



# «کارآموزی مهندسی یادگیری ماشین»

گزارش کارآموزی

ارائه شده به:

سرکار خانم دکتر منیره عبدوس

توسط:

محمد هاشمی

دانشکده مهندسی و علوم کامپیوتر

تابستان سال ۱۴۰۰

## پیش‌گفتار

تأثیر جادویی علم شگفت‌انگیز هوش مصنوعی<sup>۱</sup> زندگی ما را راحت‌تر و متمایزتر از قبل کرده است. اخیراً نیز علاقه زیادی در زمینه یادگیری ماشین<sup>۲</sup> بوجود آمده است و افراد زیادی متوجه گستره کاربردهای جدید ایجاد شده توسط رویکرد یادگیری ماشین می‌شوند. این تکنیک که جزو زیر شاخه های علم بزرگ هوش مصنوعی محسوب می‌شود یک نقشه‌ی راه برای ارتباط با دستگاه ایجاد می‌کند و آن را برای پاسخ به دستورالعمل‌ها و دستورات ما قابل فهم می‌سازد. در این گزارش مهمترین و پرکاربردترین روش‌ها و الگوریتم‌های یادگیری ماشین مورد مطالعه و بررسی قرار می‌گیرد. همچنین درباره کاربردهای شاخص و واقعی که امروزه در دنیای واقعی در صنعت استفاده می‌شوند بحث می‌شود. بدلیل کاربردهای بسیار متنوعی که این علم در صنعت پیدا کرده‌است امروزه شرکت‌های نوپا<sup>۳</sup> و همچنین شرکت‌های بزرگ و قدیمی بسیاری هستند که الگوریتم‌های یادگیری ماشین را به مرحله پیاده‌سازی و بهره‌وری رسانده‌اند. فرصت‌های شغلی‌ای نظیر دانشمند داده<sup>۴</sup> و مهندس یادگیری ماشین دلایلی هستند که بنده را مجاب کرده‌اند تا یک دوره کامل مهندسی یادگیری ماشین به صورت پروژه محور را در تابستان ۱۴۰۰ در شرکت رهنما کالج بگذرانم. از صمیم قلب از تمامی منتورها و هم‌گروهی پرتلاشم سرکار خانم پرستو فلک‌افلاکی که مرا در به پایان رساندن این دوره همراهی کردند صمیمانه متشکرم.

محمد هاشمی

دانشجوی کارشناسی مهندسی کامپیوتر

شهریور ۱۴۰۰

---

<sup>1</sup> Artificial intelligence

<sup>2</sup> Machine learning

<sup>3</sup> Start-up

<sup>4</sup> Data scientist

## فهرست مطالب

۱	مقدمه	۱
۲	۱-۱ معرفی سازمان محل کارآموزی	۲
۲	۱-۲ شرح فعالیت و اهداف سازمان محل کارآموزی	۲
۲	۱-۳ شرح کلی فعالیت‌های مرتبط سازمان با رشته تحصیلی	۲
۳	۲ هدف از کارآموزی	۳
۳	۲-۱ فعالیت‌های انجام‌شده	۳
۶	۳ شرح فعالیت‌های کارآموزی	۶
۶	۳-۱ هفته اول - آشنایی با کتابخانه‌های کار با داده در پایتون و API نویسی	۶
۷	۳-۲ هفته دوم - ریاضیات	۷
۱۱	۳-۳ هفته سوم - بهینه‌سازی و مقدمات یادگیری ماشین	۱۱
۱۴	۳-۴ هفته چهارم - دسته‌بندی و مقدمات یادگیری عمیق	۱۴
۱۷	۳-۵ هفته پنجم - CNN ها و RNN ها	۱۷
۱۸	۳-۶ هفته ششم - آشنایی با کلان داده‌ها	۱۸
۱۸	۳-۷ هفته هفتم - ارائه فاز اول پروژه	۱۸
۲۲	۳-۸ هفته هشتم - توسعه پروژه نهایی	۲۲
۲۵	۳-۹ هفته نهم - توسعه پروژه نهایی (ادامه)	۲۵
۲۸	۳-۱۰ هفته دهم - ارائه فاز دوم پروژه نهایی	۲۸
۳۰	۴ خلاصه	۳۰

## چکیده

در این دوره کارآموزی پایه‌ترین و اساسی‌ترین الگوریتم‌ها و روش‌های یادگیری ماشین و یادگیری عمیق مورد مطالعه قرار می‌گیرند. برای درک مفاهیم موجود این حوزه، نیاز به داشتن دانش خوبی در زمینه ریاضیات به خصوص جبرخطی و آمار و احتمالات است. این دوره شامل ۱۰ هفته می‌باشد و هر هفته تمرین و پروژه‌های مربوط به خود را دارد. در هفته آغازین، که به هفته دست‌گرمی نامیده شد، با یادگیری کتابخانه‌های مهم کار با داده در زبان برنامه نویسی پایتون، آمادگی خود را برای شروع دوره کسب نمودیم. همچنین با کتابخانه Flask برای API نویسی نیز آشنا شدیم و یک API ساده طراحی و پیاده‌سازی نمودیم. مهارت‌های ورژن‌سازی به کمک Git نیز آموزش دیده شد تا بتوانیم برای پروژه نهایی دوره از آن برای کد نویسی بصورت گروهی استفاده نماییم. در هفته دوم آموزش در قسمت ریاضیات، با شروع از مفاهیم جبرخطی شروع شد. این هفته تمرین عملی ای وجود نداشت و صرفاً مطالب تئوری در ویدئوهای آموزشی دوره یاد گرفته شد. در ادامه این هفته مطالب مهم در آمار و احتمال مانند متغیرهای تصادفی، انواع توزیع‌های آماری و تخمین زدن به کمک MLE و MAP آموزش داده شد. در هفته سوم نیز مباحث مهم بهینه‌سازی به کمک الگوریتم‌های Gradient Descent و Lagrange بحث و بررسی شد. مبانی یادگیری ماشین و دو الگوریتم پایه‌ی این بخش یعنی رگرسیون و دسته‌بندی دو موضوعی بودن که در ادامه این هفته تدریس شد و همچنین تمرین‌های عملی آن نیز در زبان پایتون نوشته و تحلیل شد. همچنین مباحث مهمی اعم از نحوه مواجهه با داده‌های imbalanced نیز مورد مطالعه قرار گرفت. در هفته بعد مباحث مربوط به الگوریتم‌های مهم غیر نظارتی و Mixture Model ها مورد بررسی قرار گرفت. متریک‌های مهم در تمام الگوریتم‌های نظارتی و همچنین غیرنظارتی نیز تدریس و باهم مقایسه شدند. در این هفته مبحث بسیار مهم یادگیری عمیق شروع به تدریس شد. ابتدا کاربردهای آن و دلیل روی آوردن به الگوریتم‌های یادگیری عمیق به جای روش‌های سنتی یادگیری ماشین در عصر حاضر مورد بررسی قرار گرفت. مباحث مربوط به شبکه‌های عصبی و نحوه یادگیری در این مدل‌ها به کمک روش Back propagation به صورت دقیق مورد مطالعه قرار گرفت. یک تمرین عملی نیز در قسمت پیاده‌سازی شبکه‌های عصبی به کمک فریمورک TensorFlow و Keras نیز انجام شد. در هفته پنجم ابتدا شبکه‌های عصبی پیچشی (CNN) و فواید آن نسبت به شبکه‌های عصبی ساده مورد مطالعه قرار گرفت. معماری‌ها و ساختارهای متفاوتی این نوع شبکه عصبی و نحوه پیاده‌سازی آنها تنها به کمک TensorFlow آموزش داده شد. به همین منظور یک تمرین عملی دیگر برای پیاده‌سازی این نوع شبکه‌ها روی دیتاست ارقام MNIST انجام شد. در انتهای این هفته RNN ها به عنوان یکی دیگر از پرکاربرترین شبکه‌های عصبی بررسی و دلایل استفاده از آن و کاربردهای وسیع آن مورد مطالعه قرار گرفت. در ادامه اتوانکودرهای به عنوان یک روش غیرنظارتی مورد بررسی قرار گرفت و یک تمرین عملی از آن نیز زده شد. در هفته ششم نیز، مباحث مقدماتی کلان داده‌ها از قبیل Hadoop, Map Reduce و Spark در حد آشنایی و معرفی تدریس شدند. از هفته هفتم الی دهم پروژه نهایی این دوره کارآموزی تعریف و انجام شد. شرکت رهنماکالج با همکاری شرکت سنجا، پروژه ای با عنوان تشخیص ناهنجاری و خزنده در وب تعریف نمودند. دیتای این پروژه مجموعه‌ای از لاگ‌های سرور وبسایت سنجا است که هدف پیاده‌سازی یک مدل غیرنظارتی برای تشخیص خزنده یا ناهنجاری بودن یا نبودن درخواست‌هایی است که به این وبسایت زده می‌شود. ما در این پروژه روش‌های متعددی اعم از Local Outlier Factor, Isolation Forest،

Autoencoder ها و ... را پیاده سازی و نتایج آنها را با هم مقایسه نمودیم. اتوانکودر ها بهترین نتایج را به همراه داشتند از آن برای فاز production استفاده نمودیم. این پروژه در دو فاز انجام شد که فاز اول اغلب اختصاص به تحلیل داده و EDA به منظور استخراج و مهندسی ویژگی داده شد. در آخر فاز اول نیز ارائه ای از کارهایی که صورت گرفت به تیم منتور و داوران رهنماکالج داده شد. در فاز دوم نیز مدل های مختلف غیر نظارتی پیاده سازی و مقایسه شدند و در آخر به کمک فریمورک های Flask و ReactJS یک وب اپلیکیشن برای تست بهتر و کاربرپسندانه تر پیاده سازی شد. همچنین ارائه دوم نیز از کارهای انجام شده در هفته دهم صورت گرفت.

به علت محدودیت ۳۰ صفحه ای که برای گزارش کارآموزی در نظر گرفته شده است، امکان بررسی تمامی مطالب فراگرفته شده در این دوره نیست. اما مهمترین نکات و بخصوص قسمت هایی که دارای تمرین های عملی داشته است در این گزارش قابل نوشته شده اند. تمامی کدهای این دوره در ریپازیتوری<sup>۵</sup> گیتهاب بنده قابل مشاهده هستند. علاوه بر آن، کدها و مستندات مربوط به پروژه نهایی در این ریپازیتوری<sup>۶</sup> وجود دارند. برای درک بهتر روند پروژه و کارهایی که برای آن انجام شده به شدت پیشنهاد می شود تا اسلایدهای [فاز اول](#) و [فاز دوم](#) را مطالعه کنید. نمودارها و توجیحاتی که برای کارهایی که انجام شده در آن است که در این گزارش موجود نیست.

---

<sup>5</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege>

<sup>6</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection>

## ۱ مقدمه

یادگیری ماشین یکی از زیرمجموعه های هوش مصنوعی است که به سیستم‌ها این امکان را می‌دهد تا بدون برنامه‌نویسی به صورت صریح توانایی یادگیری و پیشرفت داشته باشند. تمرکز اصلی یادگیری ماشینی بر توسعه برنامه‌های رایانه‌ای است که بتوانند به داده‌ها دسترسی پیدا کنند و از آن‌ها برای یادگیری خود استفاده کنند. یادگیری ماشین نه تنها یک تکنولوژی جدید و کاربردی است بلکه تغییرات بسیار زیادی در دنیا ایجاد کرده‌است و تأثیرات عمیقی در دنیا و بخصوص حوزه اقتصاد و کسب‌وکارها داشته‌است. همچنین با افزایش اهمیت کلان داده‌ها<sup>۷</sup> و پیچیده شدن فرآیند تجزیه و تحلیل آنها، یادگیری ماشین به یک روش اصلی برای رسیدگی به مسائل مربوط به کلان داده‌ها تبدیل شده‌است. در این مواقع یادگیری ماشین می‌تواند به کارهای زیر رسیدگی کند:

- پردازش تصویر، تشخیص چهره، بینایی کامپیوتر و تشخیص شیء.
- پردازش زبان طبیعی، برنامه‌های تشخیص و شناسایی صدا و زبان.
- پیش‌بینی قیمت بازار و سهام بورس.
- زیست‌شناسی محاسباتی، تشخیص تومورهای سرطانی، تعیین توالی DNA و کشف دارو.

به همین دلیل یادگیری ماشینی به خاطر آوردهایی که برای سرمایه‌گذاران دارد و همچنین تحولاتی که در سایر حوزه‌ها می‌تواند ایجاد کند، به یک موضوع داغ تبدیل شده‌است. این مساله باعث شده‌است در سراسر جهان بسیاری از افراد خواستار تحصیل و کسب تجربه در رشته هوش مصنوعی و یادگیری ماشین شوند. در این دوره کارآموزی، ابتدا اساسی‌ترین و محوری‌ترین موضوعات در این علم مورد مطالعه و بررسی قرار می‌گیرد. به دلیل عجین شدن یادگیری ماشین با ریاضیات و آمار و احتمال، پیش از بررسی روش‌ها و الگوریتم‌های یادگیری ماشین، به ریاضیات به خصوص، جبر خطی و آمار و احتمال پیشرفته پرداخته می‌شود. در ادامه مفاهیم یادگیری ماشین و انواع روش‌ها و الگوریتم‌های یادگیری ماشین اعم از یادگیری نظارتی<sup>۸</sup> و بدون نظارت<sup>۹</sup> هم به صورت تئوری و هم عملی مورد بررسی قرار گرفته می‌شود. سپس یکی از مهمترین و به روزترین زیرشاخه‌ای از یادگیری ماشین به نام یادگیری عمیق<sup>۱۰</sup> معرفی می‌شود. به پیچیده شدن توابعی که قرار است توسط ماشین یاد گرفته شوند، بسیاری از مواقع الگوریتم‌های یادگیری ماشین که برای یادگیری نیازمند به ویژگی‌های از پیش تعریف شده توسط انسان هستند و همچنین پیچیده بودن توابع هدف، این الگوریتم‌ها به مرور زمان از کارایی‌شان کاسته شد. در این دوره، یادگیری عمیق و کاربرهای آن به صورت پروژه‌ای و پیشرفته نیز مورد مطالعه قرار گرفته‌است. در انتها نیز، گریزی به موضوع کلان داده‌ها و نحوه‌ی پیاده‌سازی سیستم‌های یادگیری که قرار است برای داده در مقیاس بزرگ استفاده شود، زده می‌شود. یکی دیگر از ویژگی‌هایی که این دوره را متمایز از دیگر دوره‌های کارآموزی می‌کند این است که علاوه بر مهارت‌های فنی و عملی، مهارت‌های نرم اعم از کارگروهی، رزومه‌نویسی و شیوه ارائه تدریس

<sup>7</sup> Big data

<sup>8</sup> Supervised learning

<sup>9</sup> Unsupervised learning

<sup>10</sup> Deep learning

می‌شود ولی در این گزارش این موارد به دلیل دور بودن از مباحث اصلی، مورد بررسی قرار نمی‌گیرد.

## ۱-۱ معرفی سازمان محل کارآموزی

رهنماکالج<sup>۱۱</sup> یک فضای کارآموزی حرفه‌ای است که در بازه‌ی زمانی چند هفته‌ای (برای دوره مهندسی یادگیری ماشین ۱۰ هفته) تجربه‌ی کار واقعی را در پروژه‌های واقعی به شکل تیمی یا انفرادی رقم می‌زند. در هر دوره بهترین متخصص‌های هر رشته، تجربه‌هایشان را به شکل عملیاتی به کارآموزها انتقال می‌دهند، در انجام پروژه پایانی همراهی‌شان می‌کنند و فوت‌وفن‌های کار در دنیای واقعی را آموزش می‌دهند. در این دوره‌های کارآموزی علاوه بر دانش فنی، برای کارآموزها کارگاه‌های مهارت نرم برگزار می‌شود تا بتوانند تکنیک‌های رفتاری و نحوه‌ی تعامل تیمی را تمرین کنند. رهنما اولین دوره‌ای برگزار کرد در سال ۱۳۹۵ به عنوان دوره برنامه‌نویسی بود. از آن پس، با توجه به عملکرد خوب این شرکت، دوره‌های دیگری اعم از مهندسی نرم‌افزار، طراحی UI/UX، دیجیتال مارکتینگ، منابع انسانی و یادگیری ماشین را برگزار می‌کند.

## ۱-۲ شرح فعالیت و اهداف سازمان محل کارآموزی

هدف از دوره‌ی مهندسی یادگیری ماشین رهنماکالج، ایجاد کردن فرصت تجربه واقعی برای اجرای پروژه‌های هوش مصنوعی که در صنعت اتفاق می‌افتد است. برای ورود به این دوره کارآموزی، یک آزمون ورودی که شامل مباحث تئوری (جبرخطی و آمارواحتمال) و یک آزمون عملی<sup>۱۲</sup> (پروژه کوچک یادگیری ماشین) گرفته می‌شود. در این دوره از دوره‌های یادگیری ماشین رهنماکالج، از بین بالای ۷۰۰ نفر شرکت‌کننده تنها ۳۰ نفر وارد این دوره کارآموزی شدند. یکی دیگر از اهداف این سازمان، آماده‌سازی کارآموزان برای ورود به مسیر شغلی مهندسی تحلیل داده و یادگیری ماشین در شرکت‌ها و سازمان‌های معتبر است. با توجه به گفته‌های مدیر رهنماکالج، ورود پیدا کردن کارآموزان به شرکت‌های معتبر داخلی پس از اتمام دوره کارآموزی، از افتخارات و دلایل پیشرفت این سازمان است. همچنین از دیگر فعالیت‌های این سازمان می‌توان به معرفی کردن کارآموزان توانمند به شرکت‌های معتبر اشاره کرد. با توجه به پیشرفت و عملکرد بسزایی که رهنماکالج در این چند سال اخیر داشته است، شرکت‌های معتبری درخواست نیروی انسانی از این سازمان می‌کنند.

## ۱-۳ شرح کلی فعالیت‌های مرتبط سازمان با رشته تحصیلی

به طور کلی اغلب فعالیت‌های رهنماکالج حاوی جنبه‌های آموزشی است. در واقع کار اصلی رهنماکالج، کمک به افراد با استعداد و با انگیزه برای ورود به مسیر شغلی موردعلاقه‌شان است. به طور خاص در دوره کارآموزی یادگیری ماشین، کارآموزان مهارت‌های لازم برای کار کردن در مشاغل مربوط به هوش مصنوعی را فرا می‌گیرند و طبق آمار و ارقامی که این سازمان اعلام کرده است، ۹۵ درصد از کارآموزان پس از اتمام دوره در یک شرکت معتبر داخلی استخدام می‌شوند.

<sup>11</sup> <https://rahnemacollege.com/>

<sup>12</sup> <https://github.com/mohammadhashemii/Justice-In-Work>

## ۲ هدف از کارآموزی

هدف اصلی، علاقه بسیار زیاد به مباحث یادگیری ماشین و یادگیری عمیق است! امروزه با توجه به فراگیر شدن هوش مصنوعی و به خصوص یادگیری ماشین در تمامی حوزه‌ها، کمتر کسی وجود دارد که به این حوزه علاقه پیدا نکند. بنده به طور خاص مطالعات خود را در این زمینه از اواخر ترم سوم (زمستان ۹۹) شروع کردم و شروع به انجام چند پروژه تحقیقاتی به کمک اساتید دانشگاه کردم. دلیل انتخابم از مهندسی یادگیری ماشین به عنوان دوره کارآموزی، آمادگی برای وارد شدن به صنعت و انجام پروژه‌هایی است که در زندگی واقعی مان در حال رخ دادن می‌باشند. به طور خاص، دوره یادگیری ماشین رهنماکالج، دارای محتوایی است که تمام نیازها برای یک مهندس یادگیری ماشین یا دانشمند داده را برطرف می‌کند. مهارت‌های تخصصی که در این دوره کارآموزی توسط کارآموز کسب می‌شوند عبارتند از:

- جبر خطی: تسلط بر جبر خطی مورد نیاز و تجزیه‌ها
- آمار و احتمالات در یادگیری ماشین: تسلط بر توابع احتمالاتی و محاسبه MLE و MAP
- بهینه سازی<sup>۱۳</sup>: شناخت توابع بهینه‌سازی و نحوه عملکرد آنها
- مبانی هوش مصنوعی: شناخت مسائل مرتبط با هوش مصنوعی و یادگیری ماشین و ارتباط آنها
- الگوریتم‌های یادگیری ماشین: بررسی الگوریتم‌های مختلف، روش حل و نحوه پیاده‌سازی
- شبکه‌های عصبی<sup>۱۴</sup> و یادگیری عمیق: درک شبکه‌های عصبی و عمیق، الگوریتم‌های یادگیری آنها و تسلط بر مسائل مختلف
- آشنایی با کلان‌داده‌ها: آشنایی با عملگرهای اصلی در پردازش کلان‌داده‌ها، Hadoop و Spark
- کتابخانه‌های پایتون: آشنایی با کتابخانه‌های مختلف پایتون برای کار کردن در حوزه یادگیری ماشین.

در ادامه در هر قسمت به یکی از موضوعات گفته شده پرداخته می‌شود. در این دوره قسمت ریاضیات برای من از اهمیت زیادی برخوردار بود. در هیچ کارآموزی شرکت دیگری، ریاضیات مورد نیاز برای یادگیری ماشین آموزش داده نمی‌شود در صورتی که پایه و اساس الگوریتم‌ها و روش‌های حل مساله مسائل مربوطه، نیازمند فهم عمیقی از مباحث جبر خطی و آمار و احتمال است. از دیگر اهداف من از این دوره‌ی کارآموزی این بود که چگونه مسائل واقعی ای که در صنعت وجود دارند با مباحث تئوری موجود در یادگیری ماشین ارتباط برقرار می‌کنند. پروژه نهایی این دوره کارآموزی پاسخ خوبی به این سوال بود.

## ۲-۱ فعالیت‌های انجام شده

اغلب زمان صرف شده در این دوره کارآموزی صرف یادگیری مفاهیم و پیاده سازی الگوریتم‌های یادگیری ماشین شده است. در ۲ الی ۳ هفته آخر این دوره، پیاده سازی پروژه نهایی و

<sup>13</sup> Optimization

<sup>14</sup> Artificial Neural Network (ANN)



استفاده از مفاهیم فراگرفته شده در یک پروژه واقعی به صورت عملی وظیفه‌ی ما کارآموزان بود. در جدول زیر به تفکیک هفته، موضوعات مباحث آموزشی و وظایفی که بر عهده کارآموزان بود، قابل مشاهده است:

هفته	موضوع	فعالیت‌ها و وظایف
اول	آشنایی با کتابخانه‌های کار با داده در پایتون و API نویسی	<ul style="list-style-type: none"> <li>مرور و یادگیری کتابخانه Pandas و Scipy به همراه تمرین عملی.</li> <li>یادگیری استفاده از ابزارهای مصورسازی<sup>۱۵</sup> در پایتون.</li> <li>یادگیری کتابخانه Flask برای API نویسی.</li> <li>یادگیری و تسلط بر Git.</li> </ul>
دوم	ریاضیات	<ul style="list-style-type: none"> <li>مشاهده ویدئوهای مربوط به آموزش ریاضیات مورد نیاز در یادگیری ماشین. اعم از جبرخطی و آمارو احتمال.</li> <li>تاکید به تسلط به Decomposition ها در جبرخطی.</li> <li>درک عمیق از مفاهیم مربوط به PCA ها.</li> <li>تخمین زدن پارامترها توسط روش‌های MLE و MAP.</li> </ul>
سوم	بهینه‌سازی و مقدمات یادگیری ماشین	<ul style="list-style-type: none"> <li>یادگیری مفاهیم پایه و ضروری بهینه‌سازی یعنی Gradient Descent و Lagrange.</li> <li>آشنایی با مفاهیم اولیه یادگیری ماشین.</li> <li>یادگیری اولین الگوریتم یادگیری ماشین به نام رگرسیون<sup>۱۶</sup> و پیاده‌سازی آن در پایتون.</li> <li>یادگیری تئوری و مفاهیم بسته‌بندی در یادگیری ماشین.</li> <li>یادگیری متریک‌های یادگیری ماشین و نحوه استفاده از آنها برای دیتاهای imbalanced.</li> </ul>

<sup>۱۵</sup> Visualization

<sup>۱۶</sup> Regression

<ul style="list-style-type: none"> <li>• ادامه مبحث دسته‌بندی و پیاده سازی یک پروژه کوچک در پایتون.</li> <li>• یادگیری مفاهیم مربوط به Mixture Model ها.</li> <li>• آشنایی با مفاهیم پایه یادگیری عمیق و دلایل پیدایش آن.</li> <li>• شبکه‌های عصبی و یادگیری Backpropagation به عنوان روش یادگیری آنها.</li> <li>• پیاده‌سازی پروژه تشخیص دستخط روی دیتاست MNIST.</li> </ul>	دسته‌بندی <sup>۱۷</sup> و مقدمات یادگیری عمیق	چهارم
<ul style="list-style-type: none"> <li>• یادگیری و تمرین با Framework های یادگیری عمیق: Tensorflow</li> <li>• پیاده سازی Autoencoder ها به عنوان یک مدل غیرنظارتی.</li> <li>• آشنایی با LSTM ها به عنوان یکی از قویترین مدل های RNN</li> <li>• CNN ها و پیاده سازی به کمک کتابخانه Keras.</li> <li>• آشنایی با GAN ها.</li> <li>• آشنایی با مفاهیم مقدماتی Reinforcement Learning</li> </ul>	CNN ها و RNN ها	پنجم
<ul style="list-style-type: none"> <li>• Map Reduce</li> <li>• Hadoop ,Spark</li> </ul>	آشنایی با کلان داده ها	ششم
<ul style="list-style-type: none"> <li>• تعریف شدن پروژه پایانی.</li> <li>• تحلیل اکتشافی داده<sup>۱۸</sup></li> <li>• تمیزسازی داده.</li> <li>• آماده‌سازی اسلاید برای ارائه فاز اول.</li> </ul>	ارائه فاز اول پروژه	هفتم
<ul style="list-style-type: none"> <li>• تعریف و آموزش دادن مدل‌های غیرنظارتی پلایه روی ویژگی‌های استخراج یافته.</li> <li>• مقایسه مدل‌های مختلف و ارزیابی<sup>۱۹</sup>.</li> </ul>	توسعه پروژه نهایی	هشتم و نهم
<ul style="list-style-type: none"> <li>• ساخت یک API برای تست کردن مدل آموزش دیده.</li> <li>• آماده‌سازی اسلاید برای ارائه نهایی.</li> </ul>	ارائه فاز دوم پروژه نهایی	دهم

<sup>17</sup> Classification

<sup>18</sup> Exploratory Data Analysis (EDA)

<sup>19</sup> Evaluation

جدول ۱-۲: چشم‌اندازی از فعالیت‌ها و وظایف کارآموزی به تفکیک هفته

### ۳ شرح فعالیت‌های کارآموزی

در این قسمت، به ترتیب هفته، فعالیت‌هایی که در این دوره کارآموزی انجام شده‌است به صورت خلاصه شرح داده می‌شود:

#### ۳-۱ هفته اول – آشنایی با کتابخانه‌های کار با داده در پایتون و API نویسی

در این هفته آمادگی‌های لازم برای ورود به مباحث آموزشی اصلی کسب شد. ابتدا با چند کتابخانه پایتون برای کار و تحلیل داده، آموزش‌های لازم برای تمیزسازی داده<sup>۲۰</sup> و تحلیل اکتشافی داده (EDA) داده شد. در انتهای این هفته، یک تمرین برنامه‌نویسی (البته اختیاری) برای پیاده‌سازی یک API به کمک کتابخانه Flask انجام شد. شرح جزئی و نکاتی مهمی که از مطالب این هفته استخراج شده است را در ادامه می‌بینیم:

#### ۳-۱-۱ آشنایی با Pandas<sup>۲۱</sup>

در این قسمت مروری بر دستورات مهم در این کتابخانه خواهیم داشت. مطالب این قسمت از یک دوره آنلاین از وبسایت دیتاکمپ<sup>۲۲</sup> گرفته شده است.

کتابخانه Pandas بر روی دو کتابخانه دیگر پایتون به نام‌های Numpy<sup>۲۳</sup> و Matplotlib<sup>۲۴</sup> که خود از مهمترین کتابخانه‌های محاسبات عددی و ماتریسی و تصویرسازی الگوریتم‌های یادگیری ماشین محسوب می‌شوند. در پژوهش‌های حوزه علوم کامپیوتر و پروژه‌های صنعتی، اغلب، داده‌های جدولی<sup>۲۵</sup>، برای توصیف و نمایش داده به کار می‌روند. به همین دلیل Pandas نیز این قالب برای نمایش داده استفاده می‌کند که به آن DataFrame گفته می‌شود. در ادامه به چند مورد از پرکاربردترین و اساسی‌ترین دستورات تعریف شده برای دیتافریم‌ها می‌پردازیم:

دستور	توضیح و عملکرد دستور
<code>df.head()</code>	نمایش دادن چند سطر ابتدایی دیتافریم.
<code>df.info()</code>	نمایش اسامی ستون‌ها، جنس داده ذخیره شده در هر ستون، و اینکه آیا دارای

<sup>20</sup> Data cleaning

<sup>21</sup> <https://pandas.pydata.org/>

<sup>22</sup> <https://www.datacamp.com/courses/data-manipulation-with-pandas>

<sup>23</sup> <https://numpy.org/>

<sup>24</sup> <https://matplotlib.org/>

<sup>25</sup> Tabular data

Missing value هست یا خیر.	
محاسبه اساسی‌ترین روابط آماری برای دیتافریم مانند میانگین، میانه و ...	df.describe()
مقادیر ذخیره شده در ستون‌ها که به صورت یک آرایه نامپای دوبعدی برمی‌گردد.	df.values

جدول ۳-۱: چند مورد از پرکاربردترین دستورات Pandas

همچنین توابع و روش‌هایی نیز برای پیاده‌سازی توابع دلخواه ما وجود دارد مانند مرتب‌سازی، ادغام کردن دیتافریم‌ها و ... که این روش‌ها در ریپازیتوری<sup>۲۶</sup> گیت‌هاب بنده قابل مشاهده هستند.

## ۳-۱-۲ آشنایی با Flask<sup>۲۷</sup>

هدف از این قسمت آشنایی و تمرین برنامه‌نویسی برای API نویسی به زبان پایتون است. ابتدا نحوه‌ی نصب و راه‌اندازی کتابخانه Flask در [اینجا](#) آمده‌است. یک پروژه ساده به عنوان تمرین این بخش در ریپازیتوری<sup>۲۸</sup> بنده نیز قابل مشاهده است.

## ۳-۲ هفته دوم – ریاضیات

در این قسمت به مهمترین مباحث ریاضیات در یادگیری ماشین یعنی جبرخطی و آمارو احتمال می‌پردازیم. قبل از شروع جبرخطی، خوب است که کاربردهای آنرا در یادگیری ماشین و هوش مصنوعی بدانیم. این کاربردها در مواردی مانند رگرسیون خطی، پردازش تصویر، Regularization<sup>۲۹</sup>، PCA<sup>۳۰</sup> و SVD<sup>۳۱</sup> و سایر موارد که در ادامه به آنها می‌پردازیم مشهود هستند.

### ۳-۲-۱ جبرخطی

جزوه‌ای از مباحث تدریس شده از این قسمت نوشته شده است که در ریپازیتوری<sup>۳۱</sup> بنده قابل مشاهده است. در اینجا چکیده‌ای از این مطالب را مرور می‌کنیم.

#### ۳-۲-۱-۱ فضای برداری<sup>۳۲</sup>

تعریف: یک فضای برداری روی اعداد حقیق تشکیل شده است از: الف) یک مجموعه  $V$  که اعضای آن بردار نامیده می‌شوند. ب) عملگر باینری جمع<sup>۳۳</sup> روی بردارها شروط زیر را ارضا کند:

<sup>۲۶</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-0#1-pandas>

<sup>۲۷</sup> <https://flask.palletsprojects.com/en/2.0.x/>

<sup>۲۸</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-0#2-flask>

<sup>۲۹</sup> Principal Component Analysis

<sup>۳۰</sup> Singular-Value Decomposition

<sup>۳۱</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-1>

<sup>۳۲</sup> Vector space

<sup>۳۳</sup> Addition operation

۱. خاصیت جابجایی.

۲. خاصیت شرکت پذیری.

۳. وجود داشته باشد یک بردار صفر که برای هر بردار  $v \in V: v + 0 = v$ .

۴. به ازای هر بردار  $v \in V$  یک بردار یکتای  $-v$  وجود داشته باشد که  $v + (-v) = 0$ .

ج) عملگر دیگری به نام عملگر ضرب اسکالر<sup>۳۴</sup>، اگر  $a \in \mathcal{R}$  و  $v \in V$  آنگاه  $a.v \in V$ . همچنین شروط زیر را ارضا کند:

۱. برای هر  $v \in V$  باید  $1.v = v$

۲.  $(\alpha_1 . \alpha_2). v = \alpha_1 . (\alpha_2 . v)$

۳.  $\alpha(v_1 + v_2) = \alpha.v_1 + \alpha.v_2$

۴.  $(\alpha_1 + \alpha_2). v = \alpha_1.v + \alpha_2.v$

### ۳-۲-۱-۲ وابستگی خطی<sup>۳۵</sup>

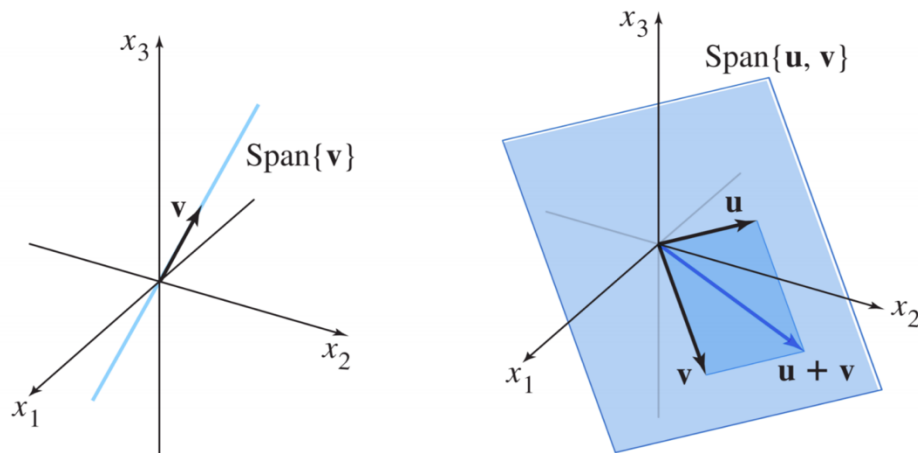
تعریف ۱: بردار  $v$  را ترکیب خطی بردارهای  $\{v_1, v_2, \dots, v_n\}$  می نامیم اگر:

$$v = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n, \alpha_i \in \mathbb{R}$$

تعریف ۲: مجموعه‌ی همه بردارهای  $v$  ملند بالا را پوشش خطی<sup>۳۶</sup> مجموعه

$\{v_1, v_2, \dots, v_n\}$  می نامیم و بصورت زیر نمایش می دهیم:

$$\text{span}(v_1, v_2, \dots, v_n) = \langle v_1, v_2, \dots, v_n \rangle$$



شکل ۳-۱: مثال هایی از  $\text{span}$

<sup>34</sup> Scalar multiplication

<sup>35</sup> Linear dependency

<sup>36</sup> Linear span

نکته: هر  $\text{span}(v_1, v_2, \dots, v_n)$  خود یک فضای برداری روی  $\mathbb{R}$  به حساب می‌آید. به این گونه مجموعه‌ها زیرفضای  $v$ <sup>۳۷</sup> نیز گفته می‌شود.

نکته:  $\text{span}(v_1, v_2, \dots, v_n)$  کوچکترین زیرفضا برای  $v$  است که شامل  $\{v_1, v_2, \dots, v_n\}$  است.

تعریف وابستگی خطی: مجموعه  $\{v_1, v_2, \dots, v_n\} \subseteq V$  را وابسته خطی می‌نامیم اگر  $\{\alpha_1, \alpha_2, \dots, \alpha_n\} \in \mathbb{R}$  ای وجود داشته باشد که هیچ کدام صفر نباشند و:

$$0 = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n$$

به عبارت دیگر حداقل یکی از بردارهای  $v_i$  بتواند بصورت ترکیب خطی بقیه بردارها نوشته شود.

تعریف استقلال خطی: مجموعه  $\{v_1, v_2, \dots, v_n\} \subseteq V$  را مستقل خطی می‌نامیم اگر  $\{\alpha_1, \alpha_2, \dots, \alpha_n\} \in \mathbb{R}$  ای وجود داشته باشد که:

$$0 = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n$$

آنگاه همه‌ی  $\alpha_i$  ها صفر باشند.

تعریف پایه: هر مجموعه مستقل خطی  $\{v_1, v_2, \dots, v_n\} \subseteq V$  که  $\text{span}$  را  $V$  کند، یک پایه برای فضای برداری  $V$  به حساب می‌آید.

به علت حجم شدن متن گزارش امکان بررسی همه‌ی مطالب گفته شده در این قسمت وجود ندارد. جزوه کاملتر آن در ریپازیتوری<sup>۳۸</sup> گیت‌هاب بنده قابل مشاهده است.

## ۳-۲-۲ آمار و احتمال

مباحث مختلفی در این قسمت در دوره کارآموزی تدریس داده شده است. اما در این گزارش به طور مختصر به یکی از مهمترین قسمت‌های آن‌ها کاربرد بسیار زیادی در مباحث تئوری یادگیری ماشین دارد می‌پردازیم.

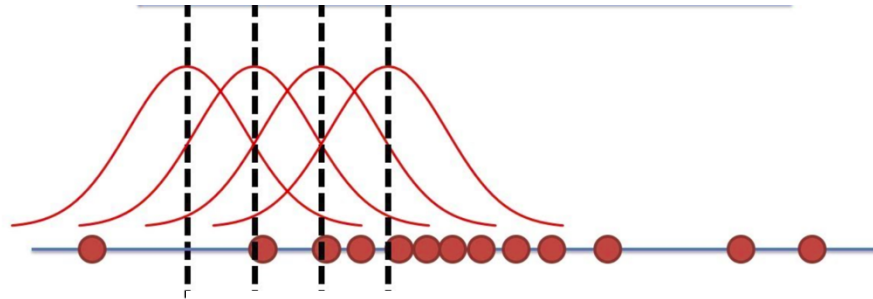
### Maximum Likelihood Estimation (MLE) ۳-۲-۲-۱

این روش برای تخمین زدن پارامترهای یک متغیر تصادفی است به کمک وجود تعداد

<sup>37</sup> Subspace

<sup>38</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-1>

محدودی از نمونه‌های آن متغیر تصادفی به کار می‌رود. یعنی نمونه‌هایی که داریم دارای توزیع مشابه با توزیع متغیر تصادفی است که قصد به تخمین زدن آن داریم. برای درک بهتر این روش، آنرا با یک مثال توضیح می‌دهیم. فرض کنید  $D = (X_1, X_2, \dots, X_n) \sim \mathcal{N}(\mu, \sigma^2)$  تا از نمونه‌هایمان هستند که دارای توزیع احتمال نرمال با میانگین  $\mu$  و انحراف معیار  $\sigma$  است. همچنین مقدار  $\sigma^2$  معلوم مساله است. هدف ما تخمین زدن پارامتر  $\mu$  می‌باشد.



نمودار ۱-۳: حالت‌های متفاوت برای پارامتر  $\mu$

در نمودار ۱-۳ می‌توان حالات متفاوت برای اندازه  $\mu$  را مشاهده کرد که به طبع برای هر  $\mu$ ، یک توزیع متفاوت برای تابعی که می‌خواهیم تخمین بزنیم خواهیم داشت. دایره‌های قرمز نشان‌دهنده نمونه‌هایی است داریم. اگر پارامترهای مساله را در یک برداری مثل بردار  $\theta$  قرار دهیم آنوقت داریم:

$$p_{\theta}(x) := p(x|\theta)$$

$$p_{\mu, \sigma^2}(x) := p(x|\mu, \sigma^2)$$

که بهترین پارامتر، پارامتری است که احتمال  $p(D|\theta)$  را بیشینه کند.

تعریف تابع *likelihood*

$$\mathcal{L}(\theta|x) := p(D|\theta) = p(x_1, x_2, \dots, x_n|\theta) = \prod_{i=1}^n p(x_i|\theta)$$

بردار  $\theta$  ای که این تابع را بیشینه کند به عنوان  $MLE$  از  $\theta$  شناخته می‌شود و به عنوان  $\theta_{MLE}$  نشان داده می‌شود. برای اینکه محاسبات راحت‌تر شود معمولاً از تابع بالا یک لگاریتم گرفته می‌شود:

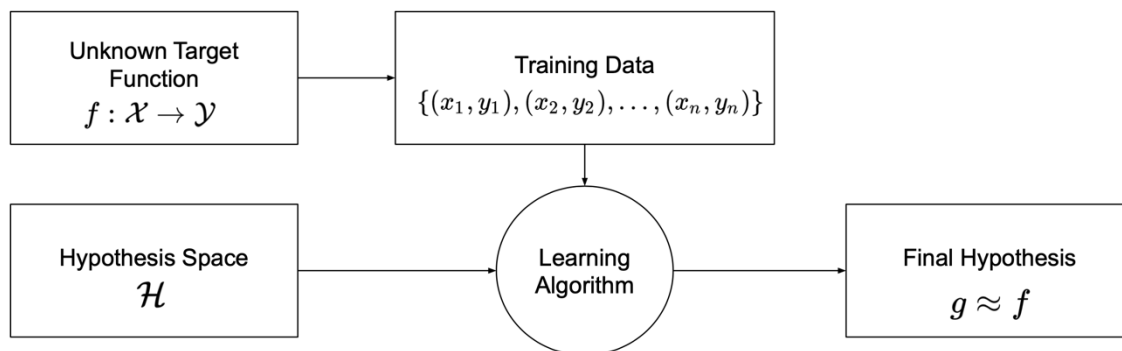
$$\log(\mathcal{L}(\theta|x)) = \sum_{i=1}^n \log(p(x_i|\theta))$$

برای پیدا کردن بهترین پارامتر، باید از تابع بالا مشتق گرفت و برابر صفر قرار داد. مثال حل

شده آن در اسلایدهای تدریس شده در دوره کارآموزی موجود می‌باشد.

### ۳-۳ هفته سوم - بهینه‌سازی و مقدمات یادگیری ماشین

اگر بخواهیم تعریف کلی از یادگیری ماشین داشته باشیم باید بگوییم هدف این علم، جستجو برای بهترین تابعی است که ورودی‌های داده‌های آموزش را به خروجی‌هایشان نگاشت می‌کند. همچنین روی داده آزمون نیز تعمیم پذیر باشد. چگونه به دنبال این تابع بگردیم؟ باید ابتدا یک فضای فرضیه<sup>۳۹</sup> که مجموعه‌ای از توابع است را انتخاب کنیم. یعنی به عنوان مثال مجموعه توابع خطی به فرم  $f = wx + b$ . برای اینکه ارزیابی‌ای روی توابع این مجموعه داشته باشیم باید یک تابع هزینه<sup>۴۰</sup> تعریف می‌کنیم تا مجموع خطای نگاشت هر نمونه از دیتاست را حساب کنیم. طبیعتاً تابعی که کمترین هزینه (مجموع خطا) را داشته باشد، مقصود ماست.



شکل ۳-۲: فرآیند کلی در الگوریتم‌های یادگیری ماشین

مسائل موجود در یادگیری ماشین به دسته‌های زیر تقسیم می‌شوند.

- یادگیری با نظارت (supervised): برچسب داده‌ها موجود است.
- یادگیری غیرنظارتی (unsupervised): برچسب داده‌ها موجود نیست.
- یادگیری تقویتی (reinforcement): تعامل با محیط و یادگیری آن رفتار کردن.

البته دسته‌های دیگری نیز موجود هستند ولی پایه و اساس یادگیری ماشین بر روی همین ۳ دسته می‌چرخد. به دلیل حجیم‌شدن متن گزارش امکان بررسی تمامی الگوریتم‌ها بصورت جزئی در اینجا نیست اما در ادامه به برخی از آنها و پیاده‌سازی‌هایشان در زبان پایتون به صورت خلاصه می‌پردازیم:

#### ۳-۳-۱ رگرسیون<sup>۴۱</sup>

در دسته‌ی الگوریتم‌های با نظارت قرار می‌گیرد و فضای خروجی نیز پیوسته است. یکی از

<sup>۳۹</sup> Hypothesis space

<sup>۴۰</sup> Cost function

<sup>۴۱</sup> Regression



مثال‌های معروف آن تشخیص قیمت خانه از روی ویژگی‌های آن است. در جدول پایین ویژگی‌ها و در ستون آخر قیمت هر خانه به عنوان برچسب آمده است.

X1 (size)	X2 (age)	X3 (rooms)	Y (price)
100	10	3	200
200	2	4	300
40	10	1	40
50	0	1	70
...	...	...	...

جدول ۲-۳: دیتاست مربوط به تشخیص قیمت خانه

مانند همه‌ی الگوریتم‌های یادگیری ماشین، باید یک فضای فرضیه برای حل مساله رگرسیون انتخاب شود که می‌توان آنرا به صورت  $y = w^T x + w_0$  نوشت. به این معنی که اگر نمونه‌ها دارای یک ویژگی باشند، فضای فرضیه یک خط، اگر دارای دو ویژگی یک صفحه و برای بیشتر از دو ویژگی یک ابرصفحه<sup>۴۲</sup> است. هدف پیدا کردن بردار وزن‌های  $w$  و عرض از مبدا  $w_0$  به عنوان مجهول‌های مساله است. به همین دلیل، با استفاده از روش پیمایشی<sup>۴۳</sup>، ابتدا بردار وزن‌ها و عرض از مبدا را مقداردهی اولیه می‌کنیم و به کمک تعریف یک تابع خطا<sup>۴۴</sup>، خطای تابع تخمین‌زده شده، آنرا ارزیابی و در صورت نیاز وزن‌ها را به روز رسانی می‌کنیم. تابع خطایی که معمولاً برای رگرسیون در نظر گرفته می‌شود تابع مجموع مربعات خطا است که به اختصار به  $MSE$  معروف می‌باشد. نحوه‌ی به‌روزرسانی وزن‌ها به کمک بهینه‌سازی تابع خطا با روش Gradient Descent صورت می‌گیرد که بدلیل حجیم شدن متن گزارش جزئیات و محاسبات ریاضی آن در اینجا نمی‌پردازیم.

به عنوان یک تمرین عملی برای بخش رگرسیون روی دیتاست Boston برای تشخیص قیمت خانه، یک پیاده‌سازی رگرسیون خطی به کمک کتابخانه Sickit-learn در ریپازیتوری گیت‌هاب<sup>۴۵</sup> بنده قابل مشاهده است.

## ۲-۳-۳ متریک‌های ارزیابی برای رگرسیون

متریک‌های متفاوتی برای ارزیابی مدل رگرسیون وجود دارد که هرکدام مزایا و عیب‌های خود را دارد.

۱.  $MSE$ : میانگین مجموع مربعات خطا است. می‌توان گفت رایج‌ترین است اما عیبی که دارد این است که شهودی به ما از فاصله نمی‌دهد. (زیرا فرمول فاصله یک جذر هم دارد)

<sup>۴۲</sup> Hyper space

<sup>۴۳</sup> Iterative

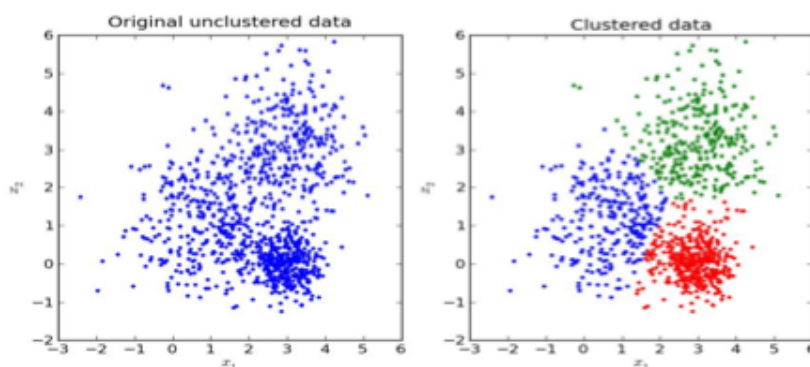
<sup>۴۴</sup> Loss function

<sup>۴۵</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-2>

۲. RMSE: جذر میانگین مجموع مربعات خطا است که این بار دقیقاً تعریف فاصله اقلیدوسی است که شهود دقیقتری از مفهوم فاصله برایمان ایجاد می‌کند
۳. MAE: میانگین قدرمطلق خطا است.
۴. MAPE: میانگین درصدی قدرمطلق خطاست. یعنی به صورت درصدی قدرمطلق اختلاف مقدار مطلوب با مقدار پیش‌بینی شده حساب می‌گردد که بدلیل درصدی بودن شهود بهتری از عملکرد مدل می‌توان گرفت.
۵. R2-score: یکی از عیوب متریک MSE این است که لزوماً نمی‌توان نتیجه گرفت که اگر MSE یکبار ۱۰ شد بهتر از موقعی است که ۱۰۰ شده است. زیرا میزان بزرگی و scale داده‌ها در آن لحاظ نمی‌شود. ولی در R2-score دقیقاً این موضوع لحاظ می‌شود.

### ۳-۳-۳ خوشه بندی<sup>۴۶</sup>

این نوع الگوریتم‌ها از روش‌های غیرنظارتی محسوب می‌شوند به عبارتی نیاز به وجود برچسب برای داده‌ها وجود ندارد. و وظیفه مدل این است که با توجه به توابع شباهتی که تعریف می‌شوند، شباهت بین نمونه از روی فضای ورودی را محاسبه و داده‌های را دسته‌های مختلف خوشه بندی کند.



شکل ۳-۳: قبل و بعد از اجرای الگوریتم خوشه بندی

برای خوشه‌بندی الگوریتم‌های متعددی معرفی شده‌اند که یکی از معروفترین‌های آنها الگوریتم K-means است. هدف این الگوریتم، خوشه بندی دیتا به  $k$  دسته است که نمونه‌های داخل یک کلاس کمترین واریانس را داشته باشند و مرکز هر دسته نسبت به دسته‌های دیگر بیشترین فاصله را داشته باشند. به همین منظور می‌توان گفت هدف این روش کمینه کردن تابع خطای زیر است. اگر نمونه  $i$  ام عضو کلاستر  $k$  باشد آنگاه  $r_{ik} = 0$  و خواهیم داشت:

$$loss = \sum_{n=1}^n \sum_{k=1}^k r_{ik} ||x_n - \mu_k||^2$$

پیدا کردن بهترین جواب یک مساله NP-hard محسوب می‌شود برای همین یک روش بازگشتی برای این

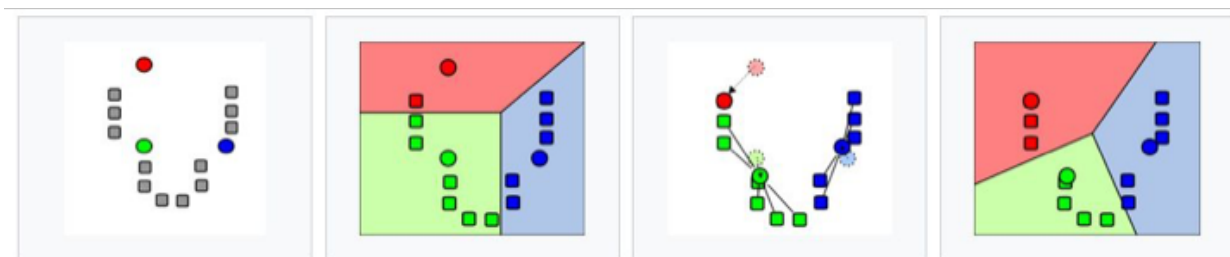
<sup>46</sup> Clustering

مساله پیشنهاد شده است که به شرح زیر می باشد:

قدم ۱: مشخص کردن  $k$  نقطه به عنوان مرکز دسته ها به صورت تصادفی.

قدم ۲: اختصاص دادن برچسب نمونه به برچسب نزدیک تری مرکز به آن نمونه.

سپس دوباره قدم ۱ را تکرار میکنیم (البته در مراحل بعدی، قدم ۱ دیگر به صورت تصادفی نخواهد بود) پس از تکرار چندین مرحله، دیگر تفاوت خاصی در مراکز دسته ها و برچسب نمونه ها صورت نمی گیرد و الگوریتم در آنجا متوقف می شود و نمونه ها به  $k$  دسته خوشه بندی شده اند.



شکل ۳-۴: مراحل الگوریتم  $K$ -means

### ۳-۴ هفته چهارم – دسته بندی<sup>۴۷</sup> و مقدمات یادگیری عمیق

دسته بندی از الگوریتم های با نظارت یادگیری ماشین محسوب می شود و این بار خروجی مساله برخلاف مساله رگرسیون، گسسته است. یعنی قرار است در فاز پیش بینی، مدل با دیدن نمونه به همراه ویژگی هایش متعلق به کدام یک از دسته های موجود است. بدلیل زیاد شدن حجم متن مقاله، بررسی تمام جزییات و مفاهیم تئوری بحث دسته بندی در اینجا امکان پذیر نیست. اما به عنوان یک پیاده سازی تمرین عملی دسته بندی روی دیتاست Titanic در ریپازیتوری<sup>۴۸</sup> گیتهاب بنده قابل مشاهده است.

#### ۳-۴-۱ آشنایی با یادگیری عمیق

هدف اصلی الگوریتم های یادگیری عمیق، تخمین تابع به کمک ترکیب تابع های ساده تر و ساخت تابع های پیچیده تر (معمولا غیر خطی) است. در یادگیری ماشین می بایستی فضای فرضیه در ابتدای کار مشخص شود ولی در الگوریتم های یادگیری به کمک شبکه های عصبی، فضای فرضیه از روی ترکیب توابع ساده تر غیر خطی به صورت خودکار ساخته می شود. به طور مثال اگر فرض کنیم ۳ تابع غیر خطی داریم، آنگاه فضای فرضیه به صورت زیر نوشته می شود:

<sup>47</sup> Classification

<sup>48</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-3>

$$\mathcal{F} = \{f(X) = W_3 g(W_2 g(W_1 X))\}$$

که در آن  $g$  یک تابع ساده غیرخطی مانند  $ReLU$ ،  $tanh$  و یا ... است.

اگر بخواهیم دلایل روی آوردن به الگوریتم‌های یادگیری عمیق به جای روش‌های سنتی یادگیری ماشین متذکر شویم، می‌توان به موارد زیر اشاره نمود:

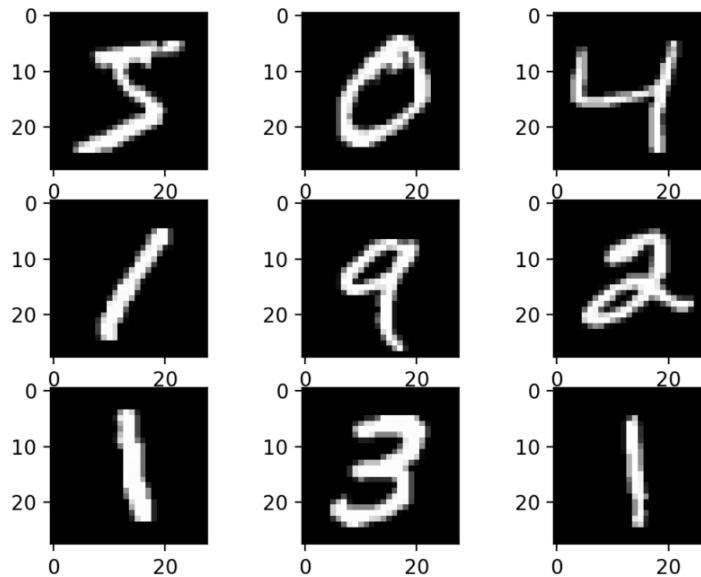
- عدم نیاز به *Feature selection* و همچنین *Feature extraction*: یکی از مراحل که در اغلب تسک‌های یادگیری ماشین صورت می‌گیرد بحث استخراج ویژگی توسط انسان است. قبل از آموزش مدل‌های یادگیری، باید مجموعه‌ای از ویژگی‌ها توسط کاربر مشخص شود تا فرآیند آموزش صورت بگیرد. همچنین برخی از این ویژگی‌ها بدلائل ممکن است تاثیر خاصی رو فرآیند آموزش نداشته باشد یا حتی موجب بیش‌برازش<sup>۴۹</sup> شود که توسط روش‌های موجود در حوزه *feature selection* مانند محاسبه وابستگی<sup>۵۰</sup> ویژگی‌ها به خروجی اغلب با رابطه پیرسون، برخی از ویژگی‌ها حذف می‌شوند. این دو فرآیند در الگوریتم‌های یادگیری عمیق عملاً وجود ندارند.
- هزینه‌بر بودن الگوریتم‌های یادگیری ماشین: همین بحث استخراج ویژگی در اغلب موارد فرآیند هزینه‌بر و زمانبری صورت می‌شود. زیرا باید به صورت دستی و توسط انسان تعریف شوند و همچنین در دیتاست‌های با مقیاس بزرگ گاهی مواقع این فرآیند نشدنی است.

### ۱-۴-۳ پیاده سازی شبکه‌های عصبی

متأسفانه توضیح الگوریتم‌های شبکه‌های عصبی و نحوه‌ی آموزش آنها با روش‌هایی مثل back propagation در این گزارش با توجه به محدودیت صفحات گزارش امکان‌پذیر نیست. اما در این قسمت پیاده‌سازی یک شبکه عصبی با لایه‌های fully connected به کمک فریمورک‌های TensorFlow و Keras برای تشخیص ارقام روی عکس‌های دستخط از دیتاست معروف MNIST خواهیم داشت.

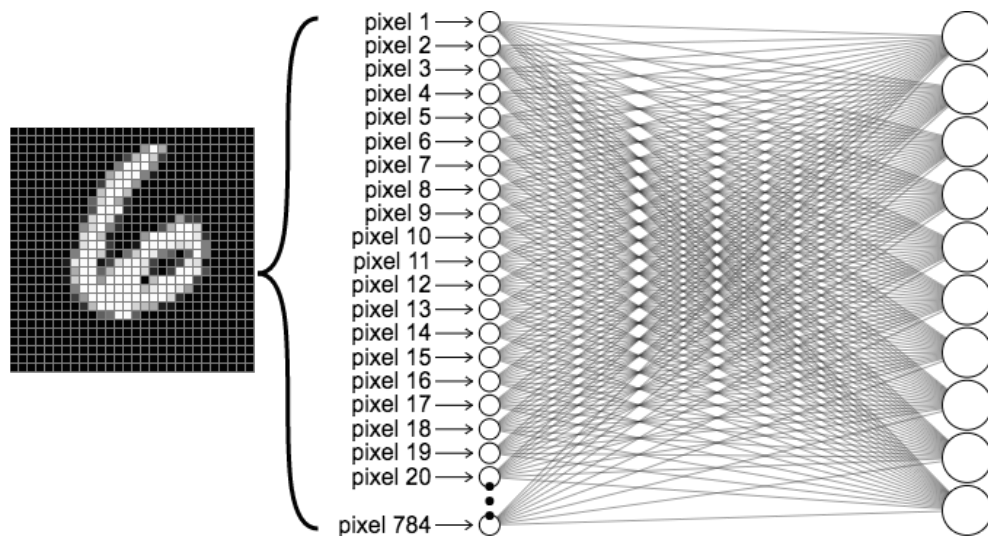
<sup>49</sup> Overfitting

<sup>50</sup> Correlation



شکل ۵-۳: دیتاست ارقام دستنویس *MNIST*

دیتاست *MNIST* دارای ۶۰۰۰۰ نمونه به عنوان داده آموزش و ۱۰۰۰۰ داده به عنوان داده آزمون است. سائز تصاویر  $28 \times 28$  پیکسل است و تصاویر همانطور که در شکل ۳-۳ می‌بینید، به صورت سیاه و سفید هستند به عبارت دیگر ابعاد دیتاست آموزش به صورت (60000, 28, 28)



شکل ۶-۳: نحوه اتصال نورون‌ها در شبکه‌عصبی با یک لایه مخفی

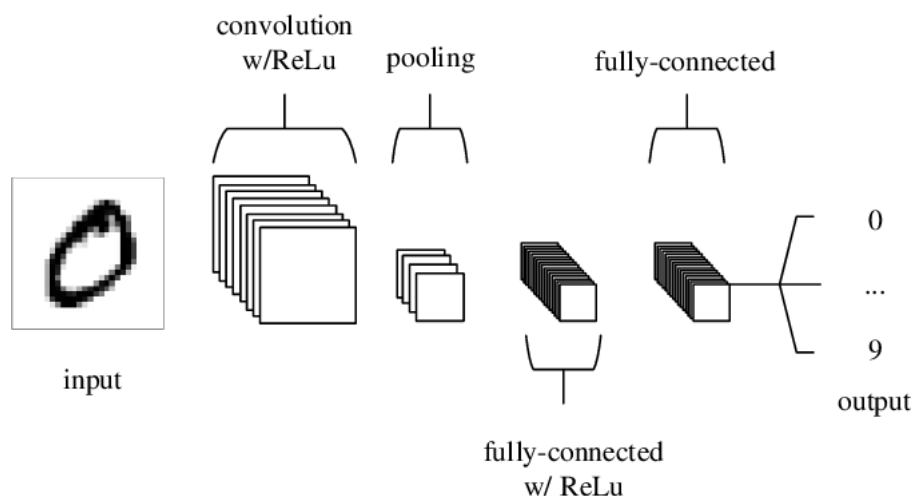
می‌باشد. مقادیر درایه‌های ماتریس عکس‌های اعداد بین ۰ تا ۲۵۵ که نشان‌دهنده شدت روشنایی برای هر نقطه می‌باشد. برای عملکرد بهتر و سریعتر همگراشدن مدل به هنگام بهینه‌سازی، عمل نرمال‌سازی صورت می‌گیرد تا همه مقادیر در بازه بسته ۰ تا ۱ قرار بگیرند. در ادامه باید معماری شبکه‌عصبی مورد نظر را مشخص کنیم. در این تمرین، صرفاً از یک شبکه‌عصبی با یک لایه مخفی

استفاده می‌کنیم. برای این منظور ابتدا ماتریس  $28 \times 28$  را به یک بردار  $784 = 28 \times 28$  تبدیل می‌کنیم. این ۷۸۴ تا نورون ورودی را به تمام نورون‌های لایه مخفی وصل می‌کنیم.

در آخر یک لایه خروجی با ۱۰ نورون با تابع softmax برای پیش‌بینی قرار می‌دهیم. همچنین تابع غیرخطی ای که استفاده می‌کنیم ReLu و برای جلوگیری از بیش‌برازش از dropout با احتمال ۰.۲ استفاده می‌کنیم تا در هر مرحله از آموزش ۰.۲ از نورون‌های لایه مخفی به صورت تصادفی را خاموش کند. نتایج ارزیابی و دقت مدل آموزش دیده را می‌توانید در کدی که برای این قسمت زده شده در ریپازیتوری<sup>۵۱</sup> گیت‌هاب بنده مشاهده کنید.

### ۳-۵ هفته پنجم – CNN ها و RNN ها

در این هفته آموزش‌های لازم در رابطه با شبکه‌های عصبی پیچشی<sup>۵۲</sup> و شبکه‌های عصبی



شکل ۳-۷: معماری یک شبکه عصبی پیچشی برای آموزش روی دیتاست MNIST

همروند<sup>۵۳</sup> داده شد و در انتها یک تمرین عملی از پیاده‌سازی یک CNN روی دیتاست MNIST و مقایسه نتایج آن با نتایج شبکه عصبی ساده (Fully connected) انجام شد. دقتی که در تمرین هفته قبل برای شبکه‌های fully connected با لایه‌های fully connected گرفته شد حدود ۹۴ درصد روی داده آزمون بود اما برای شبکه عصبی پیچشی ۹۸ درصد. کد این تمرین نیز در ریپازیتوری<sup>۵۴</sup> گیت‌هاب بنده قابل مشاهده می‌باشد.

<sup>۵۱</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-4>

<sup>۵۲</sup> Convolutional Neural Networks

<sup>۵۳</sup> Recurrent Neural Networks

<sup>۵۴</sup> <https://github.com/mohammadhashemii/ML-RahnemaCollege/tree/master/Week-5>

## ۳-۶ هفته ششم – آشنایی با کلان داده‌ها

در این هفته، مفاهیم پایه و اساسی برای کلان داده‌ها<sup>۵۵</sup> از قبیل Hadoop, Map reduced و spark تدریس شده‌اند. به علت نزدیک نبودن این مبحث به مفاهیم یادگیری ماشین در این گزارش خلاصه‌ای از آن نوشته نشده است.

## ۳-۷ هفته هفتم – ارائه فاز اول پروژه

از این هفته به بعد تمام تمرکز کارآموزان روی انجام پروژه پایانی این دوره بود. صورت این پروژه با همکاری شرکت داخلی سنجاق<sup>۵۶</sup> که یک بازار آنلاین خدمات به حساب می‌آید تعریف شده است. صورت کامل این پروژه در [اینجا](#) آمده است. اما اگر بخواهیم مختصر آنرا شرح دهیم، هدف پروژه راه اندازی یک سیستم آفلاین برای تشخیص ناهنجاری و خزنده بر بستر وب<sup>۵۷</sup> به کمک الگوریتم‌های غیرنظارتی است. دیتایی در فاز آموزش مدل در اختیار داشتیم، لاگ سرور سایت سنجاق است. این دیتا در طول مدت زمان یک هفته ای جمع آوری شده بود و نکته مهمی که این پروژه را چالشی می‌کرد، عدم وجود برچسب برای هر رکورد از دیتاست بود. در این فاز، پیش پردازش های لازم روی دیتای خام و فرآیند EDA بصورت متمرکز انجام شد. قبل از بررسی، پیشنهاد می‌شود فایل ارائه فاز اول را در اینجا مطالعه شود. همچنین تمامی کدها و منابع پروژه در ریپازیتوری<sup>۵۸</sup> گیت‌هاب بنده قابل مشاهده است.

## ۳-۷-۱ تحلیل اکتشافی داده (EDA)

دیتاستی که در این پروژه استفاده شده متشکل از لاگ NGINX سرور وبسایت سنجاق می‌باشد. تعداد رکوردهای این دیتاست، ۱۲۶۰۰۳۳ است. نمونه ای از این دیتاست در [اینجا](#) قابل مشاهده است. ویژگی‌های هر رکورد از این دیتاست عبارتند از:

ویژگی	جنس	توضیحات
ip	string	به عنوان نمونه: 207.213.193.143
time	string	به عنوان نمونه: 2021-05-12 05:06:00+04:30
method	string	Get/Put/Post/Options/Head
status_code	integer	وضعیت یک درخواست را نشان می‌دهد

<sup>55</sup> Big data

<sup>56</sup> <https://sanjagh.pro/tehran/>

<sup>57</sup> Web crawler detection

<sup>58</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection>

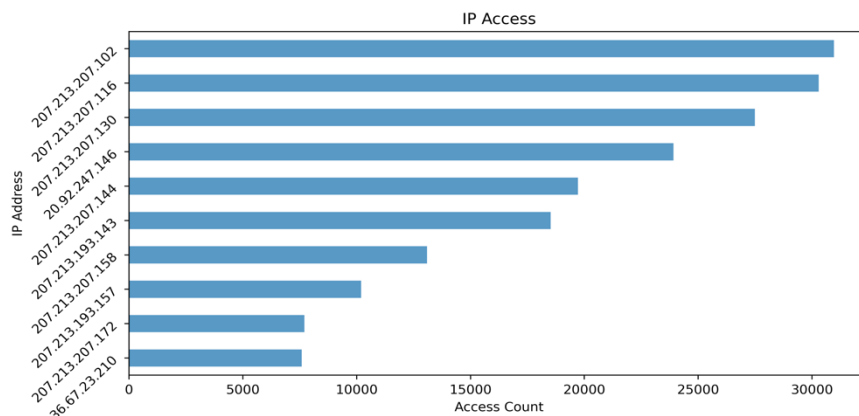
cdn/articles/1148001967 به عنوان نمونه:	string	path
حجم اطلاعات موجود در آن درخواست.	integer	response_length
Googlebot-Image/1.0 به عنوان نمونه:	string	user_agent
زمان پاسخ سرور به کلاینت	float	response_time

جدول ۳-۳: ویژگی‌های دیتاست استفاده شده

تحلیل اکتشافی داده یا همان EDA برای هر پروژه‌ای که با داده سروکار دارد الزامی است. EDA راهی برای مصورسازی، خلاصه‌سازی و تفسیرسازی اطلاعات مخفی و آشکار یک مجموعه دادگان است. EDA قدمی اساسی در علم داده است که با انجام آن یک دید خوبی از ویژگی‌های آماری یک دیتاست کسب می‌کنیم. با اتمام این مرحله می‌توان از ویژگی‌های بدست آمده در یک روش نظارتی یا غیرنظارتی بهره برد. قبل از بررسی این مرحله توصیه می‌شود فایل نوت‌بوک نوشته شده در ریپازیتوری<sup>۵۹</sup> گیت‌هاب بنده مطالعه شود.

ابتدا چند نمودار و شکل را رسم می‌کنیم تا شهود خوبی از دیتا کسب کنیم. به مرور می‌توان از روی همین نمودارها برخی از ویژگی‌های درخواست‌های جعلی<sup>۶۰</sup> را فهمید.

۱. بیشترین IP های مشاهده شده:



نمودار ۳-۲: تعداد درخواست برای هر IP.

سپس رکوردهای دیتاست را براساس بیشترین درخواست از آی‌پی‌ها مرتب می‌کنیم:

<sup>59</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection/tree/master/notebooks>

<sup>60</sup> Anomaly requests

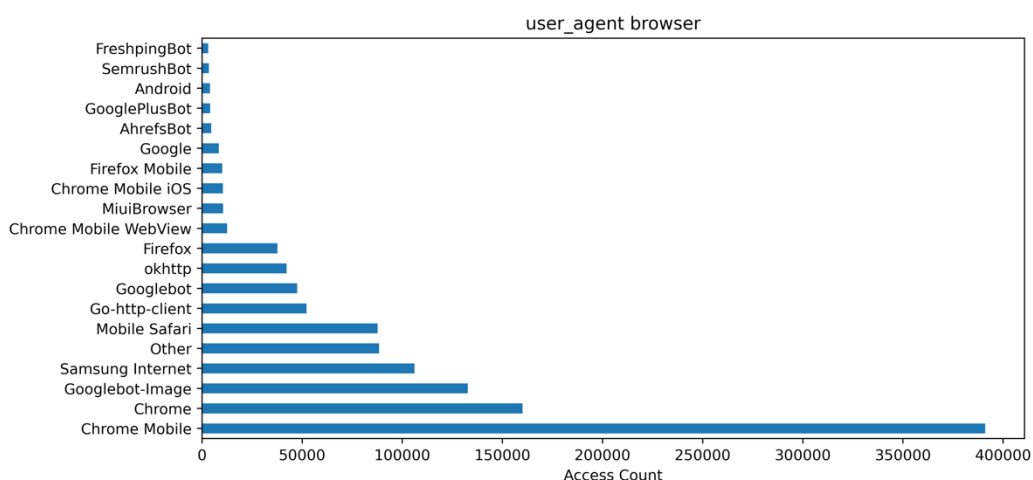


	ip	time	method	status_code	path	response_length	user_agent	response_time
0	207.213.207.102	2021-05-12 06:25:56+04:30	Get	304	cdn/articles/1148001967	0	Googlebot-Image/1.0	16.0
30968	207.213.207.116	2021-05-12 07:40:46+04:30	Get	304	cdn/profiles/1074674108	0	Googlebot-Image/1.0	16.0
61258	207.213.207.130	2021-05-12 09:25:34+04:30	Get	200	amp/price/1313296747	125767	Mozilla/5.0 (Linux; Android 6.0.1; Nexus 5X Bu...	28.0
88748	20.92.247.146	2021-05-12 05:08:46+04:30	Get	200	js/profession.c67de06df71c34fc126d.js	107970	sentry/21.4.1 (https://sentry.io)	16.0
112660	207.213.207.144	2021-05-12 09:02:56+04:30	Get	200	amp/price/252451961	121639	Mozilla/5.0 (Linux; Android 6.0.1; Nexus 5X Bu...	28.0
132378	207.213.193.143	2021-05-12 05:06:00+04:30	Get	304	cdn/profiles/1026106239	0	Googlebot-Image/1.0	32.0
150901	207.213.207.158	2021-05-12 09:25:55+04:30	Get	304	cdn/articles/2121333045	0	Googlebot-Image/1.0	16.0
163992	207.213.193.157	2021-05-12 05:07:04+04:30	Get	200	amp/blog/article/1197238235	101087	Mozilla/5.0 (Linux; Android 6.0.1; Nexus 5X Bu...	144.0
174189	207.213.207.172	2021-05-12 09:28:50+04:30	Get	304	cdn/pro_photo_gallery/1540103160	0	Googlebot-Image/1.0	28.0
181890	36.67.23.210	2021-05-12 05:06:00+04:30	Head	200	877499224	0	Go-http-client/2.0	28.0

مشاهده می‌شود که اتفاقاً IP هایی که بیشترین درخواست‌ها از آنها آمده محتمل به جعلی بودن و یا برخاسته از یک خزنده باشند هست. زیرا با نگاه کردن به user\_agent های آنها شاهد نمونه هایی از بات ها مانند GoogleBot و ... هستیم.

۲. بیشترین مرورگرهای دیده شده: برای اینکه بفهمیم مرورگر کاربر چه بوده است از کتابخانه پایتونی به نام <sup>61</sup>user\_agents استفاده کردیم. با استفاده از این کتابخانه می‌توان مرورگر، اینکه درخواست از یک بات بوده یا نه، اینکه درخواست از یک کامپیوتر یا موبایل بوده و ... از روی user agent درخواست بدست‌آوریم.

در نمودار ۲-۳ می‌بینیم که برخی از مرورگرهایی که یک بات هستند ترافیک زیادی برای سایت ایجاد کرده‌اند. به علت محدود بودن حجم گزارش، برای مطالعه دیگر بررسی‌ها و نمودارهایی که



نمودار ۳-۳: تعداد درخواست برای مرورگر هر user agent

<sup>61</sup> <https://pypi.org/project/user-agents/>

به شهود ما از دیتاست و درخواست‌های جعلی کمک می‌کند، به [اسلایدهای ارائه فاز اول](#) رجوع شود.

## ۳-۷-۲ مهندسی ویژگی<sup>۶۲</sup>

تا این قسمت شهود خوبی از ویژگی‌ها و عوامل تاثیر گذار بر جعلی بودن یک درخواست از سمت کلاینت گرفته‌ایم. با توجه به دشوار بودن تشخیص anomaly بودن یا نبودن از روی تک درخواست، تصمیم بر این گرفتیم تا دیتاست را براساس IP و User agent گروه‌بندی کنیم. با این کار گویی داریم یک نشست<sup>۶۳</sup> تعریف می‌کنیم. پس در این مرحله سعی بر تعریف و مهندسی‌سازی ویژگی برای هر نشست داریم. تعریف ویژگی‌ها از یک مقاله<sup>۶۴</sup> پژوهشی الهام گرفته شده است و در جدول ۳-۳ لیست کامل آنرا می‌توانید مشاهده کنید.

ویژگی برای نشست	توجیه و دلیل استفاده از این ویژگی
تعداد درخواست‌ها	زیادبودن تعداد درخواست در یک نشست موجب افزایش احتمال خزنده یا بات بودن آن کاربر می‌شود.
انحراف معیار عمق درخواست‌ها <sup>۶۵</sup>	کاربرهای انسان، درخواست‌هایی که معمولا در یک نشست می‌زنند دارای عمق path با طول‌های متفاوتتر نسبت به خزنده‌ها و بات‌ها است.
درصد درخواست‌های از خانواده ۴۰۰	معمولا خزنده‌ها به پاسخ‌هایی از سمت سرور با خطایی از خانواده ۴۰۰ بیشتر مواجه می‌شوند.
درصد درخواست‌های از خانواده ۳۰۰	خانواده ۳۰۰ به معنای redirect شدن به یک صفحه دیگر است. خزنده‌ها و بات‌ها معمولا بیشتر با این پاسخ روبرو می‌شوند.
درصد درخواست‌های با متد HEAD	در خزنده‌ها یا بات‌ها بیشتر است زیرا احتمال مواجهه با صفحات پاک‌شده یا تاریخ گذشته در این نوع از کاربران زیادتر است.
درصد درخواست‌های عکس	خزنده‌ها و بات‌ها معمولا به عکس‌ها درخواست نمی‌دهند.
جمع و میانگین response_length و response_time	بدلیل اینکه کاربرهای انسان از مرورگر برای دسترسی به صفحات وب استفاده می‌کنند، وقتی به یک صفحه درخواست می‌زنند، نشست مجبور به دریافت منابع و عکس‌های متفاوتی است و همین باعث زیاد شدن زمان پاسخ و حجم درخواست می‌شود.
ویژگی‌های استخراج شده از user_agent	مرورگر-سیستم‌عامل-بات بودن یا نبودن-از یک PC بودن یا نبودن.

<sup>62</sup> Feature engineering

<sup>63</sup> session

<sup>64</sup> [https://link.springer.com/chapter/10.1007/978-3-642-21916-0\\_52](https://link.springer.com/chapter/10.1007/978-3-642-21916-0_52)

<sup>65</sup> STD of path's depth

تعداد درخواست‌های به robtots.txt	خزنده‌ها برای اینکه متوجه شوند سرور سایت امکان خزش به چه صفحاتی را داده است، به robots.txt درخواست می‌زنند. اینگونه اطلاعات در آنجا نوشته شده است.
-------------------------------------	--

جدول ۳-۴: ویژگی‌های مهندسی شده و استخراج شده از دیتاست برای هر تشست

### ۸-۳ هفته هشتم – توسعه پروژه نهایی

در این هفته از ویژگی‌های تعریف شده در قسمت قبلی استفاده می‌کنیم تا مدل‌های پایه<sup>۶۶</sup> را آموزش دهیم. به دلیل نداشتن برچسب برای دادگان، باید از مدل‌های غیرنظارتی بهره می‌بریم. مدل‌های غیرنظارتی زیادی در علم یادگیری ماشین برای تشخیص داده پرت<sup>۶۷</sup> تا به الان شناخته شده‌اند. از معروفترین آنها می‌توان به موارد زیر اشاره نمود:

- Isolation Forest
- Local Outlier Factor
- One-class SVM
- Robust covariance

خوشبختانه کتابخانه `Scikit-learn`<sup>۶۸</sup> این الگوریتم‌ها را پشتیبانی و پیاده‌سازی کرده است و در این پروژه از آن استفاده کردیم. در این هفته دو مدل پیاده‌سازی شده است. کدها و نوت‌بوک‌های مربوط به این قسمت نیز در ریپازیتوری<sup>۶۹</sup> گیت‌هاب بنده قابل مشاهده هستند.

#### ۳-۸-۱ پیاده‌سازی مدل Isolation Forest

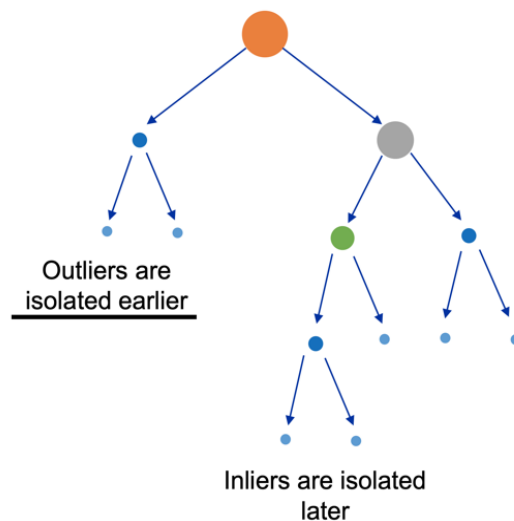
این روش، تکنیکی برای تشخیص داده‌های پرت به روش غیرنظارتی است. در این روش مشاهده‌ها، با انتخاب ویژگی‌ها به صورت تصادفی و جداسازی مقدار آن به بیشترین و کمترین مقادیر ویژگی انتخابی، ایزوله می‌شوند. به دلیل خاصیت بازگشتی بودن این روش، این روش با یک ساختار درختی قابل نمایش است. بدلیل اینکه مقادیر ویژگی‌های داده‌های پرت به طرز قابل توجهی با بقیه دادگان تفاوت دارد، داده‌های پرت زودتر در درخت تصمیم ایزوله می‌شوند (شکل ۱-۳). به عبارتی با تنظیم کردن یک آستانه که خود یک ابرپارامتر برای تعداد جداسازی از بالا تا پایین (برگ) درخت می‌توان داده‌های پرت را شناسایی کرد. پیاده‌سازی ای که در `Scikit-learn` برای این روش شده است، یک امتیاز بین ۱- و ۱ با استفاده از تعداد جداسازی‌ها به هر نمونه از دیتا می‌دهد. هرچه امتیاز به ۱- نزدیک تر باشد، به معنای پرت بودن (در اینجا خزنده یا بات بودن کاربر) می‌باشد.

<sup>۶۶</sup> Baseline

<sup>۶۷</sup> Outlier detection

<sup>۶۸</sup> [https://scikit-learn.org/stable/modules/outlier\\_detection.html](https://scikit-learn.org/stable/modules/outlier_detection.html)

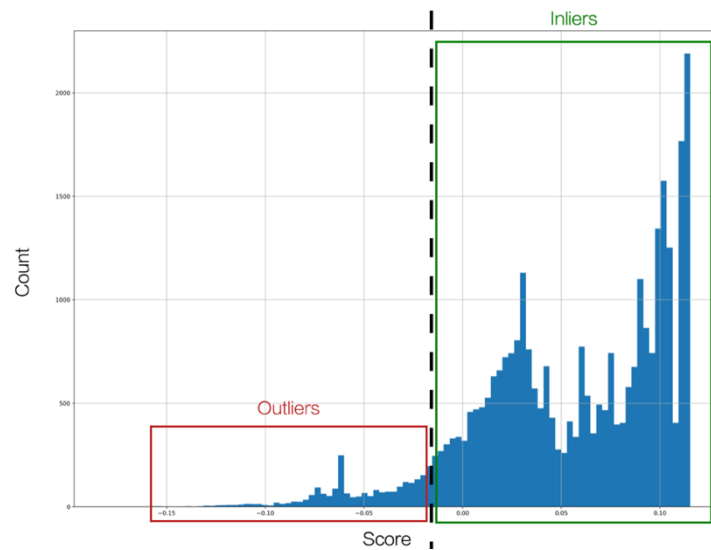
<sup>۶۹</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection/tree/master/notebooks>



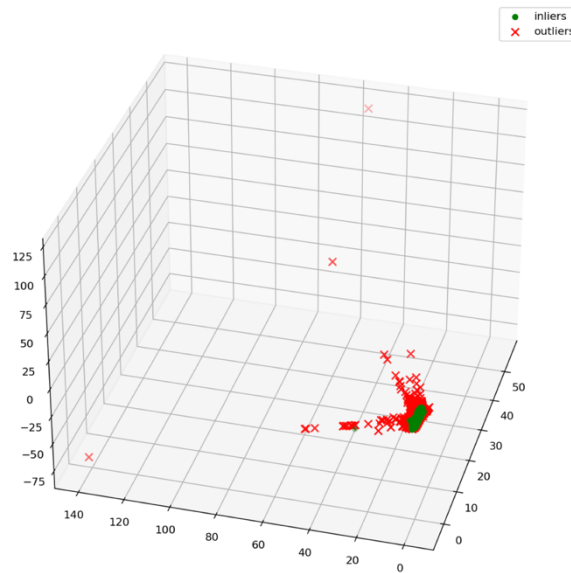
شکل ۸-۳: داده‌های پرت زودتر به برگ درخت می‌رسند.

همانطور که گفته شده ابرپارامتر آستانه باید توسط ما تعیین شود و با توجه به نداشتن برچسب برای دادگان، امکان fine-tune کردن این پارامتر وجود ندارد. پس در این مرحله فعلا بصورت تجربی و آزمون و خطا این آستانه تعیین شد. در هفته بعد، ارزیابی دقیقتری به عمل خواهد آمد. در نمودار ۳-۳ هیستوگرام امتیازهای نمونه‌ها مشاهده می‌شود. در این حالت، نمونه‌های با امتیاز کمتر از ۰ به عنوان داده پرت در نظر گرفته شده‌اند.

با این آستانه گذاری از ۳۱۵۴۱ نشستی که با گروه‌بندی به کمک IP و User agent بدست آمده بود، ۶۱۵ تای آنها داده پرت تشخیص داده شدند. برای اینکه بفهمیم چقدر این مدل درست آموزش دیده، به کمک کتابخانه user\_agent و ویژگی is\_bot بودن آن دیدیم که ۶۰۳ تا از ۶۱۵ تایی که داده پرت تشخیص داده شدند، بات هستند. این به این معنا بود که نسبتا مدل به عنوان یک مدل پایه عملکرد خوبی داشته است. اما هنوز به صورت قطعی بدلیل عدم وجود برچسب برای دادگان نمی‌توان دقت برای روش و مدل به کار گرفته شده مشخص نمود.



نمودار ۳-۴: تفکیک داده‌های پرت و غیرپرت با آستانه گذاری بصورت تجربی



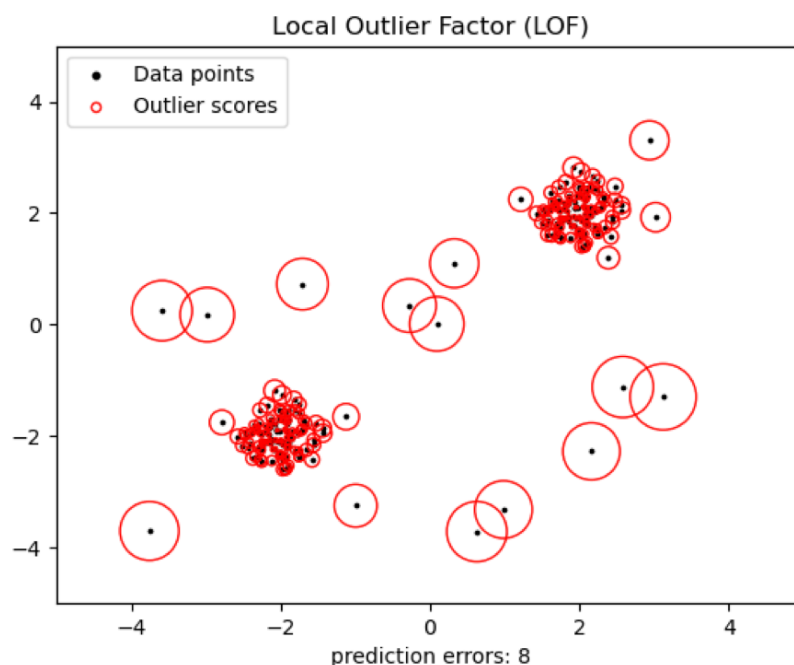
نمودار ۳-۵: مصورسازی دادگان پرت و نرمال به کمک PCA

## ۳-۸-۲ پیاده‌سازی مدل Local Outlier Factor

در روش LOF، انحراف محلی چگالی ۷۰ هر نمونه نسبت به همسایگان خود محاسبه می‌شود. پرت بودن یک داده به میزان ایزوله و تنها بودن آن نمونه نسبت به همسایگان خود سنجیده می‌شود. همچنین چگالی محلی به کمک متریک فاصله‌ای که در روش KNN محاسبه می‌شود صورت

<sup>70</sup> Local deviation of density

می‌گیرد. در نمودار ۵-۳ نیز مشاهده می‌کنید که هرچه یک نمونه همسایگان کمتری داشته باشد،



نمودار ۶-۳: نحوه عملکرد الگوریتم LOF

احتمال داده پرت بودن آن نیز بیشتر است.

اما این روش نتایج خوبی برای دیتاستی که ما داشتیم به همراه نداشت زیرا حتی محتمل‌ترین نمونه‌هایی که به عنوان داده پرت در نظر گرفته بود گاهی داده نرمال بودند. به همین دلیل از این روش در ادامه استفاده‌ای نشد.

### ۳-۹ هفته نهم – توسعه پروژه نهایی (ادامه)

در این هفته، روش‌ها و الگوریتم‌های دیگری نظیر اتوانکودرها<sup>۷۱</sup> و الگوریتم‌های مبتنی بر پردازش زبان‌های طبیعی<sup>۷۲</sup> پیاده‌سازی و بررسی شدند. همچنین ارزیابی بهتر و نسبتاً دقیق‌تری روی خروجی مدل‌ها انجام شد تا از صحت درستی عملکرد الگوریتم‌ها مطمئن شویم. کدها و توضیحات این هفته نیز در ریپازیتوری<sup>۷۳</sup> گیت‌هاب بنده قابل مشاهده است.

#### ۳-۹-۱ پیاده‌سازی الگوریتم Auto-encoder

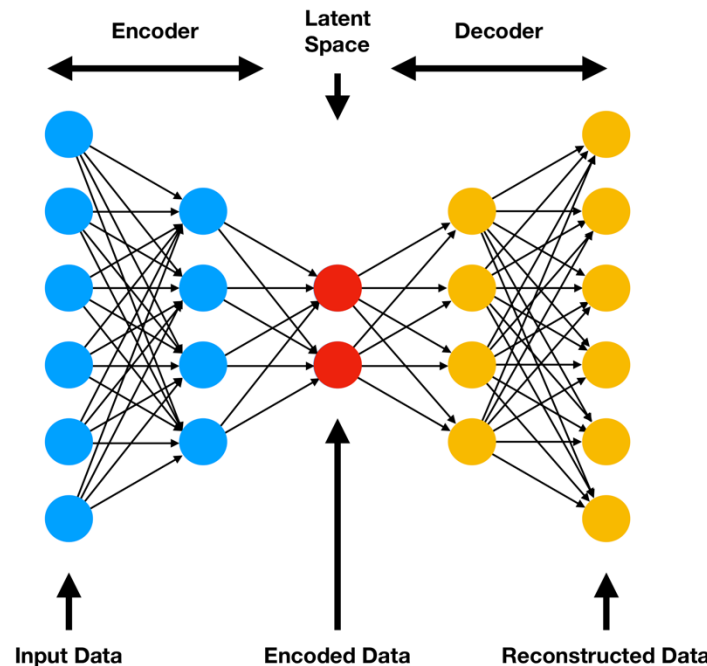
اتوانکودرها یک نوع خاصی از شبکه‌های عصبی هستند که مقادیر نورون‌های ورودی را در خروجی کپی می‌کنند. در فرآیند آموزش این نوع شبکه، دیگر نیازی به برچسب داده‌ها وجود ندارد برای همین می‌توان اتوانکودرها را جزو الگوریتم‌های غیر نظارتی دانست. به طور معمول لایه‌های

<sup>۷۱</sup> Auto-encoder

<sup>۷۲</sup> Natural Language Processing (NLP)

<sup>۷۳</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection/tree/master/notebooks>

مخفی در این نوع شبکه‌ها، دارای تعداد نورون کمتری نسبت به نورون‌های ورودی و خروجی دارند. به همین دلیل، لایه‌های مخفی اطلاعات ضروری را در خود ذخیره می‌کنند و از نویزها صرف نظر می‌کنند. در فرآیند آموزش این نوع شبکه دو قسمت مهم **encode** و **decode** وجود دارد. در مرحله‌ای **encoding**، مقادیر ورودی را فشرده‌سازی می‌کند و به فضای برداری لایه مخفی می‌برد. در مرحله



شکل ۹-۳: قسمت‌های اصلی *auto-encoder* ها

**decoding**، اطلاعات فشرده‌شده، بازسازی<sup>۷۴</sup> می‌شوند. در شکل ۲-۳ نیز این مراحل مصور شده است.

اتوانکودرها کاربردهای وسیعی در زمینه پردازش تصویر و بینایی ماشین دارند همچنین نتایج درخشانی در زمینه تشخیص ناهنجاری نیز داشته‌لند. در فرآیند **decoding**، هنگامی که عمل بازسازی صورت می‌گیرد می‌توان خطای یکسان نبودن داده خروجی با داده ورودی اولیه را حساب کرد. به عبارت دیگر، لایه مخفی سعی بر یادگیری یک **embedding** از داده‌های ورودی است و می‌خواهد داده‌های ورودی را تنها با داشتن ویژگی‌های موجود در لایه مخفی بازسازی کند. به طبع اگر داده پرتی که اختلاف زیادی از نظر مقادیر ویژگی‌ها با داده‌های دیگر وجود داشته باشد، خطای بازسازی بسیار زیاد و متفاوت است. پس با این رویکرد می‌توان از اتوانکودرها برای تشخیص ناهنجاری نیز استفاده نمود. در فرآیند آموزش از بهینه‌ساز<sup>۷۵</sup> **Adam** و تابع خطای **MSE** برای محاسبه خطا استفاده کردیم. همچنین از **ReLU** به عنوان تابع غیرخطی‌ساز استفاده کرده‌ایم. دیتاست با توزیع

<sup>74</sup> Reconstruction

<sup>75</sup> Optimizer

۸۰-۲۰ به دو داده آموزش و تست تقسیم شد و در جدول مقادیر خطا پس از آموزش برای شبکه‌های با معماری‌هایی متفاوت قابل مشاهده است.

تعداد نورون‌ها	خطای داده آموزش	خطای داده تست
[15, 7, 15]	0.42	0.48
[15, 3, 15]	0.28	0.39
[15, 7, 3, 7, 15]	0.29	0.43
[15, 7, 7, 7, 15]	0.31	0.42

جدول ۵-۳: مقادیر خطای *auto-encoder* با معماری‌های گوناگون

پس از آموزش مدل، برای فاز پیش‌بینی ناهنجاری بودن یا نبودن، مشابه الگوریتم‌های قبلی، باید یک حد آستانه برای خطای MSE نیز تعریف کرد. با توجه به دانش قبلی و بررسی بات‌بودن یا نبودن نشست‌های داخل دیتافریم براساس *user agent* شان، حدود ۵ درصد از دیتاست به واضح بات بودند، بنابراین آستانه خطای MSE به طوری انتخاب شد که ۵ درصد از نمونه‌ها به عنوان ناهنجاری تشخیص داده شوند.

## ۳-۹-۲ ارزیابی مدل نهایی

اتوانکودر با توجه به اینکه عملکرد بهتری روی دیتاست داشت به عنوان مدل نهایی انتخاب شد و تصمیم بر این شد تا این مدل به فاز *production* برود. قبل از آن با اینکه برچسب دیتاها وجود نداشت، ارزیابی نسبتاً دقیقی توانستیم ارائه دهیم.

### ۳-۹-۲-۱ وجود برچسب برای بعضی از نمونه‌ها

شرکت سنجاق، برای اینکه از عملکرد مدل‌های یادگیری ماشین کارآموزان مطلع شود، از روی عمد با استفاده از دو خزنده به وبسایت سنجاق درخواست زده بود و تیم فنی سنجاق توقع داشت مدل‌های یادگیری ما توانایی تشخیص این دو خزنده را داشته باشند که خوشبختانه مدل اتوانکودر به خوبی این دو نشست را جعلی یا همان خزنده پیش‌بینی کرد. (توضیحات کامل در [اسلایدها](#))



## ۳-۹-۲-۲ انتظارات ما از مدل

در هفته هفتم، در فاز EDA چندین ویژگی آماری تعریف و مهندسی شدند. به طور مثال یکی از آنها تعداد درخواست‌های نشست بود و ما توقع داشتیم که برای نشست‌های جعلی، تعداد درخواست‌ها به طور قابل توجهی بیشتر باشد. پس از پیش‌بینی مدل اتوانکودر، میانگین تعداد درخواست‌ها در نشست‌های نرمال و همچنین در نشست‌های جعلی را مقایسه نمودیم و دقیقاً همان چیزی که انتظار می‌رفت رخ داد. در جدول زیر این مقایسه‌ها قابل مشاهده هستند که فقط دو تا از ویژگی‌های تعریف شده نتیجه مطلوبی از خود نشان ندادند:

Average of	# of requests	Path length STD	Percentage of 4xx	Percentage of 3xx	Percentage of HEAD reqes	consecutive repeated requests	robots.txt requests	Percentage of image requests
<b>Outliers</b>	231	0.43	3.21%	9.33%	0.34%	0.81	0.08	9.74%
<b>Inliers</b>	25	0.39	0.68%	26%	0.00003	0.62	0.0	28.16%

جدول ۳-۶: میانگین مقادیر ویژگی‌های به برای هر ویژگی برای هر نشست و مقایسه آنها

## ۳-۹-۲-۳ گزارش دسته‌بندی<sup>۷۶</sup>

برای اینکه ارزیابی دیگری داشته باشیم، ۱۵۰ تا از مطمئن‌ترین پیش‌بینی‌های مدل را به صورت دستی برچسب زدیم و ارزیابی‌های کلاسیک مدل‌های نظارتی را برای آنها اعمال نمودیم:

Accuracy	90%
Precision	85.71%
Recall	100%
F1-score	92.30%

جدول ۳-۷: گزارش دسته‌بندی برای مطمئن‌ترین پیش‌بینی‌ها

## ۳-۱۰ هفته دهم – ارائه فاز دوم پروژه نهایی

در هفته آخر مدل نهایی به فاز production رسید. با استفاده از تکنولوژی Flask برای

<sup>76</sup> Classification report

بک‌اند و ReactJs برای فرانت‌اند، یک صفحه‌ی وب طراحی کردیم تا کاربر بتواند با آپلود کردن یک فایل *output.log* از لاگ سرور خود (مثال)، ناهنجاری‌ها (خزنده‌ها و بات‌ها) را به کمک مدل هوش مصنوعی شناسایی کند. در آخر [اسلایدهای](#) ارائه فاز نهایی تهیه و تنظیم شد. کدهای مربوط به API در ریپازیتوری<sup>۷۷</sup> گیت‌هاب بنده نیز قابل مشاهده است. همچنین برای دسترسی به صفحه دمو طراحی از این [لینک](#) استفاده کنید. (البته ممکن است گاهی اوقات سرور بالا نباشد و آنوقت برای اجرا گرفتن دمو روی سیستم محلی خود، طبق توضیحات [اینجا](#) عمل کنید).

**Web Crawler Detection**

This is a demo of the [Rahnema College](#) Machine Learning internship final project. You can upload a log file (e.g. *output.log*) which contains the server logs in order to detect the crawlers. [\[Source Code\]](#)

A sample log is shown here. Therefore, the logs in the file you upload must be the same as the sample below: [\[Help\]](#)

```
207.213.193.143 [2021-5-12T5:6:0.0+0430] [Get /cdn/profiles/1026106239] 304 0 [[Googlebot-Image/1.0]] 32
```

Upload

IP	User agent	MSE score	Is crawler?
113.97.153.110	Mozilla/5.0 (Linux; Android 6.0.1; SM-A500H) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/89.0.4389.105 Mobile Safari/537.36	0.233	0
207.213.207.116	Googlebot-Image/1.0	0.433	1
35.109.154.235	Mozilla/5.0 (Linux; Android 5.1; HUAWEI CUN-L21) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.101 Mobile Safari/537.36	0.964	1
35.154.147.116	Mozilla/5.0 (Linux; U; Android 9; fa-ir; Redmi Note 8 Build/PKQ1.190616.001) AppleWebKit/537.36 (KHTML, like Gecko) Version/4.0 Chrome/79.0.3945.147 Mobile Safari/537.36 XiaoMi/MiuiBrowser/12.10.5-go	0.135	0

شکل ۱۰-۳: تصویری از محیط وب‌اپلیکیشن طراحی شده

<sup>77</sup> <https://github.com/mohammadhashemii/Web-Crawler-Detection/tree/master/App>

## ۴ خلاصه

هفته اول تا ششم این دوره کارآموزی اغلب به آموزش و فراگرفتن مطالب مهم در یادگیری ماشین سپری شد. مباحث مهم جبرخطی و آماراحتمال در ریاضیات از اساسی ترین مقدمات و قسمت‌هایی است که هر مهندس یادگیری ماشین یا دانشمند داده نیاز دارد تا بتواند بهترین مدل‌ها و الگوریتم‌ها را برای یک تسک هوش مصنوعی انتخاب و پیاده‌سازی کند. اینکه یک مهندس توانایی تفکیک بین الگوریتم‌های یادگیری ماشین متعددی که وجود دارند قائل شود داشته باشد، از شروط قبولی در تمام آزمون‌های استخدامی برای مشاغلی از قبیل دانشمند داده، تحلیلگر داده و مهندس یادگیری ماشین می‌باشد. بنابراین می‌بایستی با تمامی الگوریتم‌های این زمینه به اندازه کافی اخت گرفت که تمرین‌های عملی و تدریس‌هایی که در این دوره کارآموزش صورت گرفت، یک قدم ما را به این هدف نزدیکتر کرد. همچنین تسلط به الگوریتم‌های یادگیری عمیق با توجه به حجم بودن دیتای روزمره و پیچیده بودن توابعی که قرار است کار تخمین زدن را برای ما انجام بدهند لازم و ضروری است.

در هفته‌های هفتم الی دهم، بهترین و معروفترین الگوریتم‌های غیر نظارتی که در سال‌های اخیر نتایج بسزایی در تمامی زمینه‌ها از خود نشان دادند مورد مطالعه و پژوهش قرار گرفتند. توانستیم یک سیستم که قابلیت تشخیص جعلی بودن یا نبودن درخواست‌هایی که توسط کاربران به یک وبسایت زده می‌شود را با دقت خوبی داشته باشد، پیاده‌سازی کنیم. در این پروژه از اتوانکودرها به عنوان یک شبکه‌عصبی با یک رویکرد غیرنظارتی در عین حال دقت بالا، استفاده نمودیم.

### نظرات کارآموز جهت بهبود کارآموزی:

شرکت رهنماکالج به قطع یکی از بهترین مراکز برای آمادگی کارآموز برای ورود به صنعت هوش مصنوعی در ایران است. تمامی اساتید و منتورهای این دوره، از مجرب‌ترین‌ها در شرکت‌های داخلی معتبر می‌باشند. تنها نکته ای که می‌توانم پیشنهاد کنم به عنوان کسی که این دوره کارآموزی را گذرانده این است مباحث تئوری ریاضیات یعنی جبرخطی و آماروا احتمال این دوره را به شدت جدی گرفته، زیرا کمتر جایی این مباحث به صورت حرفه‌ای تدریس می‌شود.