می دهد، این روش با یک سری از محدودیت ها روبه روست. استفاده از این مدل هنگامی که دیتاست سازگار و نرمال باشد، مناسب نیست؛ چون عملکرد آن هنگامی که پروژه های درون دادگان زیاد باشد، شبیه شبکه های عصبی و پروژه های کم شبیه ABE است. و استفاده از آن وقت تلف کردن است. بنابراین قبل از استفاده از مدل پیشنهادی باید موقعیت دادگان، به دقت بررسی شود. آن ها قصد

1. Artificial Bee Colony  
 2. Differential Evolution  
 3. Artificial Immune System  
 4. Bacterial foraging optimization  
 5. Lukasz Radlinski

6. Return Of Investment  
 7. Open Source Software

ضمیمه الف. شرح اجمالی روش‌های استفاده‌شده در مقالات اخیر

ردیف	مقاله	سال نشر	نوع نشر	روش	دادگان	معیار سنجش دقت
۱	[۱]	۲۰۱۳	ژورنالی	ANN	COCOMO	-
۲	[۲]	۲۰۱۳	کنفرانس	کاهش انحراف تلاش با مدل COCOMOII	ISBSG	-
۳	[۵]	۲۰۱۲	ژورنالی	SVR-MPSOMODEL-INTERMEDIATE COCOMO	COCOMO	MARE-PRED
۴	[۹]	۲۰۱۲	کنفرانس	ANN-LINEAR REGRESSION-MP5	COCOMO81	MMRE-PRED
۵	[۱۱]	۲۰۱۳	ژورنالی	MLP-SVR-ANFIS-PCA-GA	COCOMO-ALBRECHT-MIAZAKI-MAXWELL-DESHARNAIS	MMRE-PRED-EF
۶	[۱۲]	۲۰۱۱	ژورنالی	RBNN-GP-FAZY LOGIC-PSOA	NASSA	MMRE-PRED
۷	[۱۳]	۲۰۰۷	ژورنالی	SVM و DT، PCA، RBF، MLP	Nassa-ucs	MMRE-PRED
۸	[۱۴]	۲۰۱۲	ژورنالی	OLS-OLS+LOG-OLS+BC-RUBOST REGRESSION-LMS-MARS-RIDGE REGRESSION-CART-M5-MLP-RBFN-CBR-LS SVM-	COCOMO81-CSC-DESHARNAIS	MRE-MMRE-MDMRE-PRED-
۹	[۲۰]	۲۰۱۰	ژورنالی	AODE-BN-NB-MLP-RBFN-SMO-KNN-K*-LBR-.....	QQDefects-Maxwell-DESHARNAIS-COCOMO	CONFIDENCE
۱۰	[۲۲]	۲۰۱۲	ژورنالی	GA+ ANN	COCONASSA2	MRE-MMRE-PRED
۱۱	[۲۳]	۲۰۱۲	ژورنالی	MULTIPLE LINEAR REGRESSION-GP	دادگان ایجادشده توسط پژوهشگران آن	MRE-MER-MMRE-MMER
۱۲	[۲۴]	۲۰۰۸	ژورنالی	AR+CART	STTF-ISBSG-OTHER APPLICATION	CONFIDENCE
۱۳	[۲۵]	۲۰۱۲	ژورنالی	ANN+ABE+CLUSTERING	MAXWELL-DESHARNAIS	RE-MRE-MMRE-PRED
۱۴	[۲۶]	۲۰۰۸	ژورنالی	ANN+MIN-MAX	COCOMO81	MMRE
۱۵	[۲۷]	۲۰۱۴	ژورنالی	CART-Multipel Regression-Logestic Regression, MIN-MAX	COCOMO81-NASSA93-Baily baisili	MMRE-MDMRE-PRED
۱۶	[۲۹]	۲۰۱۰	ژورنالی	ANN	COCOMO	MMRE-PRED
۱۷	[۳۱]	۲۰۱۳	ژورنالی	NB-J48-BPNN-CLUSTERING	ISBSG-CSBSG	CONFIDENCE
۱۸	[۳۴]	۲۰۱۲	ژورنالی	ANN، COCOMO MODEL	COCOMO	MRE

- [1] Sarwar, S., Gupta, M., "Proposing effort estimation of COCOMO II through perceptron learning rule", *International Journal of Computer Applications*, 70(1), pp. 29-32, 2013.
- [۲] دولت‌آبادی نژاد، ابوذر؛ سیدمهران شرفی و ناصر نعمت‌بخش، «بررسی تأثیر کاهش انحراف توزیع تلاش بر افزایش دقت تخمین تلاش در توسعه محصول نرم‌افزاری»، پنجمین کنفرانس ملی مهندسی برق و الکترونیک ایران، گناباد، ۱۳۹۲.
- [3] Grewal, A. S., Gupta, V., Kumar, R. Emerging, estimation techniques. *International Journal of Computer Applications*, Vol. 59, No. 8, pp. 30-34, 2012.
- [4] Wen, J., Li, S., Lin, Z., Hu, Y., & Huang, C. (2012). Systematic literature review of machine learning based software development effort estimation models. *Information and Software Technology*, 54(1), pp. 41-59.
- [5] Kumari, S., & Pushkar, S. (2013). Comparison and analysis of different software cost estimation methods. *IJACSA International Journal of Advanced Computer Science and application*, 4(1), pp. 153-157.
- [6] Pandey, P. (2013, April). Analysis of the techniques for software cost estimation. In *Advanced Computing and Communication Technologies (ACCT), 2013 Third International Conference on* (pp. 16-19). IEEE.
- [7] Bardsiri, A. K., & Hashemi, S. M. (2014). Software effort estimation: A survey of well-known approaches. *International Journal of Computer Science Engineering (IJCSSE)*, 3(1), pp. 46-50.
- [8] Zaid, A., Selamat, M. H., Ghani, A. A. A., Atan, R., & Wei, K. T. (2008). Issues in software cost estimation. *IJCSNS Int J of Computer Science and Network Security*, 8(11), pp. 350-356.
- [۹] اخوان بیطرف، امیر؛ مونا یخچی و بهناز محجویی، «ارائه روشی بهینه برای تخمین هزینه نرم‌افزار مبتنی بر شبکه عصبی پرسپترون به کمک داده‌کاوی اطلاعات نرم‌افزار»، دومین کنفرانس ملی محاسبات نرم و فناوری اطلاعات، ماهشهر، ۱۳۹۰.
- [10] Keim, Y., Bhardwaj, M., Saroop, S., & Tandon, A. (2014, January). Software cost estimation models and techniques: A survey. In *International Journal of Engineering Research and Technology*, 3(2), pp.1763-1768).
- [11] Elish, M. O., Helmy, T., & Hussain, M. I. (2013). Empirical study of homogeneous and heterogeneous ensemble models for software development effort estimation. *Mathematical Problems in Engineering*, 2013, pp. 1-21.
- [12] Sehra, S. K., Brar, Y. S., & Kaur, N. (2011). Soft computing techniques for software project effort estimation. *International Journal of Advanced Computer and Mathematical Sciences*, 2(3), pp. 160-167.
- [13] Baskales, B., Turhan, B., & Bener, A. (2007, November). Software effort estimation using machine learning methods. In *Computer and information sciences, 2007. iscis 2007. 22nd international symposium on* (pp. 1-6). IEEE.
- [14] Alpaydin, E. (2014). *Introduction to machine learning*. MIT press.
- [15] Dejaeger, K., Verbeke, W., Martens, D., & Baesens, B. (2012). Data mining techniques for software effort estimation: A comparative study. *IEEE Transactions on Software Engineering*, 38(2), pp. 375-397.
- [16] Hamza, H., Kamel, A., & Shams, K. (2013, April). Software effort estimation using artificial neural networks: A survey of the current practices. In *Information Technology: New Generations (ITNG), 2013 Tenth International Conference on* (pp. 731-733). IEEE.
- [17] Han, J., Pei, J., & Kamber, M. (2011). *Data mining: Concepts and techniques*. Elsevier.
- [18] "Bayesian network - Wikipedia, the free encyclopedia." pp. 1, 2011.
- [19] Radlinski, L. (2010). A survey of bayesian net models for software development effort prediction. *International Journal of Software Engineering and Computing*, 2(2), pp.95-109.
- [20] Radlinski, Lukasz, and Wladyslaw Hoffmann. (2010). "On predicting software development effort using machine learning techniques and local data." *International Journal of Software Engineering and Computing* 2(2), pp. 123-136.
- [21] "Genetic algorithm - Wikipedia, the free encyclopedia." pp. 1, 2011.
- [22] Benala, T. R., Dehuri, S., Satapathy, S. C., & Madhurakshara, S. (2012). Genetic algorithm for optimizing functional link artificial neural network based software cost estimation. In *Proceedings of the International Conference on Information Systems Design and Intelligent Applications 2012 (INDIA 2012) held in Visakhapatnam, India, January 2012* (pp. 75-82). Springer Berlin Heidelberg.

- [23] Chavoya, A., Lopez-Martin, C., Andalon-Garcia, I. R., & Meda-Campaña, M. E. (2012). Genetic programming as alternative for predicting development effort of individual software projects. *PloS one*, 7(11), e50531.
- [24] Bibi, S., Stamelos, I., & Angelis, L. (2008). Combining probabilistic models for explanatory productivity estimation. *Information and Software Technology*, 50(7), pp. 656-669.
- [25] Bardsiri, V. K., Jawawi, D. N. A., Hashim, S. Z. M., & Khatibi, E. (2012). Increasing the accuracy of software development effort estimation using projects clustering. *IET software*, 6(6), pp. 461-473.
- [26] Ayyıldız, M., Kalıpsız, O., & Yavuz, S. (2008). A metric-Set and model suggestion for better software project cost estimation. *International Journal of Computer, Information, Systems and Control Engineering*, 2(11), pp. 262-267.
- [27] Benala, T. R., Mall, R., Srikavya, P., & HariPriya, M. V. (2014). Software effort estimation using data mining techniques. In *ICT and Critical Infrastructure: Proceedings of the 48th Annual Convention of Computer Society of India-Vol I* (pp. 85-92). Springer International Publishing.
- [28] Kocaguneli, E., Menzies, T., Keung, J., Cok, D., & Madachy, R. (2013). Active learning and effort estimation: Finding the essential content of software effort estimation data. *IEEE Transactions on Software Engineering*, 39(8), pp. 1040-1053.
- [29] Attarzadeh, I., & Ow, S. H. (2010, April). Proposing a new software cost estimation model based on artificial neural networks. In *Computer Engineering and Technology (ICCET), 2010 2nd International Conference on* (Vol. 3, pp. V3-487). IEEE.
- [۳۰] دولت‌آبادی‌نژاد، ابوذر؛ سید مهران شرفی و ناصر نعمت‌بخش، «بررسی رویکردهای استفاده از منطق فازی در تخمین هزینه و تلاش نرم‌افزار»، همایش ملی علوم و مهندسی کامپیوتر، نجف‌آباد، دانشگاه آزاد اسلامی واحد نجف‌آباد، ۱۳۹۱
- [31] Zhang, W., Yang, Y., & Wang, Q. (2011, June). A study on software effort prediction using machine learning techniques. In *International Conference on Evaluation of Novel Approaches to Software Engineering* (pp. 1-15). Springer Berlin Heidelberg.
- [۳۲] ایرانمنش، الناز؛ وحید خطیبی بردسیری، «کاربرد انواع مدل‌های الگوریتمی در تخمین تلاش پروژه‌های نرم‌افزاری»، همایش ملی پژوهش‌های مهندسی رایانه‌ای، تهران، ۱۳۹۳.
- [33] Bailey, J. W., & Basili, V. R. (1981, March). A meta-model for software development resource expenditures. In *Proceedings of the 5th international conference on Software engineering* (pp. 107-116). IEEE Press.
- [34] Wolverton, R. W. (1974). The cost of developing large-scale software. *IEEE Transactions on Computers*, 100(6), pp. 615-636.
- [35] Putnam, L. H. (1978). A general empirical solution to the macro software sizing and estimating problem. *IEEE transactions on Software Engineering*, 4(4), pp. 345.-361.
- [36] Abbas, S. A., Liao, X., Rehman, A. U., Azam, A., & Abdullah, M. I. (2012). Cost estimation: A survey of well-known historic cost estimation techniques. *Journal of Emerging Trends in Computing and Information Sciences*, 3(4), pp. 612-636.
- [37] <http://www.softwaremetrics.com/fpafund.htm>, access in 21/12/2013.
- [38] Kitchenham, B. A., Mendes, E., & Travassos, G. H. (2007). Cross versus within-company cost estimation studies: A systematic review. *IEEE Transactions on Software Engineering*, 33(5), pp. 316-329.
- [39] Mendes, E., Kalinowski, M., Martins, D., Ferrucci, F., & Sarro, F. (2014, May). Cross-vs. within-company cost estimation studies revisited: an extended systematic review. In *Proceedings of the 18th ACM International Conference on Evaluation and Assessment in Software Engineering*, p. 12.
- [40] Kaushik, A., Chauhan, A., Mittal, D., & Gupta, S. (2012). COCOMO estimates using neural networks. *International Journal of Intelligent Systems and Applications*, 4(9), pp. 22-28.
- [۴۱] بیرانوند، صبا؛ محمدعلی زارع چاهوکی، «افزایش دقت در تخمین هزینه نرم‌افزار، با انتخاب زیرمجموعه‌ای از ویژگی‌ها مبتنی بر تلفیق معیارهای ارزیابی»، هشتمین کنفرانس داده‌کاوی ایران، تهران، ۱۳۹۳
- [42] Minku, L. L., & Yao, X. (2013). Software effort estimation as a multiobjective learning problem. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 22(4), pp. 1-32.
- [43] Trendowicz, Adam, and Ross Jeffery. Software project effort estimation. DOI 10.1007/978-3-319-03629-8,# Springer International Publishing Switzerland 2014 451, 2014.
- [44] Srinivasan, K., & Fisher, D. (1995). Machine learning approaches to estimating software development effort. *IEEE Transactions on Software Engineering*, 21(2), pp. 126-137.

## ارائه یک الگوریتم بهبود یافته وب کاوی برای وب معنایی ۶۵

[۴۵] جهانشاهی، سعیده؛ وحید خطیبی بردسیری؛ رضوان حسنی سعدی و معصومه مقبلی، «کاربرد روش های محاسبات نرم در بهبود عملکرد مدل کوکومو به منظور تخمین هزینه پروژه های نرم افزاری»، همایش ملی مهندسی کامپیوتر و فناوری اطلاعات، شوشتر، دانشگاه آزاد اسلامی واحد شوشتر، ۱۳۹۲.