



دانشگاه تربیت مدرس

دانشکده مهندسی برق و کامپیوتر

تمرین پنجم درس یادگیری ماشین

ماشین‌های بردار پشتیبان و یادگیری تجمیعی

استاد محترم درس:

جناب آقای دکتر قادری

دستیاران آموزشی:

محمد مولائی

نیلوفر مقدس

- ✓ پروژه فقط با زبان برنامه نویسی پایتون قابل قبول میباشد.
- ✓ فایل تحویلی شما، یک فایل زیپ شده نهایی شامل گزارش کار (فایل pdf) و فایل کد پایتون با پسوند ipynb (Jupyter Notebook) می باشد. لطفا آن را به صورت زیر نام گذاری و ارسال نمایید.

✓

HW_۵_HMM_[LastName]_[FirstName]

- ✓ گزارش کار خود را در یک فایل pdf تحویل دهید و از گذاشتن صرفا اسکرین شات های پشت سرهم از کد در گزارش کار خودداری کنید.
- ✓ توجه داشته باشید که در فایل ارسالی پایتون، خروجی هر سلول (شامل نمودار، خروجی عددی و غیره) حتما ذخیره شده و قابل مشاهده باشد.
- ✓ لازم است حتما نتایج بدست آمده را گزارش و تحلیل کنید.
- ✓ علاوه بر مهارت حل سوالات، نوشتن پاسخ مینی پروژه ها در فرمت گزارش فنی (فصل بندی و صفحه بندی مناسب، رعایت اصول نگارش و ...) برای دانشجویان تحصیلات تکمیلی اهمیت دارد، این مورد نیز در ارزشیابی لحاظ می شود.
- ✓ در صورت فراموشی در ارسال کد پایتون، هیچ نمره ای به شما تعلق نخواهد گرفت.
- ✓ در صورت مشاهده تشابه در هر بخش از انجام پروژه، نمره هر دو نفر صفر لحاظ میگردد.

برای پاسخ به سوالات دانشجویان در مورد مینی پروژه ها، دو راه ارتباطی وجود دارد:

1. ایمیل: برای پرسش سوال از طریق ایمیل، در قسمت To، ایمیل دستیار آموزشی این مینی پروژه، آقای مولائی (mhmdmovlaie@gmail.com) را قرار دهید و در قسمت CC، ایمیل [niloofar.moghaddas\[at\]gmail.com](mailto:niloofar.moghaddas[at]gmail.com) را قرار دهید.
2. گروه تلگرام: می توانید برای پرسش سوال از مینی پروژه چهارم در گروه، آقای مولائی را با شناسه تلگرامی @mhmdmolaee، در پیام خود نام ببرید.

ایمیل دستیار آموزشی: mhmdmovlaie@gmail.com

هدف از انجام این مینی پروژه اعمال و بررسی عملکرد الگوریتم‌های SVM و Ensemble Learning بر روی مجموعه داده‌ی معرفی شده است.

در این تمرین یک فایل notebook با خروجی هر سلول در اختیار شما قرار داده شده است و شما می‌بایست تلاش کنید تا خروجی هر سلول مجدداً تولید کنید (نتایج مدل‌ها الزامی نیست).

- مجموعه داده Bank Marketing

این مجموعه داده برای کمپین بازاریابی یک موسسه‌ی بانکی پرتغالی است. هدف از طبقه‌بندی پیش‌بینی ثبت یا عدم ثبت سپرده (deposit) در این بانک است. هر یک از ویژگی‌های این مجموعه داده در این [آدرس](#) توضیح داده شده است.

- معیار ارزیابی

با توجه به مجموعه داده از معیارهای Accuracy, Precision, Recall و F_1 _score استفاده می‌کنیم.

پروژه پنجم: ماشین هسته و یادگیری تجمعی

- گام اول EDA:

در این بخش مجموعه داده را از دید آماری مورد بررسی قرار دهید و با چالش های مجموعه داده آشنا شوید. مواردی که می بایست بررسی کنید:

(الف) آیا در مجموعه داده مقادیر گمشده وجود دارد؟ این مورد منجر به چه مشکلاتی می شود؟ چگونه می توان با آن مقابله کرد؟

(ب) توزیع کلاس خروجی به چه صورت است؟

(ج) همبستگی ویژگی ها با یک دیگر به چه صورت است؟ همبستگی میان دو ویژگی از رابطه ی زیر محاسبه می شود.

$$\rho_{X,Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

این معیار نشان می دهد که آیا ارتباط خطی میان دو ویژگی وجود دارد یا خیر. آیا در صورتی که این معیار بین یکی از ویژگی ها و کلاس خروجی نزدیک به صفر بود، می توان آن را حذف کرد؟ چرا؟

(د) آیا در مجموعه داده، داده ی پرت وجود دارد؟ این مورد منجر به چه مشکلاتی می شود؟ چگونه می توان با آن مقابله کرد؟

(ه) پس از کاهش ابعاد، با رسم داده ها به چه نتیجه ای می توان رسید؟ کدام روش کاهش ابعاد مناسب تر به نظر می رسد؟ چرا؟

در نظر داشته باشید که انجام این مراحل الزامی است و در صورت تمایل می توانید دیدهای بیشتر و متفاوت تری مجموعه داده را بررسی کنید.

- گام دوم: پیش پردازش

در این مرحله، بعد از آشنایی با مجموعه داده، پیش پردازش هایی که از نظرتان الزامی هستند را اعمال کنید. انجام مراحل ذکر شده الزامی است.

(الف) تقسیم مجموعه داده به بخش های train و test.

(ب) normalize کردن داده های عددی

- گام سوم: یادگیری و انتخاب مدل

در این بخش در ابتدا الگوریتم SVM و یک مدل یادگیری تجمعی مانند Random Forest و XGBoost را بر روی داده های پیش پردازش شده اعمال کنید. نتایج هر مدل را با مشخص کردن پیش پردازش اعمال شده، در یک جدول نمایش دهید (نتایج می بایست با استفاده از Cross-validation تعیین شود). سپس چند مدل برتر از میان این مدل ها انتخاب کنید و با استفاده از روش هایی مانند رای اکثریت یک مدل واحد بسازید. نتیجه این مدل را نیز به انتهای جدول اضافه کنید. بهترین

مدل را با توجه به نتایج انتخاب کنید. توجه داشته باشید که استفاده از مدل‌های نام برده شده الزامی است ولی محدود به این مدل‌ها نیستید.

۱ مدل، پیش‌پردازش و نتایج بدست آمده

الگوریتم	پیش‌پردازش	میانگین \pm انحراف معیار Accuracy	میانگین \pm انحراف معیار Precision	میانگین \pm انحراف معیار Recall	میانگین \pm انحراف معیار F1_score
SVM (rbf)	Missing Values (Method), PCA, etc.	$70 \pm 0,5$
جنگل تصادفی	Forward Feature Selection

- گام چهارم: ارزیابی مدل

توضیح دهید که بهترین مدل را چگونه انتخاب می‌کنید و آن را بر روی مجموعه تست اعمال کنید. توجه کنید که داده‌های تست را فقط و فقط یک بار می‌توانید به مدل بدهید، در غیر این صورت نتایج معتبر نخواهند بود.

- گام پنجم: بحث و نتیجه‌گیری

نتایج بدست آمده از هر بخش را به صورت مجزا گزارش کنید و علت هر کدام از پیش‌پردازش‌هایی که اعمال کرده‌اید را شرح دهید. همچنین توضیح دهید، به نظر تان چرا بهترین نتیجه را از آن الگوریتم بدست آورده‌اید؟