



دانشگاه صنعتی شریف

دانشکده‌ی مهندسی کامپیووتر

## یادگیری ماشین

پاییز ۱۴۰۴

استاد: دکتر علی شریفی زارچی

مسئول پژوهش: امیر رضا آذری

مهلت ارسال نهایی: ۳۰ دی

پژوهش

- پژوهش، ۱+۳ نمره درس می‌باشد.

ابتدا، با دقت کامل، تمامی پژوهه‌ها را مطالعه بفرمایید، سپس در فرم قرار داده شده، اولویت خود را از یک تا نه مشخص بفرمایید تا یک پژوهه به شما تعلق بگیرد. توجه بفرمایید ظرفیت هر پژوهش، تنها ۷ تیم می‌باشد.

معیار تصمیم‌گیری نهایی برای هر تیم، زمان پرکردن فرم می‌باشد. بنابراین این کار را به زمان دیگر موکول نکنید زیرا در این صورت اولویت‌های پایانی به شما خواهد رسید و جای اعتراضی وجود نخواهد داشت.

یک نمره امتیازی پژوهش، شامل ۰.۵ نمره بخش رقابتی پژوهش و ۰.۵ نمره دیگر برای بخش خلاقیت می‌باشد.

**بخش رقابتی:** هر پژوهش، یک معیار رقابت دارد. تیم اول ۰.۵، تیم دوم ۰.۳ و تیم سوم ۰.۲ نمره اضافه بر روی نمره پژوهش و درس دریافت خواهد کرد. تیم‌های برتر بر اساس نظر مسئول مربوط به هر پژوهش انتخاب خواهند شد.

**بخش خلاقیت:** در این بخش هرگونه خلاقیت و ایده نو و جدید که قابل بحث و تاثیرگذار باشد، نمره خواهد گرفت. تشخیص میزان خلاقیت با مسئول پژوهش انتخاب شده می‌باشد.

هر پژوهش مسئول مشخص دارد که بعد از مشخص شدن پژوهه‌های تیم‌ها، معرفی خواهد شد. شما در حین انجام دادن پژوهش با ایشان در ارتباط خواهید بود و هرگونه ابهام و سوالی را می‌توانید با ایشان در میان بگذارید.

**گزارش نهایی:** گزارش نهایی هر پژوهه متفاوت می‌باشد. اما در تمامی گزارش‌ها باید میزان همکاری اثرباره هر شخص مشخص شده باشد. همچنین در صورت استفاده از ابزارهای LLM، بخشی که از ابزارها استفاده شده است باید مشخص شود و همچنین پرداخته شده ذکر شود.

هر تیم موظف است یک مخزن گیت‌هاب برای پژوهش خود ایجاد کند. نکته حیاتی این است که هر عضو تیم باید منحصر با اکانت شخصی خود فعالیت کند و تمام تغییرات و ارسال کد (Commit) خود را به صورت جزئی و مستمر در مخزن ثبت نماید. میزان همکاری اعضا در تاریخچه کامیت‌ها به دقت بررسی خواهد شد و انتظار می‌رود که مشارکت تمامی اعضای تیم در پیشبرد پژوهش عادلانه و متوازن باشد؛ چرا که این شفافیت، بخشی از فرآیند ارزیابی نهایی خواهد بود.

**توجه بفرمایید هرگونه عدم تطابق تحويل و گزارش تیم‌ها با نکات ذکر شده توسط مسئول پژوهش، سبب نمره منفی خواهد شد.**

دقت بفرمایید نمره هر پژوهش از ۱۰۰ می‌باشد، اما بعضی پژوهه‌ها بنا به سختی آن‌ها، نمرات امتیازی دارند که این نمره امتیازی، تنها بر روی خود نمره پژوهش اعمال خواهد شد. به این معنی که نمره ۱۱۰ با ۱۰۰ فرقی نخواهد داشت.

## پژوهش‌ها

### پژوهش ۱ (۱۰۵ نمره)

#### تعریف مسئله

#### سیگنال‌های الکتروانسفالوگرام (EEG)

سیگنال‌های EEG ثبت کننده فعالیت الکتریکی مغز هستند که از طریق الکترودهایی روی پوست سر اندازه‌گیری می‌شوند. این سیگنال‌ها اطلاعات ارزشمندی درباره وضعیت‌های ذهنی، فعالیت‌های شناختی و عملکرد بخش‌های مختلف مغز فراهم می‌کنند و به طور گسترده در علوم اعصاب، تشخیص‌های بالینی و واسطه‌های مغز-رایانه مورد استفاده قرار می‌گیرند.

در این پژوهش، با داده‌های EEG آشنا می‌شویم و روش‌های پیش‌پردازش، حذف نویز و تمیز سازی داده‌ها را بررسی می‌کنیم. سپس به استخراج ویژگی‌های مناسب از سیگنال‌ها پرداخته و داده‌ها را دسته بندی و خوشه‌بندی خواهیم کرد.

## تصویرسازی حرکتی (Motor Imagery)

تصویرسازی حرکتی فرآیندی است که در آن از فرد خواسته می‌شود بدون انجام یک حرکت، صرفاً انجام آن حرکت را در ذهن خود تصور کند؛ برای مثال تصور حرکت دست یا پا. در این پروژه، از داده‌های EEG ثبت شده از افراد سالم در حین انجام تصویرسازی حرکتی استفاده می‌کنیم. سیگنال‌های EEG با قرار دادن ۵۹ الکترود روی نواحی مختلف جمجمه اندازه‌گیری شده اند و تغییرات پتانسیل الکتریکی مغز را به صورت تابعی از زمان نشان می‌دهند. مجموعه داده‌های مورد استفاده، دارای سه کلاس «دست چپ»، «دست راست» و «پا» می‌باشد. در طی ثبت داده‌ها، از نشانه‌های بصری (مانند فلاش‌ها) برای مشخص کردن نوع حرکت مورد تصور استفاده شده و فرد مورد آزمایش، در بازه‌های زمانی مشخص، تصویرسازی حرکت نمایش داده شده را انجام داده است.

انتظار می‌رود هنگام تصور یک حرکت خاص، الگوهای متمایزی در فعالیت نواحی مختلف مغز ایجاد شود که بتوان آن‌ها را در سیگنال‌های EEG تشخیص داد. هدف این پروژه، پردازش سیگنال‌های EEG و پیش‌بینی نوع حرکت تصور شده در هر پنجره‌ی زمانی با استفاده از روش‌های یادگیری ماشین است.

### ساختار داده‌های ورودی

فایل دیتابست را از [این لینک](#) دانلود کنید. در ادامه از فایل BCICIV\_calib\_ds1a.mat استفاده می‌کنیم که فرآیند تصویرسازی حرکتی را توسط دو کلاس «دست چپ» و «پا» انجام داده است. برای درک ساختار داده‌های مورد استفاده به [این لینک](#) مراجعه کنید.

### ۱. پیش پردازش داده (۳۰ نمره)

#### ۱.۰۱ (۱۵ نمره) بارگذاری داده

فرآیند تصویرسازی حرکتی برای یک فرد به صورت پشت سر هم و مطابق توضیحات ارائه شده در لینک توضیحات صورت می‌گیرد. به منظور تفکیک بازه‌های زمانی متناسب با هر حرکت، نیاز است که سیگنال پیوسته‌ی ورودی را پنجره بندی کنیم. طول هر پنجره‌ی زمانی را با توجه به بازه‌ی نمایش نشانه‌های بصری و فرکانس نمونه برداری به دست آورده و از متغیر pos به عنوان نقطه‌ی شروع هر پنجره استفاده کنید.

کدی بنویسید که باخشهای مورد نیاز از دیتابست را خوانده و پردازش کند. خروجی این مرحله باید یک دیتابست در قالب معمول مسائل یادگیری ماشین باشد که فیچرها و برچسب‌های متناظر با هر نمونه را شامل می‌شود. ۷۵٪ داده‌ها را به آموزش و ۲۵٪ را به تست اختصاص دهید.

نمودار مربوط به یکی از نمونه‌های دیتابست به دست آمده را به ازای چنل‌های ۰، ۱۵، ۳۰، ۴۵ و ۵۹ به صورت تابعی از زمان رسم کنید.

#### ۲.۰۱ (۵ نمره) فیلترهای زمانی

سیگنال‌های مغزی بر اساس بازه‌های فرکانسی به چندین دسته تقسیم می‌شوند. در کاربرد فعلی، سیگنال‌های Mu، در بازه‌ی فرکانسی ۸ تا ۱۳ هرتز و سیگنال‌های Beta در بازه‌ی فرکانسی ۱۳ تا ۳۰ هرتز بیشترین کاربرد را دارند. با اعمال یک فیلتر band-pass، اطلاعات فرکانسی مربوط به این دو دسته را حفظ کرده و باقی فرکانس‌ها را حذف کنید.

#### ۲.۰۲ (۱۰ نمره) استخراج ویژگی‌ها

با توجه به نزدیکی مکانی الکترودها، سیگنال‌های چنل‌های مختلف، در سطح جمجمه با یکدیگر ترکیب می‌شوند. به این منظور از فیلترهای مکانی مختلفی برای استخراج ویژگی‌ها و کاهش بُعد استفاده می‌شود. یکی از روش‌های متداول در کاربرد فعلی الگوریتم (CSP) common spatial patterns را به صورت کامل توضیح داده و آن را برابر روی داده‌ها اعمال کنید.

توسط TSNE، نمودار scatter plot نمونه‌ها را در فضای دو بُعدی، قبل و بعد از اعمال این الگوریتم رسم کرده و عملکرد این گام را در استخراج ویژگی‌ها بررسی کنید.

#### ۲.۰۳ (۶۰ نمره)

#### ۲.۰۴ (۴۰ نمره) پیاده‌سازی Kernel SVM

در این بخش، دیتابست پردازش شده‌ی حاصل از مراحل قبل را دسته بندی می‌کنیم. به این منظور روش SVM با کرنل RBF را پیاده‌سازی کرده و توسط آن دسته بندی دو کلاسه انجام دهید. برای اطمینان از صحت پیاده سازی خود می‌توانید عملکرد نهایی مدل خود را با نسخه‌ی آماده از scikit-learn مقایسه کنید.

#### ۲.۰۵ (۱۰ نمره) ارزیابی

به منظور ارزیابی عملکرد، متريک‌های precision, recall, F1, accuracy و ROC و confusion matrix را گزارش کنید.

#### ۰.۳۰۲ (۱۰ نمره) مقاييسه با روش‌های دسته بندی دیگر

به انتخاب خود سه روش يادگيري ماشين دیگر را انتخاب کرده و توسط آنها نيز دسته بندی را انجام داده و نتایج را مقاييسه کنيد. دقت کنيد که برای اين بخش نيازی به پياده سازی وجود ندارد و می‌توانيد از مدل‌های آماده‌ی scikit-learn استفاده کنيد.

#### ۰.۳ (۱۵ نمره) خوش‌بندی

توسط الگوريتم k-means داده‌ها را خوش‌بندی کرده و scatter plot آن را در فضای دو بعدی رسم کنيد. تعداد خوش‌بندی‌هاي بهينه را توسيط رسم کردن نمودار معيارهای silhouette score و WCSS به ازاي k های مختلف به دست آوريد.

#### ۰.۴. معيار رقابت

علاوه بر فایل BCICIV\_calib\_ds1a.mat که تا به اين جا از آن استفاده کردیم، توسيط پايپلайн پياده ساري شده، فایل BCICIV\_calib\_ds1c.mat را نيز دسته بندی کرده و نتایج را مقاييسه کنيد. ميانگين دقت به دست آمده روی اين دو فایل به عنوان معيار رقابت در نظر گرفته می‌شود و به تيمی که بالاترین دقت را به دست آورده باشد نمره‌ی امتيازی تعلق خواهد گرفت.

به منظور بالاتر بردن دقت نهايی، مجاز به استفاده از هر تكينيک دلخواه برای پيش پردازش، تنظيم هايپر پaramترها و يا انتخاب مدل می‌باشيد.

## پروژه شماره ۲۵ (۱۲۰ نمره)

### هدف پروژه

هدف این پروژه، پیاده‌سازی یک سامانه‌ی تشخیص اشیاء (Object Detection) برای شناسایی و مکان‌یابی موجودات زیرآب در تصاویر، با تکیه بر معماری Faster R-CNN است. شما باید یک آشکارساز دو مرحله‌ای طراحی کنید که (۱) نواحی کاندید (Region Proposals) را تولید کند و (۲) برای هر ناحیه، کلاس شیء و مختصات جعبه‌ی مرزی (Bounding Box) را پیش‌بینی نماید. در تمام مراحل، همه چیز باید از پایه پیاده‌سازی شود و استفاده از مدل آماده/از پیش آموزش دیده مجاز نیست.

### قواعد و محدودیت‌های مهم (برای کل پروژه)

- پیاده‌سازی از پایه: همه‌ی بخش‌های اصلی سیستم (پیش‌پردازش، دیتالودر، معماری، زیان‌ها، متريک‌ها، پس‌پردازش NMS، و حلقه آموزش) باید توسط شما پیاده‌سازی شود.
- ممنوعیت استفاده از مدل آماده: استفاده از مدل‌های از پیش آموزش دیده (وزن‌های آماده) و مأذول‌های آماده‌ی تشخیص شیء (مثل torchvision.models.detection) ممنوع است؛ اما استفاده از لایه‌های پایه مانند Conv2d، ReLU و BatchNorm مجاز است.
- ممنوعیت کد آماده: استفاده از کدهای آماده‌ی اینترنتی برای بخش‌های اصلی پروژه مجاز نیست و کد باید محصول کار خودتان باشد.
- زبان برنامه‌نویسی: فقط Python مجاز است.
- تحويل کد: کدها باید در قالب ipynb باشد و همه‌ی سلول‌ها اجرا شده و خروجی‌ها ذخیره شده باشد.
- گزارش‌نویسی: تمرکز گزارش روی نتایج، تحلیل و تصمیمات طراحی است (نه توضیح خطبه‌خط کد).

### مجموعه‌دادگان (Dataset)

از مجموعه‌دادگان Underwater Object Detection استفاده کنید (قابل دریافت از Kaggle). این مجموعه شامل ۶۳۸ تصویر از موجودات زیرآب در ۷ کلاس است و برای هر شیء در تصویر، جعبه‌ی مرزی (Bounding Box) ارائه شده است.

### ساختار داده و تقسیم‌بندی

- از زیرپوشش‌های validation، train و test همان‌طور که در دیتاست وجود دارد استفاده کنید.
- فرمت Annotation (مثل XML/JSON/TXT) را شناسایی کرده و تبدیل/استانداردسازی لازم را انجام دهید.
- تمام مراحل خواندن و تبدیل برچسب‌ها باید از پایه پیاده‌سازی شود.

### تعریف مسئله (Problem Definition)

#### ورودی

یک تصویر RGB یا BGR از صحنه‌ی زیرآب.

#### خروجی

برای هر تصویر، مجموعه‌ای از پیش‌بینی‌ها شامل کلاس هر شیء (۷ کلاس)، مختصات جعبه‌ی مرزی ( $x_{min}, y_{min}, x_{max}, y_{max}$ )، confidence score (نموده اطمینان) و نمره اطمینان (Soft-NMS) باید از پایه پیاده‌سازی شود.

### مراحل انجام پروژه و روش نمره‌دهی

### مرحله ۱: پرسش‌های تشریحی و مبانی نظری (۱۰ نمره)

۱. (۳ نمره) مقایسه خانواده Region-based CNN، R-CNN و Fast R-CNN: معماری‌های R-CNN و Faster R-CNN را مختصر توضیح دهید و بگویید هر کدام چه مشکل/گلوگاهی را حل کرده‌اند.

۲. (۳ نمره) مقایسه آشکارسازهای تک مرحله‌ای و دو مرحله‌ای: مزايا/معایب Stage-one (مثل YOLO/SSD) و Stage-two (مثل Faster R-CNN) را بیان کنید و کاربردهای رایج هر کدام را مثال بزنید.

۳. (۴ نمره) مفاهیم GIoU، Soft-NMS و OHEM: هر کدام را تعریف کنید، مزايا و کاربردشان را بگویید و توضیح دهید چرا ممکن است روی داده‌های زیر آب مفید باشد.

## مرحله ۲: معرفی دیتاست، EDA و پیش‌پردازش (۱۰ نمره)

۱. نمایش نمونه‌ها: حداقل ۱۰ تصویر به همراه Bounding Box‌های واقعی و نام کلاس‌ها را نمایش دهید.
۲. آماری: موارد زیر را محاسبه و به صورت نمودار/جدول گزارش کنید: تعداد تصاویر در هر کلاس، توزیع تعداد اشیاء در هر تصویر، توزیع اندازه‌ی Bounding Box‌ها (عرض/ارتفاع/نسبت تصویر).
۳. پیش‌پردازش: تصمیمات پیش‌پردازشی (مثلاً تغییر اندازه، نرمال‌سازی، تبدیل رنگ، حذف تصاویر خراب) را توضیح دهید و علت هر تصمیم را بنویسید.

## مرحله ۳: تقویت داده (۱۰ نمره)

۱. تقویت‌های پایه: با توجه به ماهیت تصاویر زیر آب، چند روش مناسب (چرخش، برش، تغییر روشنایی/کنتراست، Color Jitter، Blur...) انتخاب کنید. برای هر روش دلیل انتخاب، نحوه اثرگذاری، و نحوه اعمال روی Bounding Box‌ها (بهروزرسانی صحیح مختصات) را بنویسید و چند نمونه قبیل/بعد را نمایش دهید.

۲. Mosaic Augmentation: ایده Mosaic را توضیح دهید (مزايا/معایب/کاربردها) و سپس آن را برای دیتاست خود پیاده‌سازی کنید. خروجی چند نمونه موزاییکی را همراه با جعبه‌های مرزی صحیح نمایش دهید.

## مرحله ۴: ساخت دیتالودر و collate\_fn اختصاصی (۱۰ نمره)

۱. به دلیل متفاوت بودن تعداد اشیاء در تصاویر، طول برچسب‌ها در یک batch یکسان نیست. یک collate\_fn که برای هر batch: لیست boxes و labels را تا طول بیشترین تعداد شیء در آن pad کند؛ mask مناسب برای تشخیص داده‌های واقعی از padding تولید کند؛ و همه چیز را به tensorهای مناسب تبدیل کند.

۲. یک batch نمونه را چاپ/نمایش دهید و نشان دهید padding درست انجام شده است.

## مرحله ۵: تعریف مسئله و طراحی معماری Faster R-CNN (۲۵ نمره)

۱. Backbone (استخراج ویژگی): یک شبکه‌ی کانولوشنی مناسب را خودتان طراحی کنید (می‌تواند سبک‌تر از ResNet باشد). استفاده از وزن‌های آماده/پیش‌آموزش دیده مجاز نیست.

۲. ROI Pooling/ROI Proposals: معماری کلی Faster R-CNN را توضیح دهید: تولید RPN + Head، سپس سری تشخیص (کلاس‌بندی + رگرسیون جعبه).

۳. استفاده از GIoU: در بخش‌های مرتبط با همپوشانی و زیان رگرسیون جعبه‌ها، از GIoU به جای IoU استفاده کنید و دقیق توضیح دهید در کدام قسمت‌ها (مثلاً برچسب‌گذاری anchors یا زیان رگرسیون).

۴. جزئیات پیاده‌سازی: توضیح دهید ابعاد ورودی/خروچی هر ماژول چیست، و چگونه tensorها را بین اجزا منتقل می‌کنید.

## مرحله ۶: طراحی و پیاده‌سازی Region Proposal Network (RPN) (۲۵ نمره)

۱. تولید Anchor Boxes و Anchor Points: روی feature map نقطه anchor ایجاد کنید و برای هر نقطه  $k$  جعبه با اندازه/نسبت‌های مختلف بسازید. مقدار  $k$  و مجموعه مقیاس‌ها/نسبت‌ها را خودتان انتخاب کنید و دلیل انتخاب را بنویسید.

۲. نمایش بصری: anchors را روی تصویر از دیتاست رسم و نمایش دهید.

۳. برچسب‌گذاری مثبت/منفی: با یک/چند آستانه مناسب، anchors را به positive و negative تقسیم کنید. سپس نمونه‌های مثبت/منفی را برای ۲ تصویر نمایش دهید و تمام پارامترهای به کاررفته (آستانه‌ها، تعداد نمونه‌برداری و ...) را گزارش کنید.

## مرحله ۷: آموزش مدل کامل وتابع زیان (۱۵ نمره)

۱. تابع زیان: درباره زیان‌های رایج Faster R-CNN تحقیق کنید و یک زیان چندبخشی پیاده‌سازی کنید (مثلاً زیان objectness برای RPN، زیان رگرسیون جعبه برای RPN (ترجیحاً مبتنی بر IoU) یا Smooth-L1 (G) یا زیان کلاس‌بندی برای ROI head، و زیان رگرسیون جعبه برای ROI head).

۲. گزارش تنظیمات: نرخ یادگیری، بهینه‌ساز، epoch، تعداد batch size، وزن‌دهی بخش‌های زیان، و هر تنظیم مهم را کامل گزارش کنید.

۳. نمودارهای آموزش: نمودار روند آموزش (زیان کل و زیان‌های اجزا، و در صورت امکان mAP روی اعتبارسنجی) را رسم کنید.

## مرحله ۸: ارزیابی و نمایش نتایج (۱۵ نمره)

۱. ارزیابی کمی: حداقل یک معیار مناسب تشخیص شیء را گزارش کنید (پیشنهاد: mAP@0.5:0.95 و/یا mAP@[0.5:0.95]). روش محاسبه را توضیح دهید و از پایه پیاده‌سازی کنید.

۲. ارزیابی کیفی: ۱۰ نمونه از تصاویر تست را نمایش دهید به همراه جعبه‌های واقعی، جعبه‌های پیش‌بینی شده (پس از (Soft-)NMS)، نام کلاس و نمره اطمینان.

۳. تحلیل خطأ: حداقل ۵ نمونه‌ی شکست/اشتباه را انتخاب کنید و دلیل احتمالی را تحلیل کنید.

## خلاصه ورودی/خروجی

توضیح	مولفه
تصویر زیر آب	ورودی
مجموعه‌ای از (کلاس، Bounding Box، نمره اطمینان)	خروجی
پیاده‌سازی شده از پایه Faster R-CNN	مدل
پیاده‌سازی شده از پایه (Soft-)NMS	پس‌پردازش

## نحوه ارزیابی و رتبه‌بندی نهایی (نسخه‌ی ساده)

### خروجی مدل.

برای هر تصویر، مدل مجموعه‌ای از تشخیص‌ها به صورت  $\{(b_i, c_i, s_i)\}$  تولید می‌کند. ارزیابی روی خروجی نهایی پس از NMS انجام می‌شود.

### IoU.

$$\text{IoU}(A, B) = \frac{\text{area}(A \cap B)}{\text{area}(A \cup B)}.$$

**قانون تطبیق.**

اگر  $\text{IoU} \geq 0.5$  باشد و آن گراندتروث قبل از تطبیق داده نشده باشد، تشخیص درست محسوب می‌شود.  
محاسبه‌ی امتیازها.

$$\text{Precision}_c = \frac{TP_c}{TP_c + FP_c}, \quad \text{Recall}_c = \frac{TP_c}{TP_c + FN_c}$$

$$\text{Score} = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} F^{\gamma_c}.$$

## پروژه شماره ۳ (۱۱۵ نمره)

### تعریف مسئله

در این پروژه، هدف طراحی و پیاده‌سازی یک مدل یادگیری عمیق مبتنی بر معماری Transformer برای تحلیل احساسات<sup>۱</sup> متون مالی است. شما باید مدلی بسازید که بتواند اخبار و جملات مالی را پردازش کرده و آنها را در سه دسته مثبت، منفی یا خنثی طبقه‌بندی کنند.

نکته کلیدی این پروژه این است که باید از کلاس‌های آماده مدل استفاده کنید، بلکه باید لایه‌های اصلی ترانسفورمر (بهویژه مکانیزم Self-Attention) را خودتان پیاده‌سازی کنید.

دیتاست: برای آموزش و ارزیابی اولیه، از دیتاست استاندارد Financial PhraseBank استفاده کنید. می‌توانید این داده‌ها را مستقیماً از لینک زیر دریافت کنید:

[Financial PhraseBank on Hugging Face](#)

**نکته مهم:** برای داده‌های Test نهایی جهت ارزیابی در سیستم داوری، یک فایل جداگانه در اختیار شما قرار می‌گیرد که برچسب (Label) ندارد. لینک دانلود این داده در بخش مربوطه آمده است.

**قسمت اول: آماده‌سازی داده برای مدل و پیش‌پردازش (۲۰ نمره)**  
در این گام باید داده‌ها را برای ورود به شبکه عصبی آماده کنید.

۱. (۳ نمره) دانلود و بازگذاری صحیح داده‌ها از فایل ZIP (با توجه به ساختار خاص فایل‌های Financial PhraseBank که شامل چندین فایل متنی است) و تجمعی آنها در یک DataFrame واحد. به همراه رسم توزیع داده‌های هر دسته.

۲. (۱۵ نمره) پیش‌پردازش درست بر داده‌ها و دسته‌بندی ۸۰-۲۰ بین داده train و validation.

۳. (۲ نمره) پیاده‌سازی کلاس Dataset در PyTorch و ایجاد DataLoaders با استفاده از یک Tokenizer استاندارد (مانند BERT Tokenizer) برای تولید input ids و attention mask.

### قسمت دوم: ساخت مدل (۴۵ نمره)

شما باید در این بخش معماری ترانسفورمر را بدون استفاده از ماثولوهاز آماده سطح بالا پیاده‌سازی کنید.

۱. (۱۵ نمره) پیاده‌سازی کلاس SelfAttention از صفر. باید محاسبات  $Q, K, V$  را انجام دهید و فرمول زیر را با استفاده از عملیات ماتریسی (مثل matmul torch.einsum یا) پیاده‌سازی کنید:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

پشتیبانی از Multi-Head Attention در این بخش الزامی است.

۲. (۱۵ نمره) پیاده‌سازی کلاس TransformerBlock شامل لایه‌های Feed For-ward Network و LayerNorm و اتصال‌های باقیمانده (Residual Connections).

۳. (۱۵ نمره) تجمعی بلوک‌ها در کلاس اصلی مدل (FinancialTransformer) شامل لایه‌های Embedding (برای کلمات و موقعیت‌ها) و یک لایه Classifier نهایی برای تولید خروجی ۳ کلاسه.

**قسمت سوم: آموزش مدل (۲۰ نمره)**  
در این مرحله مدل طراحی شده را آموزش می‌دهید.

Sentiment Analysis<sup>۱</sup>

۱. (۱۰ نمره) نوشتن حلقه آموزش (Training Loop) استاندارد با استفاده از Optimizer و تابع هزینه مناسب.
۲. (۱۰ نمره) ثبت و رسم نمودار Loss و Accuracy برای هر دو مجموعه Train و Validation در پایان هر Epoch جهت پایش فرآیند یادگیری.

#### قسمت چهارم: ارزیابی دقیق مدل بر داده‌ی validation (۱۰ نمره)

در این مرحله باید عملکرد نهایی مدل آموزش دیده را با معیارهای استاندارد روی داده‌های اعتبارسنجی بسنجید تا کیفیت مدل فراتر از نمودار خطای مشخص شود.

۱. (۴ نمره) محاسبه معیارهای ارزیابی کمی: برای کل مجموعه داده Validation، معیارهای زیر را محاسبه و گزارش کنید. با توجه به اینکه مسئله چند کلاسه (۳ کلاس) است، در صورت استفاده از کتابخانه‌هایی مثل scikit-learn، حتماً نوع میانگین‌گیری (مثلاً weighted macro) را مشخص کنید:

- Accuracy •
- Precision •
- Recall •
- F1-Score •

۲. (۶ نمره) رسم ماتریس درهم‌ریختگی (Confusion Matrix): ماتریس درهم‌ریختگی را محاسبه کرده و آن را به شکل گرافیکی رسم کنید. محورهای نمودار باید دارای برچسب‌های واضح کلاس‌ها (Negative, Neutral, Positive) باشند تا مشخص شود مدل بیشترین خطای را در تشخیص کدام کلاس‌ها دارد.

#### قسمت پنجم: ارزیابی خروجی مدل بر داده‌ی test (۵ نمره)

در این بخش عملکرد نهایی مدل خود را روی داده‌های دیده نشده می‌سنجید.

۱. (۵ نمره) تولید فایل خروجی برای مسابقه: مدل خود را روی [این جملات](#) اجرا کرده و پیش‌بینی‌ها را در یک فایل CSV ذخیره کنید. متعاقباً روش ارزیابی این بخش اعلام می‌شود و شما باید فایل csv را در مکانی که گفته می‌شود آپلود کنید تا نمره نهایی داوری و محاسبه شود.

#### قسمت ششم: خروجی بصری مکانیزم توجه (۱۵ نمره)

برای درک بهتر عملکرد مدل، باید نشان دهید مدل به کدام کلمات توجه بیشتری دارد.

۱. (۱۵ نمره) استخراج ماتریس وزن‌های توجه (Attention Weights) از آخرین لایه مدل و رسم Heatmap برای حداقل ۵ جمله اول از داده‌های تست بخش قبل. خروجی باید نشان دهد که مدل برای پیش‌بینی حس جمله، روی چه کلماتی تمرکز کرده است مثلاً کلماتی مثل "growth", "loss", "risk".

## پروژه شماره ۴ (۱۰۰ نمره)

### تعریف مسئله:

هدف این پروژه طراحی و پیاده‌سازی یک مدل یادگیری ماشین برای پیش‌بینی فروش آتی فروشگاه‌های زنجیره‌ای Rossmann است. علاوه بر پیش‌بینی دقیق فروش، دانشجویان باید بتوانند ویژگی‌های پیشرفته‌تر مدل‌سازی، مانند برآورد عدم قطعیت، دسته‌بندی فروش، و تحلیل مدل را پیاده‌سازی کرده و تحلیل دقیقی از عملکرد ارائه دهند.

### مشخصات دیتابست:

- داده‌ها از رقابت فروش فروشگاه‌های Rossmann استخراج شده است.
- تعداد نمونه‌ها: بیش از ۱ میلیون ردیف فروش روزانه برای بیش از ۱۱۰۰ فروشگاه
- ستون هدف: Sales
- فایل‌های جانبی: store.csv (حاوی مشخصات هر فروشگاه شامل نوع فروشگاه، فاصله تا رقیب، نوع سرویس‌دهی و ...)
- ویژگی‌ها: ترکیبی از عددی، باینری و زمانی (تاریخ، روز هفت، ماه و ...)

### مراحل پروژه و نمره‌بندی:

#### فاز ۱: Data Visualization - تحلیل اکتشافی و شناخت داده‌ها - ۵ نمره

- رسم نمودارهای روند فروش در طول زمان برای چند فروشگاه منتخب
- بررسی تأثیر ویژگی‌های مختلف (StateHoliday، SchoolHoliday، Promo و ...) بر فروش
- نمودار همبستگی متغیرها و تحلیل آن
- مقایسه روند فروش در فروشگاه‌های مختلف

#### فاز ۲: Feature Engineering - مهندسی ویژگی و آماده‌سازی داده - ۱۵ نمره

- ترکیب جداول sales و store ID بر اساس Store ID
- ایجاد ویژگی‌های زمانی جدید (روز هفت، ماه، تعطیلی، میانگین متحرک و ...)
- ایجاد lag features و آماره‌های متحرک (فروش روز قبل، میانگین ۷ روز اخیر، انحراف معیار و ...)
- استخراج ویژگی‌های دوره‌ای (مثل Fourier terms برای فصل‌ها) و فاصله تا تعطیلات رسمی
- نرمال‌سازی و مدیریت داده‌های گمشده
- تفکیک داده به آموزش/اعتبارسنجی/ تست با حفظ ترتیب زمانی

#### فاز ۳: Learn & Estimate - آموزش و ارزیابی مدل‌ها - ۳۰ نمره

توجه کنید در این فاز باید کامل مدل را پیاده‌سازی کرده و استفاده از کتابخانه‌ها مجاز نیست.

- پیاده‌سازی مدل baseline ساده (رگرسیون خطی یا میانگین متحرک) و ارزیابی آن
- آموزش مدل‌های کلاسیک مانند LightGBM یا XGBoost با تنظیم hyperparameter
- استفاده از اعتبارسنجی زمانی مبتنی بر فروشگاه (store-aware CV) و گزارش RMSE در هر مرحله

- تحلیل feature importance با استفاده از SHAP

• تحلیل خطای بررسی مواردی که مدل در آنها عملکرد ضعیفی داشته به طور دقیق با تفسیر و مدرک

#### فاز ۴: Uncertainty & Sales Classification - تحلیل عدم قطعیت و طبقه‌بندی فروش - ۲۵ نمره

- پیاده‌سازی پیش‌بینی چندگامی (multi-step forecasting) (فروش ۷ روز آینده) با مدیریت انتشار خطای ensemble quantile regression یا uncertainty با روش‌هایی مثل
- دسته‌بندی فروش به سه کلاس (کم، متوسط، زیاد) و آموزش مدل طبقه‌بند چندکلاسه با ROC-AUC و F1-score

#### فاز ۵: Deep - قسمت عمیق پروژه - ۲۵ نمره

- پیاده‌سازی دقیق uncertainty estimation + multi-step forecasting + SHAP
- استفاده از مدل‌های یادگیری عمیق LSTM و TCN با مستندات کامل
- بهینه‌سازی هایپرپارامترها با Optuna یا Ray Tune

#### فاز ۵ + ۱: مستندسازی و ارائه نتایج

- گزارش کامل نتایج پیش‌بینی، تحلیل خطاهای نمودارها و تفسیر نهایی
- مستندسازی مرتب، ساختارمند و قابل فهم

#### معیار رقابت بین تیم‌ها:

- معیار اصلی: Root Mean Squared Error مدل نهایی بر روی مجموعه تست
- معیارهای مکمل:  $R^2$  و Mean Absolute Percentage Error
- برای طبقه‌بندی: F1 Score و ROC-AUC

## پروژه شماره ۵ (۱۲۰ نمره)

### تعریف مسئله

در این پروژه هدف این است که یک مدل مولد مبتنی بر latent space برای تصاویر بسازید و سپس آن را بهبود دهید تا (۱) نمونه‌های تولیدی با کیفیت‌تری ایجاد کند، (۲) فضای نهفته رفتار قابل تفسیرتری داشته باشد، و (۳) بتوان تولید را با یک ورودی کنترلی (برچسب کلاس) هدایت کرد.

دیتابست این پروژه: Fashion-MNIST (۱۰ کلاس، تصاویر خاکستری  $28 \times 28$ )

### فاز اول: آماده‌سازی داده‌ها و تحلیل اکتشافی (۱۰ نمره)

هدف این فاز آماده کردن ورودی استاندارد برای مدل و ارائه یک تصویر کلی از توزیع داده‌ها است. آماده‌سازی داده و ساخت ورودی مدل (۷ نمره)

دیتابست Fashion-MNIST را دریافت کنید و ورودی مناسب برای آموزش مدل را بسازید. موارد زیر باید در نوت‌بوک شما به صورت شفاف وجود داشته باشد:

- دانلود/بارگذاری داده و گزارش تعداد نمونه‌های train/test/validation و شکل هر نمونه.
- تعریف یک pipeline پیش‌پرداز.
- ساخت train/test/validation برای DataLoader.
- یک sanity check: نمایش یک batch (حداقل ۲۵ تصویر) به همراه برچسب کلاس.
- ثبت seed برای بازتولید پذیری نتایج.

### تحلیل اکتشافی (EDA) و نمودارها (۳ نمره)

- نمودار میله‌ای توزیع کلاس‌ها (۱۰ کلاس) در مجموعه آموزش.
- هیستوگرام شدت پیکسل‌ها (برای کل داده یا چند کلاس منتخب).
- یک هیستوگرام آماری ساده در سطح تصویر (یکی از موارد زیر):
  - میانگین شدت پیکسل هر تصویر، یا
  - تعداد پیکسل‌های غیرصفر در هر تصویر.

### فاز دوم: پیاده‌سازی VAE و بهبود کیفیت خروجی (۵ نمره) پیاده‌سازی VAE از ابتدا (۳۰ نمره)

یک VAE را پیاده‌سازی و آموزش دهید (reparameterization trick + Decoder + Encoder from scratch) تابع هزینه باید شامل موارد زیر باشد:

- reconstruction loss (یکی از BCE یا MSE را انتخاب و ثابت نگه دارید).
  - یک ترم regularization روی latent space (در حالت پایه با وزن ۱).
- خروجی مورد انتظار از این بخش:
- نمایش نمونه‌های reconstruction (حداقل ۲۰ تصویر ورودی/خروجی).
  - نمایش نمونه‌های sampling (حداقل ۵۰ تصویر).
  - گزارش روی test set: مقادیر reconstruction، regularization و total regularization.

**بهبود کیفیت خروجی و معیار رقابت (۲۰ نمره)**  
حداصل یک تغییر برای بهبود کیفیت خروجی اعمال کنید (نام روش مهم نیست). یکی از موارد زیر را انجام دهید:

- تغییر convolutional architecture (مثلاً convolutional architecture یا عمیقتر).
- تغییر latent dimension.
- تغییر شدت یا برنامه‌ی regularization در طول آموزش.
- مقایسه با مدل پایه باید شامل موارد زیر باشد:
  - جدول متريک‌های test set.
  - چند نمونه sampling و reconstruction (مقایسه قبل و بعد از بهبود).

#### معیار رقابت بین تیم‌ها:

معیار اصلی رقابت استفاده از FID (Frechet Inception Distance) روی ویژگی‌های استخراج شده از classifier ارائه شده در درس است. برای این کار:

- حداصل 10,000 تصویر تولید کنید.

- FID را بین تصاویر تولیدی و تصاویر test set محاسبه کنید (با استفاده از classifier ارائه شده).  
تیم‌های برتر که کمترین FID را ثبت کنند، امتیاز اضافی دریافت می‌کنند.

**فاز سوم: کنترل رفتار latent space با تغییر وزن regularization (۲۰ نمره)**  
در این فاز هدف این است که با تنظیم دقیق شدت جریمه بر فضای نهفته، به مدلی دست یابید که ویژگی‌های بصری داده‌ها را به شکلی تغییر کنند و معنادار در ابعاد خود ذخیره کند.

**تحلیل اثر وزن دهنده روی latent space (۲۰ نمره)**  
در این فاز باید شدت ترم regularization را تغییر دهید تا رفتار latent space قابل تفسیرتر شود. منظور این است که با تغییرات کوچک و کنترل شده در بردار نهفته، تغییرات معنادار و پیوسته در تصویر خروجی مشاهده شود.  
موارد الزامی که باید در گزارش شما مستند شود:

- حداصل سه مقدار متفاوت برای وزن ترم regularization انتخاب کنید (مثلاً ۰.۵، ۱، ۴).
- برای هر مقدار، مدل را آموزش داده و روی test set reconstruction گزارش کنید: مقادیر regularization و total.
- برای هر مقدار، عملیات latent traversal را به شرح زیر انجام دهید:
  - انتخاب حداصل ۵ بعد متفاوت از بردار نهفته (latent vector).
  - برای هر بعد، تولید حداصل ۷ مقدار پیوسته (مثلاً در بازه -۳ تا +۳).
  - نمایش خروجی‌های بازسازی شده به صورت یک شبکه (grid) منظم.
- نگارش یک پاراگراف جمع‌بندی که تحلیل کند کدام وزن، بهترین توازن بین کیفیت بازسازی (reconstruction) و قابلیت تفسیر فضای نهفته را ایجاد کرده است.

**فاز چهارم: تولید کنترل پذیر با استفاده از label (۲۰ نمره)**  
در این فاز، هدف ارتقای مدل به یک خودرمزگذار متغیر شرطی (Conditional VAE) است تا فرآیند تولید تصویر تحت کنترل برچسب‌های کلاس قرار گیرد.

**تولید تصویر با کنترل کلاس (۲۰ نمره)**  
در این فاز باید کاری کنید که مدل بتواند با دریافت label (کلاس تصویر)، نمونه‌هایی مطابق همان کلاس تولید کند.  
موارد الزامی که باید رعایت شوند:

- برچسب (label) را به عنوان ورودی کنترلی به مدل اضافه کنید (طراحی دقیق معماری با تیم).
- مدل را آموزش دهید و حداقل برای هر یک از 10 کلاس:
  - حداقل 20 تصویر تولید کنید.
  - خروجی‌ها را به صورت یک شبکه (grid) نمایش دهید.
- برای سنجش کنترل‌پذیری، از classifier ارائه شده در درس استفاده کنید و گزارش دهید:
- دقت دسته‌بندی (classification accuracy) روی تصاویر تولیدی (به تفکیک هر کلاس و میانگین کل).
- یک مقایسه کوتاه با حالت بدون label ارائه دهید (شامل چند نمونه تصویری و یک پاراگراف توضیح).

**فاز پنجم: گزارش نهایی و جمع‌بندی (۲۰ نمره)**  
این مرحله شامل تجمعی یافته‌ها و ارائه مستندات نهایی پژوهه به شکلی ساختاریافته است.

**گزارش و ارائه نتایج (۲۰ نمره)**  
در این فاز باید یک گزارش نهایی تهیه کنید که تمام مراحل پژوهه را به صورت خلاصه، شفاف و قابل ارزیابی ارائه دهد.  
موارد الزامی گزارش:

- ارائه یک جدول نهایی از نتایج روی test set برای همه مدل‌ها و تنظیمات اجرا شده شامل مقادیر: reconstruction, total و regularization.
- نمایش نمونه‌های کلیدی (به صورت grid) برای هر فاز شامل:
  - بازسازی (reconstruction) و نمونه‌برداری (sampling) برای مدل پایه.
  - مقایسه قبل و بعد برای بهبود کیفیت خروجی.
  - پیمایش فضای نهفته (latent traversal) برای وزن‌های مختلف.
  - تولید کنترل‌پذیر با برچسب برای تمامی کلاس‌ها.
- جمع‌بندی کوتاه (حداکثر یک صفحه) شامل شناسایی بهترین تنظیمات از نظر کیفیت تولید و توازن با قابلیت تفسیر، به همراه بیان محدودیت‌ها و پیشنهادات برای توسعه آتی.
- کدها باید کاملاً قابل اجرا بوده و تنظیمات بازتولید‌پذیری (seed) در ابتدای نوت‌بوک درج شده باشد.

## پروژه شماره ۶ (۱۱۰ نمره)

**بیان مسئله: تحلیل احساسات جنبه‌محور (Aspect-Based Sentiment Analysis)** در بسیاری از کاربردهای واقعی، دانستن اینکه یک نظر به طور کلی «مثبت» یا «منفی» است، کافی نیست. برای مثال در جمله «قیمت محصول مناسب بود ولی کیفیت ساخت پایینی داشت»، با دو حس متضاد روبرو هستیم. هدف این پروژه پیاده‌سازی سیستمی مبتنی بر یادگیری عمیق است که بتواند احساسات را نسبت به جنبه‌های خاص موجود در متن تشخیص دهد.

### فاز اول: آماده‌سازی داده‌ها و تحلیل اکتشافی (۱۵ نمره)

مجموعه داده استاندارد ۴-SemEval-2014-Task (زیرمجموعه رستوران یا لپتاپ) را دریافت کنید.

- **استخراج داده (۱۰ نمره):** از مجموعه داده‌ی دریافتی داده‌ی ورودی مدل خود را استخراج کنید. توجه کنید که جملات ممکن است دارای چند جنبه باشند. به همین دلیل ممکن است هر جمله به چند نمونه آموزشی مجزا تبدیل شود.

- **تحلیل داده (۵ نمره):** نمودارهای زیر را برای تحلیل داده رسم کرده و در گزارش بیاورید:
  - نمودار میله‌ای توزیع کلاس‌ها (مثبت، منفی، خنثی).
  - نمودار هیستوگرام تعداد کلمات در جملات.
  - ابر کلمات (Word Cloud) برای کلمات مربوط به جنبه‌های مختلف.

### فاز دوم: طراحی و پیاده‌سازی مدل (۶۰ نمره)

در این بخش معماری مدل مبتنی بر ترنسفورم‌ها را برای شناسایی قطبیت هر جنبه طراحی می‌کنید.

- **ورودی مدل (۱۰ نمره):** از آنجا که ورودی مدل شامل «جمله» و «کلمه جنبه» است، باید آن‌ها را به گونه‌ای ترکیب کنید که ورودی قابل قبول برای مدل فراهم شود (معمولًاً در پیاده‌سازی کلاس Dataset).

- **آموزش مدل (۴۰ نمره):** یک مدل زبانی پیش‌آموزش دیده دلخواه (مانند BERT-Base یا RoBERTa) را بارگذاری کرده و لایه آخر آن را برای کلاس‌بندی ۳ کلاسه تغییر دهید. مدل را روی داده‌هایی که آماده کرده‌اید آموزش دهید.

- **تنظیمات و هایپرپارامترها (۵ نمره):** تمامی تنظیمات استفاده شده برای آموزش (مانند Learning Rate و ... ) باید دقیقاً گزارش شوند.

- **نمودارهای آموزش (۵ نمره):** نمودار تغییرات loss و accuracy را برای داده‌های آموزشی و اعتبارسنجی در طول اپاک‌های مختلف رسم کنید.

### فاز سوم: ارزیابی و تحلیل خطأ (۱۵ نمره)

مدل آموزش دیده را روی داده‌های تست ارزیابی کرده و موارد زیر را در گزارش نهایی درج کنید:

- **گزارش دسته‌بندی (۵ نمره):** محاسبه معیارهای Precision، Recall و F1-Score بودن داده‌ها، معیار Macro-F1 را نیز حتماً گزارش دهید.

- **ماتریس آشتفتگی (۵ نمره):** ماتریس آشتفتگی (Confusion Matrix) را بر اساس خروجی‌های مدل رسم کنید.

- **تحلیل خطأ (۵ نمره):** ۵ نمونه از نظراتی که مدل به غلط پیش‌بینی کرده است را انتخاب کنید. متن جمله، برچسب واقعی و پیش‌بینی مدل را بنویسید و دلیل احتمالی خطأ را تحلیل کنید.

#### **فاز چهارم: تست مدل (۲۰ نمره)**

بررسی استحکام (Robustness) مدل در مواجهه با جملات جدید و چالش برانگیز.

- **تست جملات چالش برانگیز (۱۰ نمره):** ۲۰ جمله چالش برانگیز بنویسید (ساختمانهای شرطی، منفی در منفی و ...) و خروجی مدل را گزارش کنید.

- **حمله تخاصمی ساده (۱۰ نمره):** جمله‌ای بنویسید که مدل آن را درست دسته‌بندی می‌کند، سپس با تغییرات جزئی بدون تغییر معنا، سعی کنید مدل را به اشتباہ بیندازید.

---

**ارزیابی رقابتی:** شایان ذکر است که در پایان این پژوهه، تمامی تیم‌هایی که روی موضوع یکسانی فعالیت کرده‌اند، بر اساس معیارهای ارزیابی (بهویژه Macro-F1) با یکدیگر مقایسه خواهند شد. این مقایسه نه تنها برای سنجش کیفیت مدل‌های طراحی شده، بلکه برای رتبه‌بندی نهایی و شناسایی پیاده‌سازی‌های برتر انجام می‌گیرد.

---

## پروژه شماره ۷ (۱۰۵ نمره)

### تعریف مسئله

هدف این پروژه ساخت یک مدل Binary Classification کارآمد و دقیق برای شناسایی مشتریان در معرض ریزش (Churn) در حوزه مخابرات است. تیم پروژه باید بر انتخاب بهترین ویژگی‌ها، مدیریت عدم توازن داده‌ها و مقایسه مدل‌های پیشرفته تمرکز کند. تمرکز اصلی ارزیابی مدل‌ها بر بیشینه‌سازی معیار Recall خواهد بود.

### دیتاست: Telco Customer Churn Dataset

#### فاز اول: تحلیل اکتشافی داده‌ها (EDA) — ۱۵ نمره

- هدف: شناخت رفتار داده‌ها
- شناخت اولیه (۵ نمره)

- گزارش ابعاد دیتاست (Shape)، نوع متغیرها و مدیریت داده‌های تکراری (Duplicates)

#### • تحلیل تکمتغیره (۵ نمره)

- رسم نمودار Histogram برای متغیرهای عددی (مانند Tenure)
- رسم نمودار Bar Plot برای متغیرهای دسته‌ای
- تحلیل توزیع‌ها

#### • تحلیل دومتغیره و همبستگی (۵ نمره)

- رسم Heatmap برای بررسی همبستگی متغیرهای عددی
- استفاده از نمودارهای مقایسه‌ای مانند Boxplot برای بررسی رابطه مهم‌ترین ویژگی‌های عددی با متغیر هدف (Churn)

#### فاز دوم: پیش‌پردازش داده‌ها (Preprocessing) — ۱۵ نمره

##### • هدف: تمیزکاری داده‌ها

##### • مدیریت مقادیر گمشده (۵ نمره)

- شناسایی ستون‌های دارای مقدار گمشده (NaN)
- پیاده‌سازی Imputation همراه با توجیه متنی دقیق برای روش انتخابی

##### • انکودينگ و مقیاس‌بندی (۱۰ نمره)

- اعمال Label Encoding برای متغیرهای دودویی
- اعمال One-Hot Encoding برای متغیرهای دسته‌ای
- اعمال MinMaxScaler یا StandardScaler بر تمام ویژگی‌های عددی

#### فاز سوم: مهندسی و انتخاب ویژگی (Feature Engineering & Selection) — ۲۵ نمره

##### • هدف: افزایش قدرت مدل با استفاده از ویژگی‌های بهتر

##### • مهندسی ویژگی (۵ نمره)

- ایجاد حداقل دو ویژگی جدید معنادار (مانند ترکیب ویژگی‌ها یا بازه‌بندی Tenure)
- انتخاب ویژگی (۲۰ نمره)
  - روش فیلترمحور:
- استفاده از آزمون‌های آماری مانند Chi-Squared یا ANOVA
- انتخاب ۱۰ تا ۱۵ ویژگی برتر
- روش مدل‌محور:
  - استفاده از Feature Importance در مدل‌هایی مانند:
    - Lasso Regression \*
    - Random Forest \*
- الزامی:
  - ارائه توضیح متنی مستدل برای انتخاب زیرمجموعه نهایی ویژگی‌ها

#### فاز چهارم: مدل‌سازی پیشرفته و بهینه‌سازی (Modeling & Optimization) — ۴۰ نمره

- مدیریت عدم توازن داده‌ها (۱۰ نمره)
  - اعمال SMOTE روی داده‌های آموزشی
  - مقایسه تأثیر Class Weights و
- مدل‌های پایه (۱۰ نمره)
  - Logistic Regression –
  - SVM یا KNN –
  - Random Forest –
  - (LightGBM یا XGBoost) Gradient Boosting –
- آنسامبلینگ مدل (۵ نمره)
  - پیاده‌سازی Soft Voting Classifier با ۳ مدل برتر
- تنظیم فرآپارامتر (۱۰ نمره)
  - استفاده از RandomizedSearchCV یا GridSearchCV
  - گزارش بهترین پارامترها برای حداقل دو مدل
- اعتبارسنجی پیشرفته (۵ نمره)
  - استفاده از Stratified K-Fold Cross Validation
  - گزارش میانگین و انحراف معیار Recall

#### فاز پنجم: ارزیابی و گزارش نهایی (۱۰ نمره)

- گزارش معیارها (۴ نمره)

- ارائه جدول جامع شامل معیارهای F1-Score، Precision، Accuracy و Recall برای تمامی مدل‌های پیاده‌سازی شده.

- مقایسه ROC-AUC (۳ نمره)

- گزارش مقدار عددی AUC برای چهار مدل برتر پروژه جهت سنجش قدرت تفکیک کلاس‌ها.

- تحلیل نهایی (۳ نمره)

- رسم ماتریس آشفتگی (Confusion Matrix) برای تحلیل دقیق‌تر خطاهای مدل.

- معرفی مدل برنده و استخراج مهم‌ترین ویژگی‌های اثرگذار بر خروجی.

- ارائه پیشنهاد تجاری مشخص بر اساس تحلیل نتایج به دست آمده.

**بخش مسابقه** (Bonus Challenge) بیشینه‌سازی Recall برای کلاس ریزش ( $\text{Churn} = 1$ ) مشروط به اینکه F1-Score کمتر از ۵٪ نباشد. -

## پروژه شماره ۸ (۱۱۰ نمره)

Anomaly Detection on MVTec AD

### هدف کلی پروژه

در این پروژه، شما یک سیستم Anomaly Localization و Anomaly Detection برای تصاویر بازرسی صنعتی پیاده‌سازی می‌کنید. دیتابیس پروژه MVTec AD است. شما ابتدا baseline PaDiM را تکمیل و اجرا می‌کنید و سپس روش خودتان را پیاده‌سازی کرده و با baseline مقایسه می‌کنید.

### آنچه تحويل می‌دهيد

- یک نوت‌بوک شامل تمام TODO‌های پیاده‌سازی شده که بدون خطا اجرا شود.
- خروجی‌های گزارش محور: جدول‌های نتایج، نمودارها، و visualization کیفی (overlay + GT contour).

### قوانین و محدودیت‌ها

- آموزش (fit) فقط مجاز است از تصاویر train/good استفاده کند.
- خروجی predict باید شامل نمره‌ی سطح تصویر و نقشه‌ی ناهنجاری سطح پیکسل باشد.
- کد باید قابلیت تکرار داشته باشد (ثبت seed).

### فاز بندی و نمره‌دهی

#### فاز ۱ (۶ نمره): آماده‌سازی داده و Sanity Check

- پیاده‌سازی DataLoader و MVTecDataset ها شامل:
  - بارگذاری صحیح تصویر، برچسب و ماسک GT.
  - همتراز بودن ترتیب داده‌ها و بازگشت دیکشنری استاندارد در \_\_getitem\_\_.

#### فاز ۲ (۳۰ نمره): خط مبنا PaDiM

- استخراج ویژگی (۶ نمره): پیاده‌سازی ResNetFeatureExtractor.
- استخراج embeddings (۶ نمره): پیاده‌سازی تابع extract\_embeddings.
- برآش مدل آماری (۸ نمره): پیاده‌سازی تابع fit\_padim\_gaussian (میانگین/کوواریانس مکانی).
- محاسبه‌ی score map (۶ نمره): استفاده از فاصله Mahalanobis.
- اتصال مراحل (۴ نمره): تولید خروجی استاندارد در تابع run\_padim\_category.

#### فاز ۳ (۶۰ نمره): روش دانشجو (StudentMethod)

##### fit (۴۰ نمره):

- رعایت قوانین آموزش و عدم نشت اطلاعات (۱۰ نمره).
- آماده‌سازی بخش سطح تصویر (۱۰ نمره).
- آماده‌سازی بخش سطح پیکسل (۲۰ نمره).

**۲۰) predict :**

- تولید خروجی‌های صحیح و معنادار برای سطح تصویر (۱۰ نمره).
- تولید خروجی‌های سطح پیکسل با ابعاد صحیح و قابل ارزیابی (۱۰ نمره).

**فاز ۴ (۱۴ نمره): متريک‌ها، گزارش و مقاييسه**

**۱۰) پياده‌سازي متريک‌ها :**

- سطح تصویر: AUROC, AP, best-F1, balanced accuracy

- سطح پیکسل: Pixel AUROC, best Dice

**۴) توليد گزارش‌ها :**

- جدول نتایج per-category و خلاصه کلی (macro mean)

**معيار رقابت و امتياز اضافي (Leaderboard)**

تيم‌ها بر اساس عملكردشان در سه معيار زير، رتبه‌بندی می‌شوند و يك امتيازی از هر معیار می‌گيرند. امتياز نهايی هر تيم از جمع امتيازش در هر معیار به دست می‌آيد و در نهايیت تيم‌ها بر اساس جمع امتيازشان رتبه‌بندی نهايی می‌شوند.

OverallScore	PixelScore	ImageScore	رتبه
۳+	۳+	۳+	اول
۲+	۲+	۲+	دوم
۱+	۱+	۱+	سوم

## پروژه شماره ۹ (۱۱۵ نمره)

### ۱. مقدمه و هدف پروژه

دانشجویان دانشگاه صنعتی شریف در طول دوران تحصیل خود با پرسش‌های متعددی درباره قوانین و مقررات آموزشی مواجه هستند. این مقررات به صورت رسمی در آیین‌نامه‌ها و دستورالعمل‌های منتشرشده توسط مدیریت امور آموزشی دانشگاه صنعتی شریف در دسترس قرار دارد، اما جست‌وجو و استخراج پاسخ دقیق از این متون معمولاً دشوار و زمان‌بر است.

هدف این پروژه، طراحی و پیاده‌سازی یک چتبات هوشمند مبتنی بر Retrieval-Augmented Generation (RAG) است که بتواند صرفاً بر اساس آیین‌نامه‌ها و مقررات رسمی دانشگاه صنعتی شریف به پرسش‌های آموزشی دانشجویان پاسخ دهد و از تولید اطلاعات نادرست یا خارج از متن آیین‌نامه جلوگیری کند.

### ۲. تعریف مسئله

در این پروژه، هر تیم موظف است یک چتبات راهنمای آموزشی طراحی کند که بتواند به پرسش‌هایی در حوزه‌های زیر پاسخ دهد:

- قوانین آموزشی دورهٔ کارشناسی؛
- آیین‌نامه امتحانات و مقررات غیبت؛
- کارآموزی، پروژهٔ کارشناسی و کار و آموزش پایدار (کوآپ)؛
- مهمانی، انتقال و معرفی به استاد؛
- معافیت تحصیلی و مهلت فراغت از تحصیل؛
- و سایر موضوعات مندرج در آیین‌نامه‌های رسمی آموزشی.

چتبات باید  **فقط** بر اساس آیین‌نامه‌ها و مقررات منتشرشده در وب‌سایت رسمی مدیریت امور آموزشی دانشگاه صنعتی شریف به آدرس زیر پاسخ دهد:

<https://ac.sharif.edu/rules/>

### الزامات پاسخ‌دهی پاسخ‌های چتبات باید دارای ویژگی‌های زیر باشند:

- عدم تولید قانون یا تفسیر خارج از متن آیین‌نامه؛
- ارجاع صریح به نام آیین‌نامه و در صورت امکان شماره ماده، بند یا تبصره؛
- اعلام شفاف عدم وجود پاسخ در آیین‌نامه با عباراتی مانند:
  - «در آیین‌نامه موجود مطلبی در این‌باره یافت نشد»
  - «نیاز به استعلام از آموزش دارد»

### ۳. معماری مورد انتظار سامانه

معماری مورد انتظار یک RAG تک‌مرحله‌ای (Single-hop) است که شامل مراحل زیر می‌باشد:

۱. تبدیل سؤال کاربر به embedding
۲. بازیابی بخش‌های مرتبط از آیین‌نامه‌ها؛
۳. ارسال متن بازیابی‌شده به همراه سؤال به مدل زبانی؛

۴. تولید پاسخ کوتاه، دقیق و مستند.

#### ۴. پایگاه دانش (Knowledge Base)

##### ۱۰۴. منبع داده

منبع داده رسمی و واحد برای تمام تیم‌ها، صفحه‌آیین‌نامه‌ها و مقررات مدیریت امور آموزشی دانشگاه صنعتی شریف است. این صفحه شامل آیین‌نامه‌هایی از جمله موارد زیر می‌باشد:

- آداب‌نامه استفاده از ابزارهای هوش مصنوعی؛
- آیین‌نامه آموزشی دوره کارشناسی (نسخه جاری و نسخه‌های ورودی‌های قبل)؛
- آیین‌نامه برگزاری و غیبت در امتحانات؛
- آیین‌نامه روابط پیش‌نیازی و همنیازی؛
- آیین‌نامه کارآموزی؛
- دستورالعمل پروژه کارشناسی؛
- شیوه‌نامه مهمانی و انتقال؛
- قوانین تغییر رشته؛
- مقررات معافیت تحصیلی؛
- نظام‌نامه آموزشی.

#### ۲۰۴. یکنواخت‌سازی داده

برای تضمین عدالت در ارزیابی:

- یک نسخه آفلاین (HTML) از تمام آیین‌نامه‌ها توسط تیم تدریس تهیه و در اختیار همه تیم‌ها قرار می‌گیرد؛
- استفاده از هر منبع دیگری خارج از این مجموعه مجاز نیست؛
- پایگاه دانش هر تیم باید صرفاً از این فایل‌ها ساخته شود.

#### ۳۰۴. ساخت پایگاه دانش

هر تیم باید مراحل زیر را انجام دهد:

۱. استخراج متن از فایل‌ها؛

۲. تمیزسازی حداقلی (حذف هدر، فوتر، شماره صفحه و تکرارها)؛

۳. تقسیم متن به بخش‌های کوچک (Chunking):

- ترجیحاً بر اساس ماده، بند یا تبصره؛
- یا پنجره‌های طول ثابت (حدود ۲۰۰ تا ۳۰۰ توکن فارسی).

هر chunk باید شامل متادیتای زیر باشد:

- متن؛
- نام آیین‌نامه؛

- شماره ماده یا بخش (در صورت امکان):
- تاریخ آئین نامه (اختیاری):
- شناسه یکتا.

## **Retrieval و Embedding .۵ Embedding .۱۰۵**

- استفاده از یک مدل embedding مناسب برای زبان فارسی یا چندزبانه:
- محاسبه embedding برای تمام ها:
- ذخیره در یک index (مانند FAISS یا ابزارهای مشابه).

## **Retrieval .۲۰۵ برای هر سؤال کاربر:**

1. سؤال محاسبه می شود:
  2.  $k$  بخش مرتبط (مثلاً بین ۵ تا ۲۰) بازیابی می شود:
  3. بخش های بازیابی شده برای تحلیل ذخیره و نمایش داده می شوند.
۶. مدل زبانی

- مدل: Qwen 2.5 – 7B (Instruct/Chat)
- اجرا روی Google Colab:
- استفاده از quantization (۴ یا ۸ بیت):
- بدون prompt engineering و صرفاً با fine-tuning

## **۷. رابط کاربری پیاده سازی رابط کاربری آزاد است و می تواند شامل موارد زیر باشد:**

- Notebook تعاملی (Jupyter)
- رابط خط فرمان (CLI)
- رابط وب سبک (Gardio, Streamlit)

## **۸. ارزیابی ۱۰.۸. مجموعه سؤالات مشترک**

یک مجموعه ثابت شامل حدود ۱۵ تا ۲۰ سؤال به تمام تیمها داده می شود. هر تیم موظف است این سؤالها را از طریق کد به سیستم داده و خروجی را در قالب JSON یا CSV تحویل دهد.

## **۲۰.۸ ۰.۸. معیارهای نمره دهنده پاسخ ها برای هر سؤال:**

- درستی محتوایی (۰ تا ۲ نمره)
- وفاداری به منبع (۰ یا ۱ نمره)

## ۹. نمره‌دهی کل پروژه

- ساخت پایگاه دانش: ۲۰ نمره
- پیاده‌سازی RAG: ۴۵ نمره
- کیفیت پاسخ‌ها روی مجموعه ارزیابی: ۳۰ نمره
- گزارش نهایی و دمو: ۲۰ نمره

-