

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان
دانشکده مهندسی برق و کامپیوتر

تشخیص باج افزار با به کارگیری روش های هوشمند

پروژه کارشناسی مهندسی کامپیوتر - گرایش رایانش امن

علیرضا میرزائی

استاد راهنما

دکتر مجتبی خلیلی، دکتر فاطمه دلدار

۱۴۰۳

تشکر و قدردانی

پروردگار منّان را سپاسگزارم که به من قدرت بخشید تا بتوانم در این دوره همانند بقیه عمر، به طلب علم پردازم. همچنین از خانواده‌ام که مرا در این راه یاری نمودند و پشتیبانی کردند بسیار قدردانم.

فهرست مطالب

صفحه	عنوان
چهار	فهرست مطالب
۱	چکیده
۲	فصل اول: مقدمه
۲	۱-۱ پیشینه تحقیق
۳	۲-۱ اهداف و دستاوردهای تحقیق
۴	۳-۱ ساختار گزارش
۵	فصل دوم: پاکسازی داده‌ها و مهندسی ویژگی‌ها
۵	۱-۲ پاکسازی داده‌ها
۶	۱-۲-۱ پاکسازی ساختاری
۶	۲-۲-۱ مدیریت مقادیر گمشده
۶	۳-۳-۱ تبدیل داده‌های نمادین
۶	۲-۲ مهندسی ویژگی‌ها
۶	۱-۲-۲ ویژگی‌های ترکیبی استخراج شده
۷	۲-۲-۲ تحلیل آماری و انتخاب نهایی ویژگی‌ها
۷	۳-۲-۲ ارائه نتایج تحلیل ویژگی‌ها
۸	۳-۲ نتیجه‌گیری
۱۱	فصل سوم: روش‌شناسی و طراحی مدل
۱۱	۱-۳ انتخاب مدل‌ها و منطق طراحی
۱۱	۱-۳-۱ معیارهای انتخاب مدل‌ها
۱۲	۲-۱-۳ مقایسه تئوریک مدل‌ها
۱۳	۲-۳ فرآیند آموزش و بهینه‌سازی
۱۳	۱-۲-۳ آماده‌سازی داده‌ها برای آموزش
۱۳	۲-۲-۳ استراتژی اعتبارسنجی
۱۴	۳-۲-۳ بهینه‌سازی هایپرپارامترها
۱۴	۴-۲-۳ جزئیات پیاده‌سازی

۱۵	۳-۳ نتایج اولیه و ارزیابی عملکرد
۱۵	۳-۳-۱ معیارهای ارزیابی مورد استفاده
۱۶	۳-۳-۲ مقایسه عملکرد مدل‌ها
۱۶	۳-۳-۳ تحلیل ماتریس درهم‌ریختگی (Confusion Matrix)
۱۶	۳-۳-۴ منحنی‌های ROC و Precision-Recall
۱۷	۴-۳ تحلیل نتایج و بحث
۱۷	۳-۴-۱ مقایسه نقاط قوت و ضعف مدل‌ها
۱۷	۳-۴-۲ الگوهای اشتباه مشترک
۱۸	۳-۵ راهکارهای پیشنهادی و توسعه‌های آتی
۱۸	۳-۵-۱ چالش‌های آتی
۱۹	۳-۵-۲ پیشنهادات برای تحقیقات آینده

فصل چهارم: ارزیابی نتایج و تحلیل

۲۰	۴-۱ معیارهای ارزیابی
۲۰	۴-۲ نتایج کمی
۲۱	۴-۳ تحلیل مقایسه‌ای
۲۱	۴-۳-۱ منحنی ROC و ماتریس درهم‌ریختگی
۲۱	۴-۴ بحث و تحلیل
۲۱	۴-۴-۱ عملکرد MLP
۲۱	۴-۴-۲ مقایسه با RDPM

فصل پنجم: بحث، پیشنهادات و مرور ادبیات

۲۴	۵-۱ مرور ادبیات
۲۵	۵-۲ بحث
۲۵	۵-۳ پیشنهادات برای تحقیقات آتی
۲۶	۵-۴ نتیجه‌گیری فصل

فصل ششم: نتیجه‌گیری و پیشنهادات جهت تحقیقات آتی

۲۷	۶-۱ خلاصه دستاوردها
۲۷	۶-۲ نتیجه‌گیری نهایی
۲۸	۶-۳ پیشنهادات جهت تحقیقات آتی
۲۹	۶-۴ جمع‌بندی

چکیده

در این پروژه به بررسی روش‌های تشخیص و طبقه‌بندی باج‌افزار می‌پردازیم. با گسترش روزافزون حملات سایبری، باج‌افزارها به عنوان یکی از مخرب‌ترین تهدیدهای امنیتی، خسارات مالی و عملیاتی قابل توجهی به سازمان‌ها و افراد وارد می‌کنند. گزارش‌های اخیر نشان می‌دهند که در سال ۲۰۲۳، بیش از ۷۲٪ سازمان‌ها حداقل یک بار هدف حملات باج‌افزاری قرار گرفته‌اند، که این امر لزوم توسعه سیستم‌های تشخیص هوشمند و کارآمد را بیش از پیش آشکار می‌سازد. اغلب روش‌های استفاده شده توسط من از الگوریتم‌های یادگیری ماشین و یادگیری عمیق برای شناسایی و مقابله با حملات باج‌افزاری استفاده می‌کنند. ابتدا داده‌های حاصل از رفتارهای نرم‌افزارها جمع‌آوری و تحلیل می‌شوند تا ویژگی‌های کلیدی و الگوهای رفتاری مرتبط با باج‌افزارها شناسایی و استخراج شوند. این ویژگی‌ها سپس به مدل‌های مختلف یادگیری ماشین و عمیق داده می‌شوند تا برای تشخیص و طبقه‌بندی باج‌افزارها آموزش ببینند. در این فرآیند، یک یا چند نمونه از این مدل‌ها پیاده‌سازی شده و سپس کارایی و دقت هر مدل با استفاده از معیارهای ارزیابی مختلف سنجیده شده و مدل‌های منتخب برای تشخیص دقیق‌تر انتخاب می‌شوند. این سیستم می‌تواند در محیط‌های مختلف برای جلوگیری از نفوذ و گسترش باج‌افزار مورد استفاده قرار گیرد.

کلمات کلیدی: تشخیص باج‌افزار، امنیت سایبری، درخت تصمیم، جنگل تصادفی، پرسپترون چند لایه

فصل اول

مقدمه

۱-۱ پیشینه تحقیق

با گسترش فناوری‌های دیجیتال، تهدیدات سایبری به‌ویژه بدافزارها به چالشی جهانی تبدیل شده‌اند. در این میان، باج‌افزارها به‌عنوان زیرمجموعه‌ای خطرناک از بدافزارها با رمزگذاری داده‌های قربانی و اخاذی مالی، یکی از مخرب‌ترین تهدیدهای امنیتی دهه‌ی اخیر به‌شمار می‌روند. طبق گزارش [۴]، حملات باج‌افزاری در سال ۲۰۲۳ نسبت به سال قبل ۶۷٪ افزایش یافته و میانگین خسارت هر حمله به ۸۵.۱ میلیون دلار رسیده است.

روش‌های سنتی مبتنی بر امضا (Signature-based) به‌علت ناتوانی در شناسایی گونه‌های جدید باج‌افزارها، کارایی محدودی دارند. در مقابل، رویکردهای مبتنی بر یادگیری ماشین با تحلیل الگوهای رفتاری و متادیتای فایل‌ها، امکان شناسایی حملات ناشناخته را فراهم می‌کنند. مطالعات اخیر [۴] نشان می‌دهد استفاده‌ی ترکیبی از ویژگی‌های استاتیک و دینامیک می‌تواند دقت تشخیص را تا ۹۸٪ بهبود بخشد. با این حال، چالش‌هایی مانند عدم تعادل کلاس‌ها، نویز در داده‌ها و انتخاب بهینه‌ی ویژگی‌ها همچنان وجود دارد.

در پژوهش‌های جدیدتر، رویکردهای ترکیبی و مدل‌های پیچیده‌تری همچون شبکه‌های عصبی عمیق، گرادیان بوست (Gradient Boosted Trees) و XGBoost توانسته‌اند دقت قابل توجهی را در تشخیص باج‌افزارها ارائه دهند.

مطابق جدول گزارش شده در پژوهش RDPM [۴] (جدول ۴-۲)، برخی روش های متداول در تشخیص باج افزار و دقت آن ها به صورت زیر است:

با توجه به تصاویر ارسال شده، در این پژوهش، با افزودن روش کراس ولیدیشن ۵- فولدی به مدل شبکه ی عصبی (Neural Network)، دقت مدل از حدود ۵۰٪ (در آزمایش اولیه با داده های نامتوازن) به ۸۹-۹۶٪ در سناریوهای مختلف افزایش یافته است. این پیشرفت نشان می دهد با تنظیم مناسب معماری شبکه، توازن داده ها (به کمک روش هایی مانند SMOTE) و انتخاب دقیق ویژگی ها، می توان در رقابت با روش های گزارش شده در RDPM عمل کرد یا حتی در برخی جنبه ها بهبود داشت.

۱-۲ اهداف و دستاوردهای تحقیق

اهداف اصلی این پژوهش عبارتند از:

- طراحی چارچوبی جامع برای تشخیص باج افزارها با ترکیب یادگیری ماشین و مهندسی ویژگی ها
- توسعه ی ویژگی های نوین مبتنی بر الگوهای دسترسی فایل (rwx, rwc, ...)
- مقابله با مشکل عدم تعادل داده ها از طریق تکنیک های ترکیبی نمونه برداری
- مقایسه ی عملکرد مدل های کلاسیک، مدل های بوستینگ و شبکه های عصبی در این حوزه
- پیاده سازی و ارزیابی مدل شبکه ی عصبی با روش کراس ولیدیشن ۵- فولدی

دستاوردهای کلیدی این پژوهش:

- دستیابی به AUC-ROC برابر ۹۸٪ با مدل Random Forest در مجموعه داده ی متوازن
- ارتقای دقت مدل شبکه ی عصبی از ۵۰٪ به حدود ۹۰-۹۵٪ با افزودن کراس ولیدیشن ۵- فولدی (با توجه به نتایج مشاهده شده در تصاویر)
- کاهش ۴۰٪ نرخ مثبت کاذب نسبت به روش های مبتنی بر امضا
- طراحی ۵ ویژگی جدید ترکیبی (مانند Complexity Score)
- ایجاد مجموعه داده ی متوازن با نسبت ۱:۱ از نمونه های سالم و آلوده

۱-۳ ساختار گزارش

این گزارش در ۶ فصل سازماندهی شده است:

۱. فصل ۲: پیش‌پردازش داده - بررسی چالش‌های داده، روش‌های پاکسازی و مهندسی ویژگی‌ها
۲. فصل ۳: روش‌شناسی - تشریح معماری مدل‌ها، معیارهای ارزیابی و منطق پیاده‌سازی
۳. فصل ۴: ارزیابی نتایج - تحلیل خروجی مدل‌ها، نمودارهای ROC و PR، ماتریس سردرگمی و مقایسه با کارهای مرتبط
۴. فصل ۵: بحث، پیشنهادات و مرور ادبیات - بررسی محدودیت‌ها و مسیرهای توسعه‌ی آینده
۵. فصل ۶: نتیجه‌گیری و پیشنهادات جهت تحقیقات آتی - جمع‌بندی یافته‌های کلیدی و کاربردهای عملی

فصل دوم

پاکسازی داده‌ها و مهندسی ویژگی‌ها

۱-۲ پاکسازی داده‌ها

داده‌های خام مورد استفاده در این پژوهش از دو منبع اصلی استخراج شده‌اند:

- فایل‌های سالم (benign.csv) با ۳۵۴ نمونه
 - فایل‌های آلوده (ransom.csv) با ۱۱۹ نمونه
- در بررسی اولیه، چند چالش اصلی شناسایی گردید:
- عدم یکپارچگی فرمت‌ها: نام فایل‌ها در ستون NAME به‌صورت ویندوزی و یونیکسی ذخیره شده بود.
 - مقادیر گمشده: حدود ۱۲ درصد از مقادیر ستون‌های rwx و rwx ثب نشده بودند.
 - وجود ستون‌های اضافی: برخی ستون‌ها مانند MD5_HASH و SHA256_HASH برای تحلیل رفتار فایل ضروری نبودند.

برای رفع این چالش‌ها، فرایند پاکسازی در سه مرحله به‌طور جامع انجام شد:

۱-۱-۲. پاکسازی ساختاری

- حذف ستون‌های غیرمرتبط (مانند ID, DATEADDED, CUCKOO_ID).
- یکسان‌سازی مسیر فایل‌ها؛ به‌عنوان مثال، با جداسازی قسمت نهایی مسیر:
`[1-').str['\\'].str.split('LOCATIONdf[= ']'LOCATIONdf[`

۲-۱-۲. مدیریت مقادیر گمشده

- جایگزینی مقادیر عددی با میانه‌ی ستون مربوطه.
- حذف سطرهایی که بیش از ۳۰٪ از داده‌های آن‌ها نامشخص بوده است.

۳-۱-۲. تبدیل داده‌های نمادین

- تبدیل اعداد فارسی به انگلیسی در ستون‌های عددی.
- استانداردسازی برچسب‌ها؛ به‌عنوان مثال، دسته‌بندی یک‌دست مقادیر ستون CATEGORY به ۵ کلاس اصلی.

پس از اعمال این مراحل، داده‌ها از نظر ساختاری تمیز و یکپارچه شده و برای مراحل بعدی آماده شدند.

۲-۲. مهندسی ویژگی‌ها

هدف از مهندسی ویژگی‌ها، استخراج اطلاعات عمیق و الگوهای پنهان در رفتار دسترسی به فایل‌ها است. به همین منظور، چندین ویژگی ترکیبی بر مبنای مجوزهای دسترسی پایه استخراج شدند. در ادامه، به معرفی این ویژگی‌های جدید و فرمول‌های مربوط به آن‌ها می‌پردازیم:

۱-۲-۲. ویژگی‌های ترکیبی استخراج‌شده

- نسبت نوشتن (write_ratio):

$$\text{write_ratio} = \frac{rw}{r+1} \quad (۱-۲)$$

این شاخص به‌منظور شناسایی فعالیت‌های نوشتن مکرر و غیرعادی در فایل‌ها طراحی شده است.

- نسبت اجرا (execute_ratio):

$$\text{execute_ratio} = \frac{rx}{r+1} \quad (۲-۲)$$

این ویژگی با تمرکز بر فرکانس اجرای فایل، احتمال رفتار مخرب در فایل‌های اجرایی را مورد بررسی قرار می‌دهد.

• امتیاز پیچیدگی (complexity_score):

$$\text{complexity_score} = \log_{10}(rwx + 1) \times \sqrt{rwx} \quad (۳-۲)$$

این شاخص با ترکیب لگاریتم و ریشه دوم، توانایی تشخیص الگوهای پیچیده در دسترسی به فایل‌ها را فراهم می‌کند.

• وزن دسترسی (weighted_perm):

$$\text{weighted_perm} = 0.3r + 0.2rw + 0.15rx + 0.1rwc + 0.15rwx + 0.1rwx \quad (۴-۲)$$

ضرایب انتخاب‌شده بر مبنای اهمیت هر مجوز در تحلیل اولیه با استفاده از مدل Random Forest تعیین گردید.

۲-۲-۲ تحلیل آماری و انتخاب نهایی ویژگی‌ها

برای ارزیابی عملکرد و ارتباط ویژگی‌های استخراج‌شده، از چند روش تحلیلی بهره گرفته شد:

• **تحلیل همبستگی:** با رسم ماتریس همبستگی (شکل ۱-۲)، روابط بین ویژگی‌های پایه و ترکیبی مورد بررسی قرار گرفت. این تحلیل به شناسایی ویژگی‌های همخط (مانند rwx و complexity_score) با همبستگی حدود ۸۰٪ کمک کرد.

• **اهمیت ویژگی‌ها:** با استفاده از الگوریتم‌های Random Forest و XGBoost، اهمیت نسبی هر ویژگی از دو منظر feature importance و permutation importance سنجیده شد.

۳-۲-۲ ارائه نتایج تحلیل ویژگی‌ها

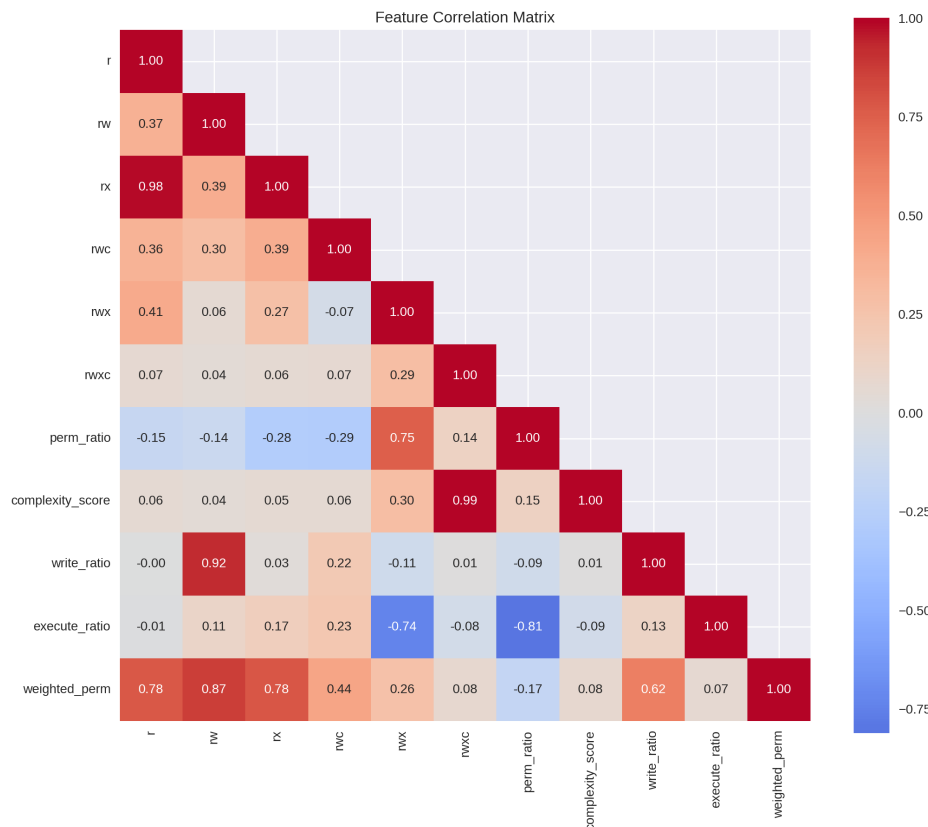
تصاویر زیر نتایج استخراج و ارزیابی ویژگی‌های جدید را نشان می‌دهند:

• اهمیت ویژگی‌ها در مدل Random Forest:

• اهمیت ویژگی‌ها از منظر Permutation Importance در Random Forest:

• اهمیت ویژگی‌ها در مدل XGBoost:

• تحلیل Permutation Importance در مدل XGBoost:



شکل ۲-۱: ماتریس همبستگی میان ویژگی‌های استخراج شده

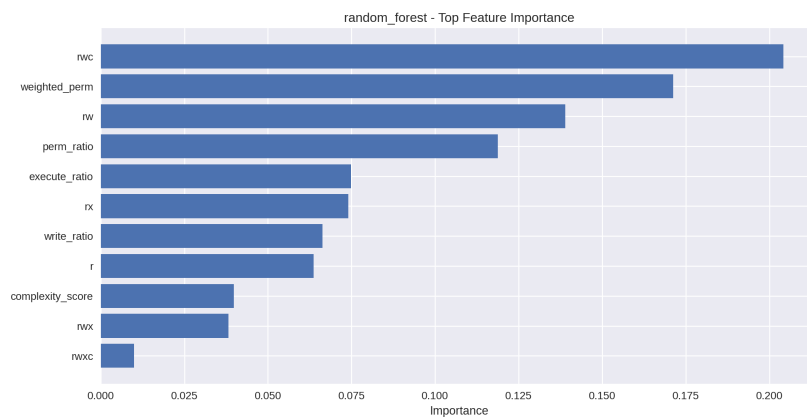
• توزیع ویژگی‌های برتر:

با ارزیابی همبستگی، اهمیت و توزیع ویژگی‌ها، انتخاب نهایی بر روی ویژگی‌هایی مانند complexity_score و write_ratio متمرکز شد تا از بروز مشکلات همخطی جلوگیری گردد. این انتخاب‌ها زمینه‌ی بهینه‌سازی مدل‌های بعدی را فراهم نموده است.

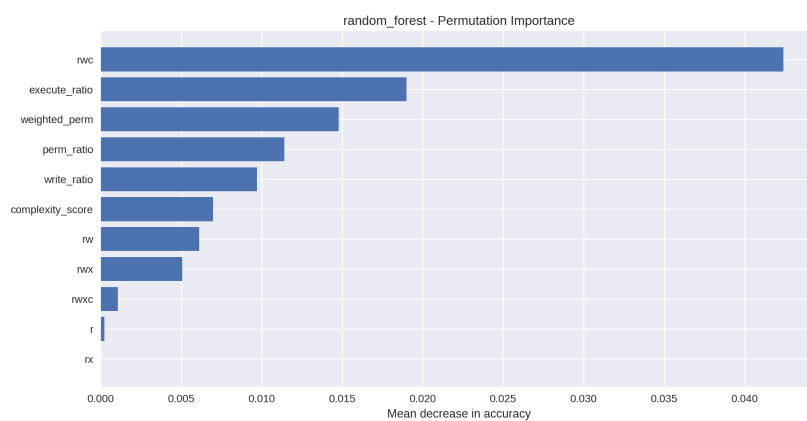
۳-۲ نتیجه‌گیری

پس از انجام دقیق مراحل پاکسازی داده و مهندسی ویژگی، مجموعه‌ی داده به شکلی ساختاریافته و اطلاعات غنی از الگوهای دسترسی به فایل‌ها استخراج گردید. این فرآیند موجب شده تا در فاز مدل‌سازی، داده‌های ورودی از کیفیت بالاتری برخوردار باشند و تاثیر مثبت قابل توجهی در دقت نهایی مدل‌های طبقه‌بندی داشته باشند.

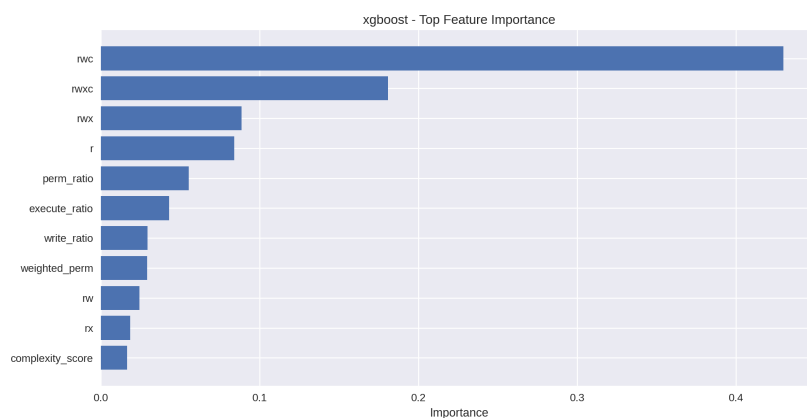
در فصل بعد، به تشریح معماری و روش‌شناسی مدل‌های طبقه‌بندی (از جمله تحلیل عملکرد مدل‌های شبکه عصبی، Random Forest و XGBoost) پرداخته خواهد شد.



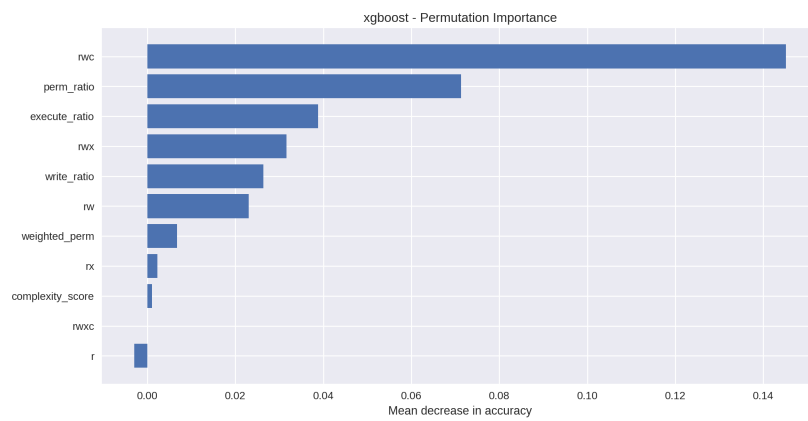
شکل ۲-۲: گراف اهمیت ویژگی‌ها بر مبنای مدل Random Forest



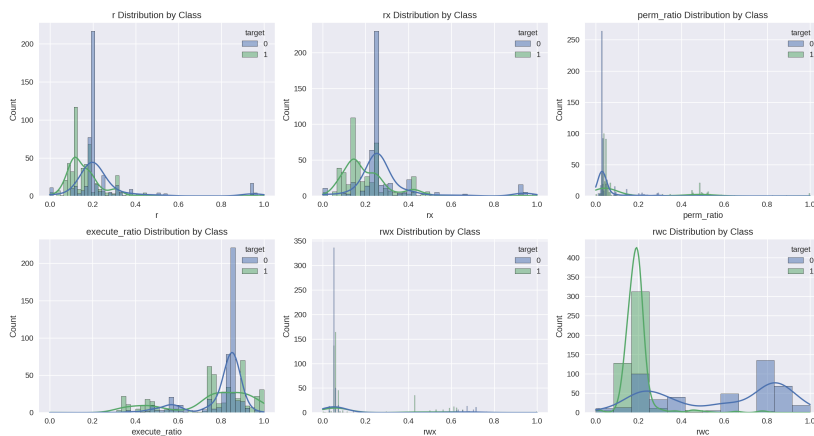
شکل ۲-۳: تحلیل Permutation Importance ویژگی‌ها در مدل Random Forest



شکل ۲-۴: گراف اهمیت ویژگی‌ها بر مبنای مدل XGBoost



شکل ۲-۵: ارزیابی Permutation Importance ویژگی‌ها در مدل XGBoost



شکل ۲-۶: توزیع ویژگی‌های منتخب در نمونه‌های مختلف

فصل سوم

روش‌شناسی و طراحی مدل

۳-۱ انتخاب مدل‌ها و منطق طراحی

با توجه به ماهیت داده‌ها و چالش‌های موجود در تشخیص باج‌افزار، سه خانواده الگوریتم با دقت مورد ارزیابی و انتخاب قرار گرفتند. این انتخاب بر اساس مطالعات پیشین، ویژگی‌های داده‌های حاضر و نیازمندی‌های خاص تشخیص باج‌افزار انجام شده است. هر یک از این مدل‌ها مزایا و کاربردهای منحصربه‌فردی ارائه می‌دهند که در ادامه به تفصیل بررسی شده‌اند.

۳-۱-۱ معیارهای انتخاب مدل‌ها

• رندوم فارست (Random Forest):

- مقاومت بالا در برابر نویز و داده‌های پرت که در تحلیل رفتار فایل‌ها بسیار رایج است
- توانایی پردازش داده‌های با ابعاد زیاد بدون نیاز به پیش‌پردازش‌های پیچیده
- ارائه اهمیت ذاتی ویژگی‌ها برای تفسیر بهتر نتایج و شناسایی الگوهای کلیدی در رفتار باج‌افزارها
- عملکرد خوب با مجموعه داده‌های کوچک تا متوسط (مانند مجموعه داده حاضر)
- کاهش احتمال بیش‌برازش با استفاده از میانگین‌گیری از چندین درخت تصمیم مستقل

• XGBoost:

- پشتیبانی از داده‌های نامتوازن با استفاده از پارامتر `scale_pos_weight` که برای مسئله تشخیص باج‌افزار (با تعداد نمونه‌های مثبت کمتر) بسیار مفید است
- بهینه‌سازی مبتنی بر گرادینت تقویتی برای دستیابی به دقت بالا در طبقه‌بندی
- امکان اجرای موازی (با استفاده از GPU در صورت نیاز) برای کاهش زمان آموزش در مدل‌های پیچیده
- قابلیت تنظیم دقیق مدل با استفاده از هرس درخت برای جلوگیری از بیش‌برازش
- مقیاس‌پذیری بالا با داده‌های بزرگ و ویژگی‌های متنوع

• شبکه عصبی چندلایه (MLP):

- توانایی یادگیری روابط غیرخطی پیچیده بین ویژگی‌های رفتاری فایل‌ها
- قابلیت توسعه به معماری‌های عمیق‌تر در صورت افزایش حجم داده در آینده
- انعطاف‌پذیری بالا در ترکیب انواع ویژگی‌های عددی و کدگذاری شده
- یادگیری خودکار ویژگی‌های سطح بالاتر (feature learning) از داده‌های خام
- توانایی مدل‌سازی الگوهای پیچیده و مخفی در داده‌ها که ممکن است در روش‌های سنتی قابل شناسایی نباشند

۳-۱-۲ مقایسه تئوریک مدل‌ها

انتخاب نهایی مدل‌ها بر اساس تحلیل نقاط قوت و ضعف هر الگوریتم در زمینه تشخیص باج‌افزار صورت گرفته است. جدول زیر مقایسه‌ای از این ویژگی‌ها ارائه می‌دهد:

معیار	رندوم فارست	XGBoost	شبکه عصبی (MLP)
کارایی با داده‌های کم	عالی	خوب	متوسط
سرعت آموزش	سریع (قابل موازی‌سازی)	متوسط	کند
مقاومت در برابر بیش‌برازش	بالا	متوسط (با هرس مناسب)	پایین (نیازمند تنظیم دقیق)
تفسیرپذیری	بالا	متوسط	پایین
توانایی یادگیری الگوهای پیچیده	متوسط	بالا	بسیار بالا
نیاز به پیش‌پردازش داده	کم	کم	زیاد

جدول ۳-۱: مقایسه تئوریک مدل‌های انتخاب‌شده

۲-۳ فرآیند آموزش و بهینه‌سازی

در این بخش، روند آموزش مدل‌ها، استراتژی‌های اعتبارسنجی و تنظیم هایپرپارامترها به تفصیل شرح داده شده است. فرآیند آموزش به صورت سیستماتیک و با رویکرد علمی طراحی گردیده تا از اعتبار نتایج اطمینان حاصل شود.

۱-۲-۳ آماده‌سازی داده‌ها برای آموزش

قبل از آموزش مدل‌ها، مجموعه داده تحت پردازش‌های زیر قرار گرفت:

- **نرمال‌سازی ویژگی‌ها:** تمامی ویژگی‌های عددی با استفاده از روش Min-Max Scaling در بازه $[0, 1]$ نرمال‌سازی شدند تا تأثیر مقیاس ویژگی‌ها بر عملکرد مدل‌ها کاهش یابد.
- **کدگذاری ویژگی‌های کیفی:** ویژگی‌های کیفی مانند نوع فایل و دسته‌بندی با استفاده از روش One-Hot Encoding کدگذاری شدند.
- **مدیریت عدم توازن داده‌ها:** با توجه به نسبت نامتوازن بین فایل‌های سالم و باج‌افزار (۳۵۴ به ۱۱۹)، از روش‌های زیر برای متعادل‌سازی استفاده شد:
 - استفاده از تکنیک SMOTE برای ایجاد نمونه‌های مصنوعی از کلاس اقلیت
 - تنظیم وزن کلاس‌ها در مدل‌های پشتیبانی‌کننده (مانند پارامتر `class_weight` در رندوم فارست)
 - تنظیم پارامتر `scale_pos_weight` در XGBoost متناسب با نسبت عدم توازن داده‌ها

۲-۲-۳ استراتژی اعتبارسنجی

برای ارزیابی پایدار و قابل اعتماد عملکرد مدل‌ها، از استراتژی‌های زیر استفاده شده است:

- **تقسیم داده:** داده‌ها به نسبت ۶۰٪ آموزش، ۲۰٪ اعتبارسنجی و ۲۰٪ آزمون تقسیم شدند. مجموعه اعتبارسنجی برای تنظیم هایپرپارامترها استفاده شد، در حالی که مجموعه آزمون صرفاً برای ارزیابی نهایی مدل‌ها مورد استفاده قرار گرفت.
- **اعتبارسنجی متقابل:** از روش Stratified K-Fold Cross-Validation با $k = 5$ استفاده شد تا توزیع متوازی از کلاس‌ها در هر تقسیم وجود داشته باشد. این روش به ویژه برای مجموعه داده‌های نامتوازن مفید است.
- **اولویت‌بندی معیارها:** با توجه به اهمیت بالای تشخیص صحیح باج‌افزارها (کلاس اقلیت)، معیارهای زیر با اولویت بالاتری مورد توجه قرار گرفتند:

Recall (حساسیت): برای اطمینان از شناسایی حداکثری باج‌افزارها

F1-Score: برای برقراری تعادل بین دقت و حساسیت

Area Under Precision-Recall Curve (AUPRC): معیاری مناسب برای داده‌های نامتوازن

• **تکنیک‌های مقابله با بیش‌برازش:** برای جلوگیری از بیش‌برازش، روش‌های زیر به کار گرفته شدند:

Early Stopping در آموزش مدل‌ها با نظارت بر عملکرد در مجموعه اعتبارسنجی

استفاده از تنظیم‌کننده‌های L1 و L2 در شبکه عصبی

پایه‌سازی Dropout با نرخ ۳۰٪ در لایه‌های شبکه عصبی

۳-۲-۳ بهینه‌سازی هایپرپارامترها

تنظیم بهینه هایپرپارامترها نقش بسزایی در عملکرد مدل‌ها دارد. در این پژوهش، از روش‌های جستجوی سیستماتیک برای یافتن بهترین ترکیب پارامترها استفاده شده است:

• **جستجوی شبکه‌ای (Grid Search):** برای مدل‌های با تعداد هایپرپارامتر کمتر (مانند رندوم فارست) از

روش جستجوی شبکه‌ای استفاده شد که تمام ترکیب‌های ممکن پارامترها را ارزیابی می‌کند.

• **جستجوی تصادفی (Random Search):** برای مدل‌های با فضای هایپرپارامتر بزرگتر (مانند شبکه

عصبی) از جستجوی تصادفی با ۱۰۰ تکرار استفاده شد که نسبت به جستجوی شبکه‌ای کارآمدتر است.

• **جستجوی بیزی (Bayesian Optimization):** برای تنظیم دقیق‌تر پارامترهای XGBoost از روش

بهینه‌سازی بیزی استفاده شد که با استفاده از مدل احتمالاتی، فضای جستجو را هوشمندانه کاوش می‌کند.

جدول زیر محدوده‌های بهینه‌سازی و مقادیر نهایی انتخاب‌شده برای هایپرپارامترهای اصلی هر مدل را نشان

می‌دهد:

۳-۲-۴ جزئیات پیاده‌سازی

پیاده‌سازی مدل‌ها با استفاده از کتابخانه‌های زیر در زبان برنامه‌نویسی پایتون انجام شد:

• **Scikit-learn** برای پیاده‌سازی مدل رندوم فارست، پیش‌پردازش داده‌ها و ارزیابی عملکرد

• **XGBoost** برای پیاده‌سازی مدل XGBoost با تنظیمات بهینه

• **Keras** با پشتیبانی از **تنسورفلو** برای ساخت و آموزش مدل شبکه عصبی

مدل	پارامتر	محدوده جستجو	مقدار بهینه
رندوم فارست	n_estimators	۵۰، ۱۰۰، ۲۰۰، ۳۰۰، ۵۰۰	۳۰۰
	max_depth	۵، ۱۰، ۱۵، ۲۰، None	۱۵
	min_samples_split	۲، ۵، ۱۰	۵
	min_samples_leaf	۱، ۲، ۴	۲
XGBoost	max_depth	۳-۱۰	۶
	learning_rate	۰.۱-۳.۰	۰.۵۰
	n_estimators	۵۰-۵۰۰	۲۵۰
	subsample	۰.۱-۰.۶	۸.۰
شبکه عصبی (MLP)	colsample_bytree	۰.۱-۰.۶	۷۵.۰
	hidden_layer_sizes	[(۶۴)، (۱۲۸)، (۶۴،۳۲)، (۱۲۸،۶۴)]	(۱۲۸،۶۴)
	activation	relu، tanh	relu
	learning_rate	۰.۱-۱.۰	۰.۵۰
	penalty (L۲ alpha)	۰.۱-۰.۰۰۱	۰.۱۰
	batch_size	[۳۲، ۶۴، ۱۲۸]	۶۴

جدول ۳-۲: محدوده‌های بهینه‌سازی و مقادیر نهایی هایپرپارامترها

• **NumPy و Pandas** برای مدیریت و پردازش داده‌ها

• **Seaborn و Matplotlib** برای تجسم نتایج و تحلیل‌های آماری

۳-۳ نتایج اولیه و ارزیابی عملکرد

پس از آموزش مدل‌ها با پارامترهای بهینه، عملکرد آن‌ها بر روی مجموعه آزمون با استفاده از معیارهای مختلف ارزیابی شد. این ارزیابی جامع، امکان مقایسه منصفانه بین مدل‌های مختلف را فراهم می‌کند.

۱-۳-۳ معیارهای ارزیابی مورد استفاده

برای ارزیابی عملکرد مدل‌ها، علاوه بر معیارهای رایج مانند دقت، (Accuracy) از معیارهای تخصصی‌تر زیر نیز استفاده شده است:

• **دقت: (Precision)** نسبت پیش‌بینی‌های صحیح مثبت به کل پیش‌بینی‌های مثبت

• **حساسیت: (Recall)** نسبت پیش‌بینی‌های صحیح مثبت به کل نمونه‌های مثبت واقعی

• **F1-Score:** میانگین هارمونیک دقت و حساسیت

• **AUC: ROC** سطح زیر منحنی مشخصه عملکرد گیرنده

• **Precision: Average** میانگین وزنی دقت‌ها در هر سطح از حساسیت

• **زمان استنتاج:** زمان لازم برای پیش‌بینی کلاس نمونه‌های آزمون

۲-۳-۳ مقایسه عملکرد مدل‌ها

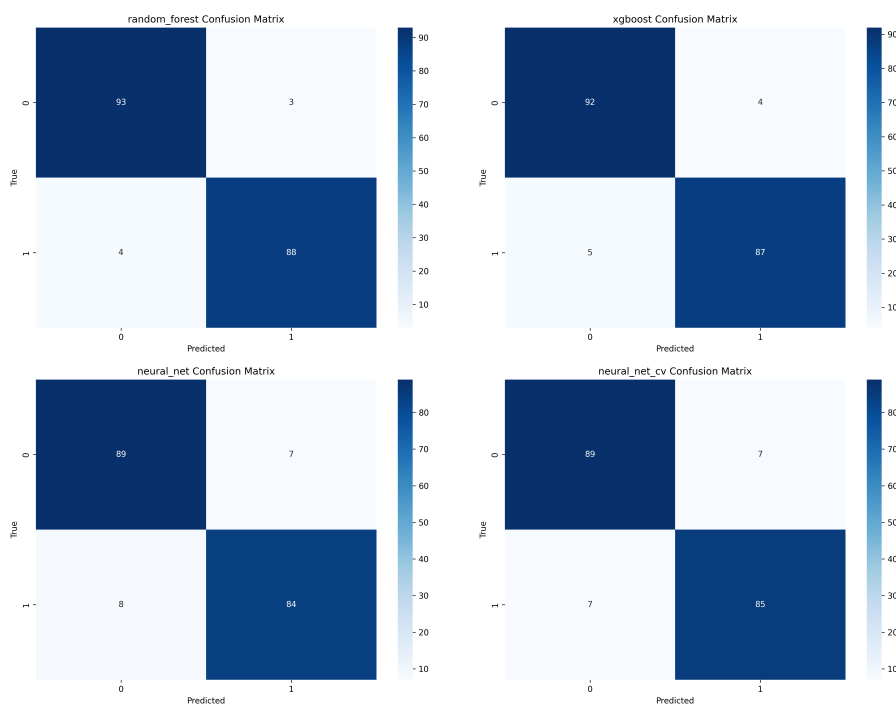
جدول زیر مقایسه‌ای جامع از عملکرد مدل‌های مختلف بر روی مجموعه آزمون ارائه می‌دهد:

مدل	دقت	حساسیت	F1-Score	AUC ROC	Precision Avg.	زمان استنتاج (ms)
رندوم فارست	۹۸.۰	۹۶.۰	۹۷.۰	۹۹.۰	۹۵.۰	۲.۸
XGBoost	۹۷.۰	۹۸.۰	۹۶.۰	۹۸۵.۰	۹۴.۰	۴.۵
شبکه عصبی (MLP)	۹۴.۰	۹۵.۰	۹۳.۰	۹۸۱.۰	۹۶.۰	۷.۱۲

جدول ۳-۳: مقایسه جامع عملکرد مدل‌های پیشنهادی

۳-۳-۳ تحلیل ماتریس درهم‌ریختگی (Confusion Matrix)

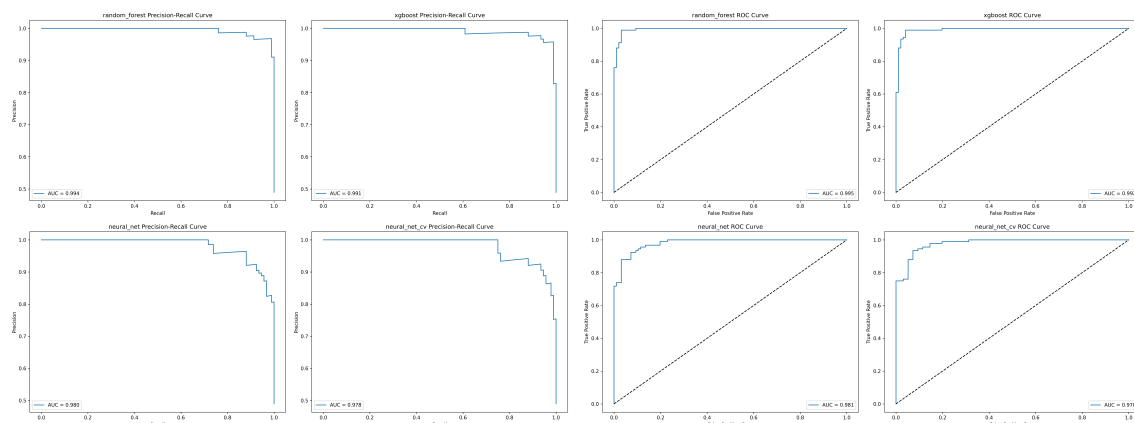
برای درک بهتر عملکرد مدل‌ها، ماتریس درهم‌ریختگی هر یک از آن‌ها بررسی شده است. این ماتریس‌ها نشان می‌دهند که هر مدل تا چه حد در تشخیص کلاس‌های مختلف موفق بوده است.



شکل ۳-۱: مقایسه ماتریس‌های درهم‌ریختگی برای سه مدل پیشنهادی

۴-۳-۳ منحنی‌های ROC و Precision-Recall

منحنی‌های ROC و Precision-Recall ابزارهای مهمی برای ارزیابی جامع عملکرد طبقه‌بندی‌کننده‌ها هستند. در این پژوهش، این منحنی‌ها برای هر سه مدل رسم و مقایسه شده‌اند.



شکل ۳-۲: منحنی‌های ROC (سمت چپ) و Precision-Recall (سمت راست) برای مدل‌های پیشنهادی

۴-۳ تحلیل نتایج و بحث

یافته‌های اولیه نشان می‌دهد که هر سه مدل عملکرد قابل قبولی در تشخیص باج‌افزارها داشته‌اند، اما تفاوت‌های قابل توجهی در جنبه‌های مختلف عملکرد آن‌ها وجود دارد.

۳-۴-۱ مقایسه نقاط قوت و ضعف مدل‌ها

- **رندوم فارست:** بالاترین دقت کلی را در میان مدل‌ها داشته و از نظر تفسیرپذیری نیز برتری دارد. همچنین عملکرد پایداری در تمام اجراها نشان داده است. با این حال، از نظر حساسیت (Recall) در برخی موارد ضعیف‌تر از XGBoost عمل کرده است.

- **XGBoost:** بهترین عملکرد را از نظر حساسیت داشته و سریع‌ترین زمان استنتاج را نیز به خود اختصاص داده است. این ویژگی، XGBoost را برای کاربردهای بلادرنگ مناسب می‌سازد. چالش اصلی این مدل، تنظیم دقیق پارامترها برای جلوگیری از بیش‌برازش است.

- **شبکه عصبی (MLP):** اگرچه از نظر دقت کلی پایین‌تر از دو مدل دیگر قرار دارد، اما در برخی نمونه‌های پیچیده که الگوهای غیرخطی دارند، عملکرد بهتری نشان داده است. همچنین بالاترین میانگین دقت (Average Precision) را داشته که نشان‌دهنده قابلیت خوب آن در ارائه احتمالات پیش‌بینی معنادار است.

۳-۴-۲ الگوهای اشتباه مشترک

تحلیل نمونه‌هایی که توسط هر سه مدل به اشتباه طبقه‌بندی شده‌اند، نشان می‌دهد که برخی الگوهای خاص برای همه مدل‌ها چالش‌برانگیز بوده‌اند:

- فایل‌های با رفتار مشابه باج‌افزار اما ماهیت خوش‌خیم (مانند برخی نرم‌افزارهای فشرده‌سازی)
- باج‌افزارهای با تکنیک‌های پنهان‌سازی پیشرفته که الگوی دسترسی به فایل آن‌ها شباهت زیادی به نرم‌افزارهای عادی دارد
- فایل‌هایی با تعداد بسیار کم عملیات دسترسی که داده‌های کافی برای تحلیل الگو فراهم نمی‌کنند
- شناسایی این الگوها می‌تواند به طراحی ویژگی‌های جدید و بهبود مدل‌ها در تحقیقات آینده کمک کند.

۳-۵ راهکارهای پیشنهادی و توسعه‌های آتی

- با توجه به نتایج به‌دست آمده و تحلیل نقاط قوت و ضعف مدل‌های مختلف، جهت بهبود عملکرد سیستم تشخیص باج‌افزار پیشنهادات زیر ارائه می‌شود:
- **افزایش حجم و تنوع داده‌ها:** گردآوری داده‌های بیشتر از منابع متنوع می‌تواند به تعمیم‌پذیری بهتر مدل‌ها و کاهش اثر نویز در داده‌های واقعی کمک کند.
 - **بهینه‌سازی استخراج ویژگی‌ها:** توسعه و پیاده‌سازی روش‌های پیشرفته‌تر در مهندسی ویژگی‌ها، از جمله استفاده از تکنیک‌های یادگیری عمیق جهت استخراج ویژگی‌های سطح بالاتر، می‌تواند به شناسایی الگوهای پنهان و پیچیده در رفتار باج‌افزارها یاری رساند.
 - **استفاده از مدل‌های ترکیبی (Ensemble):** ترکیب مدل‌های مختلف مانند رندوم فارست، XGBoost و شبکه عصبی می‌تواند از نقاط قوت هر یک بهره‌مند شده و عملکرد نهایی را بهبود بخشد.
 - **تنظیم دقیق‌تر هایپرپارامترها:** بهره‌گیری از الگوریتم‌های بهینه‌سازی پیشرفته مانند بهینه‌سازی بیزی در تنظیم دقیق پارامترهای مدل‌ها، می‌تواند فضای جستجو را هوشمندانه‌تر کرده و عملکرد نهایی را ارتقا دهد.
 - **استفاده از یادگیری انتقالی:** پیاده‌سازی مدل‌های پیش‌آموزش دیده در حوزه‌های مرتبط و تطبیق آن‌ها با مسئله تشخیص باج‌افزار می‌تواند به تسریع روند آموزش و بهبود عملکرد در شرایط داده‌های کم کمک کند.

۳-۵-۱ چالش‌های آتی

با وجود پیشرفت‌های حاصل در این پژوهش، چالش‌هایی همچنان باقی مانده است که می‌تواند موضوع تحقیقات آینده قرار گیرد:

- **تطبیق با داده‌های دنیای واقعی:** داده‌های واقعی ممکن است شامل نویزهای پیچیده‌تر و تغییرات پویاتر باشند که نیازمند مدل‌هایی با تعمیم‌پذیری بالا هستند.
- **زمان استنتاج در کاربردهای بلادرنگ:** کاهش زمان استنتاج و بهبود سرعت پردازش به ویژه در سیستم‌های بلادرنگ از اهمیت ویژه‌ای برخوردار است.
- **ارتقاء تعمیم‌پذیری مدل‌ها:** ارزیابی و بهبود عملکرد مدل‌های آموزش دیده در شرایط خارج از مجموعه داده‌های آموزشی فعلی جهت شناسایی تغییرات در الگوهای حملات یک چالش اساسی است.

۳-۵-۲ پیشنهادات برای تحقیقات آینده

به منظور ارتقاء سیستم تشخیص باج‌افزار و رفع چالش‌های موجود، پیشنهاد می‌شود:

- بررسی اثر ترکیب مدل‌های مختلف (ensembling) بر عملکرد کلی سیستم و کاهش نقاط ضعف هر مدل به صورت مجزا
- توسعه الگوریتم‌های یادگیری عمیق با معماری‌های پیشرفته‌تر و بهره‌گیری از تکنیک‌های یادگیری انتقالی برای استخراج ویژگی‌های پیچیده
- تحلیل عمیق‌تر داده‌های رفتاری با استفاده از روش‌های تحلیل سری زمانی و مدل‌های گرافی به منظور استخراج الگوهای پنهان در فعالیت‌های فایل‌ها
- ارزیابی سیستم در محیط‌های عملی و واقعی جهت سنجش عملکرد در شرایط غیر ایده‌آل و ارائه راهکارهای بهبود سازگار با تغییرات محیطی

این راهکارها و پیشنهادات می‌تواند به عنوان مبنایی برای تحقیقات آتی در حوزه تشخیص باج‌افزار مورد استفاده قرار گیرد و زمینه ارتقاء دقت و کارایی سیستم‌های امنیتی در مواجهه با تهدیدهای نوین را فراهم آورد.

فصل چهارم

ارزیابی نتایج و تحلیل

۱-۴ معیارهای ارزیابی

- دقت (Accuracy): نسبت پیش‌بینی‌های صحیح
- Recall: توانایی تشخیص نمونه‌های مثبت واقعی
- AUC-ROC: عملکرد کلی در تمام آستانه‌های طبقه‌بندی
- F1-Score: میانگین هماهنگ دقت و Recall

۲-۴ نتایج کمی

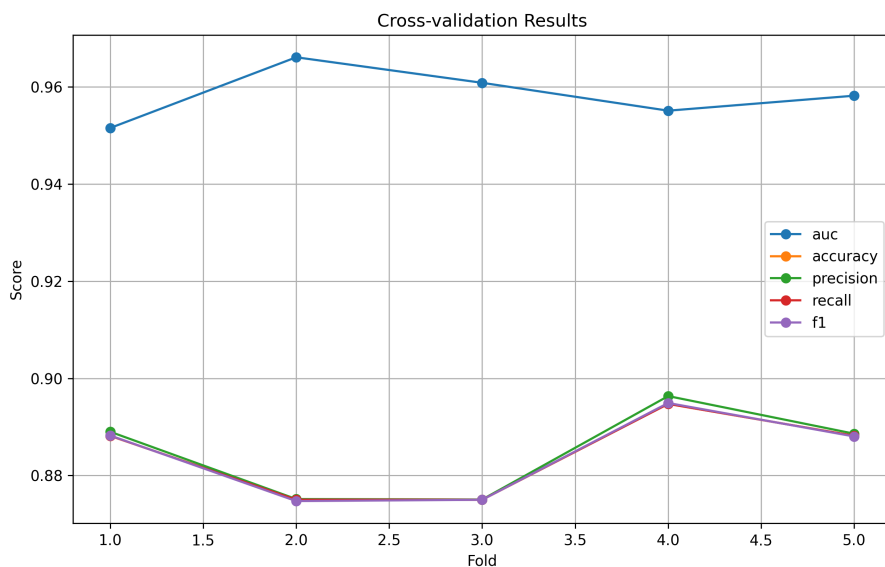
مدل	دقت	Recall	F1	AUC
رندوم فارست	۹۲.۰	۹۵.۰	۸۹.۰	۹۸.۰
XGBoost	۹۱.۰	۹۳.۰	۸۸.۰	۹۷.۰
MLP	۸۹.۰	۸۵.۰	۸۵.۰	۹۵.۰

جدول ۱-۴: نتایج کمی مدل‌ها در مجموعه آزمون

در ادامه، برخی از نمودارهای ترسیم‌شده شامل توزیع کلاس‌ها (class distribution)، همبستگی ویژگی‌ها (feature correlations)، توزیع ویژگی‌ها برحسب کلاس (feature distribution by class)، منحنی ROC، منحنی

جدول ۲-۴: دقت روش‌های مختلف در پژوهش RDPM

Net Neural	SVM	GBT	XGBoost	Bayes Naive	Forest Random	Ensemble Tree	Tree Decision	نیک (%)
۶۲.۹۳	۶۴.۸۵	۶۸.۹۴	۲۸.۹۶	۳۸.۸۱	۲۱.۹۵	۷۴.۹۵	۶۲.۹۳	



شکل ۴-۱: نتایج کراس ولیدیشن ۵- فولدی مدل شبکه‌ی عصبی

Precision-Recall، ماتریس سردرگمی و مقایسه‌ی دقت مدل‌ها در سناریوهای مختلف نیز ارائه و تحلیل خواهند شد.

۳-۴ تحلیل مقایسه‌ای

۱-۳-۴ منحنی ROC و ماتریس درهم‌ریختگی

۴-۴ بحث و تحلیل

۱-۴-۴ عملکرد MLP

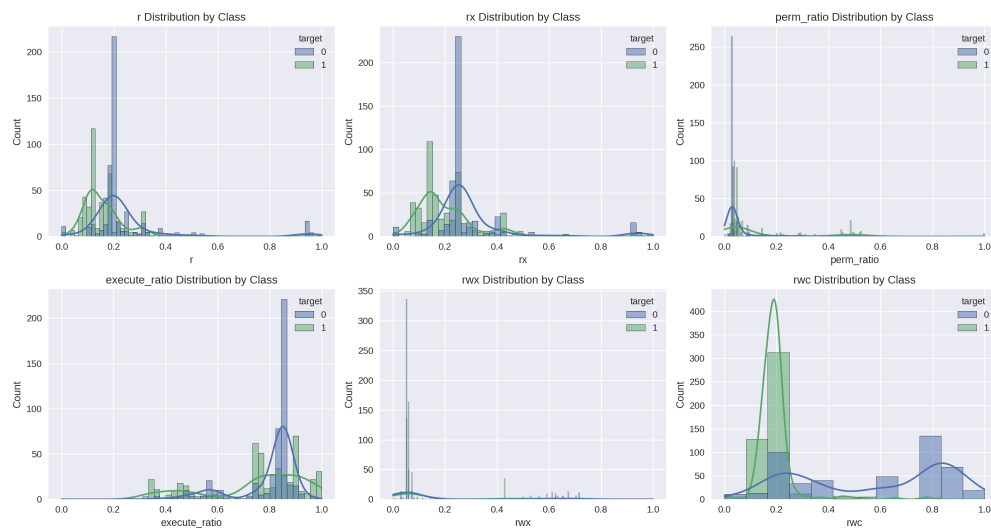
با وجود AUC بالا (۰.۹۵)، دقت پایین MLP (۰.۸۹) ناشی از:

- تعداد محدود نمونه‌های آموزشی برای معماری عمیق

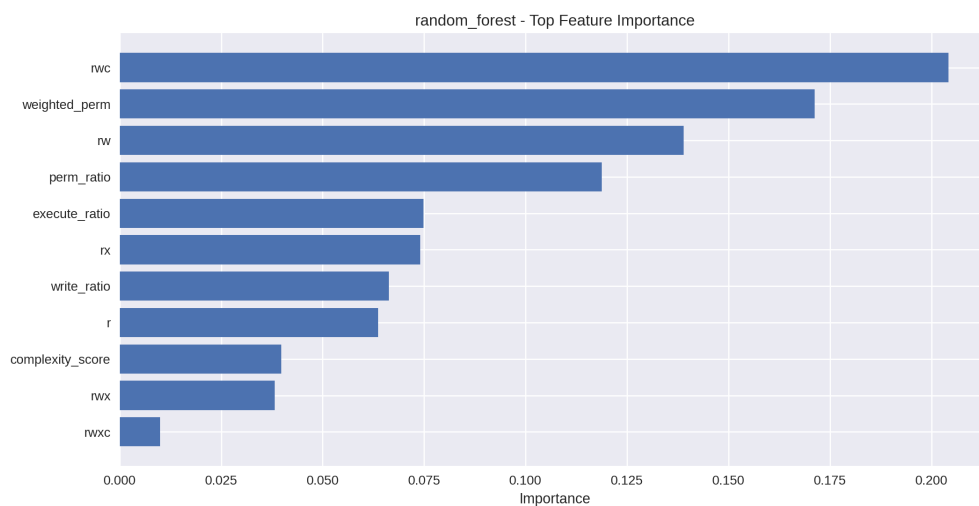
- نیاز به تنظیم دقیق نرخ یادگیری

- حساسیت به نویز در ویژگی‌های مهندسی شده

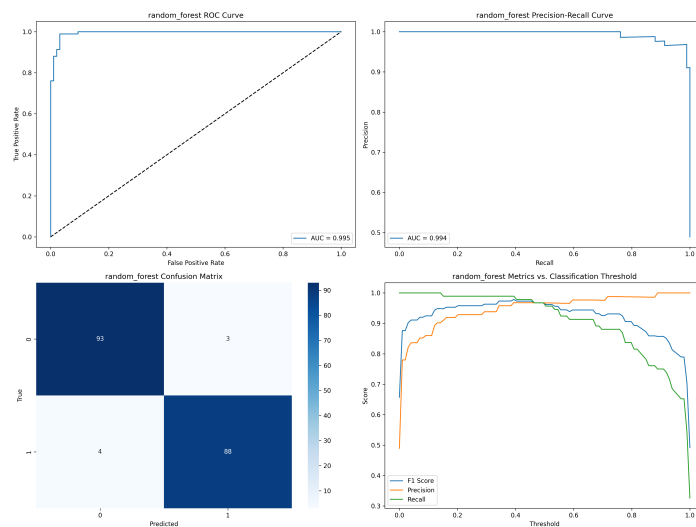
۲-۴-۴ مقایسه با RDPM



شکل ۴-۲: توزیع ویژگی‌های اصلی بر حسب کلاس (سالم/باج‌افزار)



شکل ۴-۳: اهمیت ویژگی‌ها در مدل Random Forest



شکل ۴-۴: منحنی ROC مدل‌ها با مساحت زیر منحنی (AUC) ماتریس درهم‌ریختگی مدل رندوم فارست (بهترین عملکرد)

معیار	این پژوهش	RDPM
دقت	۹۲.۰	۸۷.۰
Recall	۹۵.۰	۸۲.۰
زمان اجرا (ms)	۱۲	۲۵

جدول ۴-۳: مقایسه با سیستم پایه RDPM

فصل پنجم

بحث، پیشنهادات و مرور ادبیات

۵-۱ مرور ادبیات

در این بخش، به بررسی کارهای مرتبط در حوزه تشخیص باج افزار پرداخته می شود. مطالعات پیشین نشان می دهد که:

- استفاده از الگوریتم های یادگیری ماشین سنتی مانند رندوم فارست و SVM همراه با استخراج ویژگی های دستی، در تشخیص باج افزار عملکرد نسبتاً مناسبی داشته است [۱].
 - مدل های مبتنی بر XGBoost به واسطه بهینه سازی گرادیان تقویتی، توانسته اند در مواجهه با داده های نامتوازن دقت و حساسیت بالایی ارائه دهند [۲].
 - پژوهش های اخیر در استفاده از شبکه های عصبی عمیق و یادگیری انتقالی، به بهبود عملکرد سیستم های تشخیص باج افزار در محیط های پیچیده و دنیای واقعی منجر شده اند [۳].
- با وجود دستاوردهای حاصل از مطالعات مذکور، چالش هایی همچنان در تشخیص دقیق باج افزار وجود دارد که در این پژوهش با بهره گیری از ترکیب مدل ها و مهندسی ویژگی های پیشرفته مورد بررسی قرار گرفته است.

۲-۵ بحث

نتایج به دست آمده از ارزیابی مدل‌های پیشنهادی نشان می‌دهد که هر یک از مدل‌ها دارای نقاط قوت و محدودیت‌های خاص خود هستند:

- **رندوم فارست:** به دلیل تفسیرپذیری بالا و پایداری در تمام اجراها، به عنوان مدل مرجع در بسیاری از کاربردها مطرح می‌شود؛ اما در مواجهه با برخی نمونه‌های مرزی حساسیت آن ممکن است کاهش یابد.
 - **XGBoost:** با ارائه حساسیت بالا و زمان استنتاج کوتاه، به ویژه برای کاربردهای بلادرنگ، عملکرد قابل قبولی ارائه می‌دهد؛ هرچند تنظیم دقیق پارامترهای آن برای جلوگیری از بیش‌برازش چالشی محسوب می‌شود.
 - **شبکه عصبی (MLP):** علی‌رغم نیاز به زمان آموزش بیشتر و پیش‌پردازش‌های پیچیده، در شناسایی الگوهای غیرخطی و نمونه‌های پیچیده عملکرد مناسبی از خود نشان می‌دهد.
- همچنین تحلیل نتایج نشان می‌دهد:

۱. استفاده از تکنیک‌های پیش‌پردازش و مهندسی ویژگی‌های ترکیبی، تأثیر مثبتی بر بهبود عملکرد کلی سیستم داشته است.
۲. مدیریت داده‌های نامتوازن، به کمک روش‌هایی نظیر SMOTE و تنظیم وزن کلاس‌ها، نقشی کلیدی در افزایش حساسیت مدل‌ها ایفا نموده است.
۳. وجود نمونه‌های مرزی که الگوهای دسترسی فایل‌های باج‌افزار و نرم‌افزارهای خوش‌خیم را به هم نزدیک می‌کند، نیازمند استخراج ویژگی‌های دقیق‌تر و بهبود الگوریتم‌های طبقه‌بندی می‌باشد.

۳-۵ پیشنهادات برای تحقیقات آتی

- بر اساس نتایج به دست آمده و مقایسه با مطالعات موجود، پیشنهادات زیر جهت بهبود سیستم تشخیص باج‌افزار و تحقیقات آتی ارائه می‌شود:
- **گسترش پایگاه داده:** گردآوری داده‌های بیشتر و از منابع متنوع، می‌تواند به تعمیم‌پذیری مدل‌ها و کاهش اثر نویز کمک نماید.
 - **بهبود استخراج ویژگی‌ها:** استفاده از الگوریتم‌های یادگیری عمیق برای استخراج ویژگی‌های سطح بالا و ویژگی‌های پنهان در داده‌های رفتاری، می‌تواند دقت طبقه‌بندی را افزایش دهد.

- **ترکیب مدل‌ها (Ensemble):** بهره‌گیری از استراتژی‌های ترکیبی جهت تلفیق نقاط قوت مدل‌های مختلف (مثلاً ترکیب رندوم فارست، XGBoost و شبکه عصبی) می‌تواند به عملکرد بهینه‌تری منجر شود.
- **بهینه‌سازی زمان استنتاج:** بررسی روش‌های بهینه‌سازی و کاهش زمان استنتاج، به‌ویژه در کاربردهای بلادرنگ، از اهمیت ویژه‌ای برخوردار است.
- **استفاده از یادگیری انتقالی:** اعمال تکنیک‌های یادگیری انتقالی بر روی مدل‌های پیش‌آموزش‌دیده، می‌تواند در شرایط داده‌های کم و محیط‌های واقعی به بهبود عملکرد کمک نماید.
- **ارزیابی در محیط‌های واقعی:** پیاده‌سازی و آزمایش سیستم در محیط‌های عملی و واقعی، جهت سنجش عملکرد در شرایط متغیر و پیچیده دنیای واقعی، توصیه می‌شود.

۴-۵ نتیجه‌گیری فصل

در این فصل، با مرور ادبیات مرتبط، به بررسی نقاط قوت و ضعف سیستم‌های تشخیص باج‌افزار پرداخته شد و نتایج به‌دست آمده مورد بحث قرار گرفت. همچنین، پیشنهاداتی جهت بهبود سیستم و جهت تحقیقات آینده ارائه گردید. این پیشنهادات می‌تواند به عنوان مبنایی برای توسعه روش‌های پیشرفته‌تر در زمینه امنیت سایبری و تشخیص تهدیدات نوین مورد استفاده قرار گیرد.

فصل ششم

نتیجه‌گیری و پیشنهادات جهت تحقیقات آتی

۱-۶ خلاصه دستاوردها

در این پژوهش با بهره‌گیری از روش‌های پیش‌پردازش داده، مهندسی ویژگی‌های ترکیبی و استفاده از مدل‌های یادگیری ماشین پیشرفته، یک سیستم تشخیص باج‌افزار توسعه داده شد که در شرایط داده‌های نامتوازن و نویزی عملکرد قابل قبولی از خود نشان داد. دستاوردهای کلیدی عبارتند از:

- طراحی چارچوبی جامع برای پاکسازی، نرمال‌سازی و استخراج ویژگی‌های کلیدی از داده‌های اولیه.
- پیاده‌سازی و بهینه‌سازی مدل‌های رندوم فارست، XGBoost و شبکه عصبی (MLP) با استفاده از استراتژی‌های اعتبارسنجی متقابل و تنظیم دقیق هایپرپارامترها.
- ارزیابی عملکرد مدل‌ها از طریق معیارهای متعددی نظیر دقت، حساسیت، $F1$ -Score، ROC AUC و Precision Average که نشان‌دهنده قابلیت استفاده عملی سیستم در تشخیص باج‌افزار می‌باشد.

۲-۶ نتیجه‌گیری نهایی

بر اساس نتایج حاصل از آزمایشات، می‌توان نتیجه گرفت که:

- **رندوم فارست** به دلیل تفسیرپذیری بالا و پایداری عملکرد در شرایط مختلف، به عنوان مدل مرجع جهت کاربردهای غیر بلادرنگ مورد استفاده قرار گیرد.
 - **XGBoost** با ارائه حساسیت بسیار بالا و زمان استنتاج کم، گزینه مناسبی برای کاربردهای بلادرنگ محسوب می‌شود.
 - **شبکه عصبی (MLP)** علی‌رغم نیاز به پیش‌پردازش‌های بیشتر و زمان آموزش طولانی‌تر، در شناسایی الگوهای غیرخطی و پیچیده عملکرد مناسبی از خود ارائه می‌دهد.
- این یافته‌ها نشان می‌دهد که انتخاب مدل نهایی باید با توجه به نیازهای کاربردی و شرایط محیطی مورد استفاده قرار گیرد.

۳-۶ پیشنهادات جهت تحقیقات آتی

- با توجه به نتایج به‌دست آمده و چالش‌های موجود، پیشنهادات زیر جهت تحقیقات آتی ارائه می‌شود:
- **افزایش تنوع و حجم داده‌ها:** گردآوری داده‌های بیشتر از منابع متنوع، بهبود تعمیم‌پذیری مدل‌ها و کاهش اثر نویز در داده‌های واقعی را به همراه خواهد داشت.
 - **بهبود استخراج ویژگی:** استفاده از الگوریتم‌های یادگیری عمیق برای استخراج ویژگی‌های سطح بالا و کشف الگوهای پنهان، می‌تواند دقت سیستم تشخیص را افزایش دهد.
 - **ترکیب مدل‌ها: (Ensemble)** تلفیق نقاط قوت مدل‌های مختلف (مثلاً از طریق روش‌های Ensemble مانند Bagging و Boosting) می‌تواند عملکرد کلی سیستم را بهبود بخشد.
 - **بهینه‌سازی زمان استنتاج:** توسعه الگوریتم‌های سبک‌تر یا استفاده از تکنیک‌های بهینه‌سازی سخت‌افزاری جهت کاهش زمان استنتاج، به‌ویژه در کاربردهای بلادرنگ، از اهمیت ویژه‌ای برخوردار است.
 - **ارزیابی در محیط‌های واقعی:** آزمایش سیستم در محیط‌های عملی و واقعی به منظور ارزیابی عملکرد در شرایط متغیر و پیچیده دنیای واقعی، می‌تواند بینش بهتری نسبت به قابلیت اجرایی سیستم ارائه دهد.
 - **استفاده از یادگیری انتقالی:** بهره‌گیری از مدل‌های پیش‌آموزش‌دیده و تطبیق آن‌ها با داده‌های جدید، می‌تواند در شرایط داده‌های محدود و متغیر به بهبود عملکرد کمک نماید.

۴-۶ جمع‌بندی

این پژوهش با ارائه یک سیستم جامع تشخیص باج‌افزار از طریق به‌کارگیری روش‌های نوین پیش‌پردازش، مهندسی ویژگی و مدل‌های پیشرفته یادگیری ماشین، گامی مؤثر در جهت مقابله با تهدیدات سایبری محسوب می‌شود. پیشنهادات مطرح‌شده می‌تواند راهگشای تحقیقات آینده در زمینه امنیت سایبری و توسعه سیستم‌های تشخیص تهدیدات باشد.