

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

تمرین سری سوم یادگیری ماشین

دانشکده مهندسی کامپیوتر

استاد درس: دکتر ناظر فرد

اردیبهشت ۱۴۰۰

- تمامی مستندات شامل گزارش به همراه کدهای خود را در یک فایل فشرده با فرمت zip ذخیره کرده و با عنوان stdNum_HW2.zip که stdNum شماره دانشجویی شما است در سامانه بارگزاری کنید.
- سوالات ستاره‌دار(*) نمره اضافی داشته و انجام آن‌ها اجباری نمی‌باشد.
- مهلت انجام تمرین تا ساعت ۲۳:۵۵ روز شنبه ۸ خرداد می‌باشد.

سوال‌های تشریحی (ماشین‌های بردار پشتیبان و مدل‌های ترکیبی)

سوال ۱) صحت هر یک از موارد زیر را بررسی کرده و دلایل خود را توضیح دهید.

الف) ماشین‌های بردار پشتیبان^۱ پارامتریک‌اند.

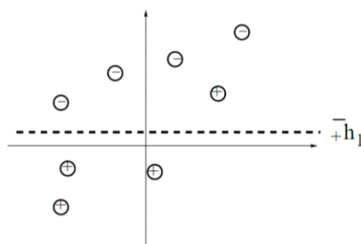
ب) مقدار حاشیه بدست آمده برای دو ماشین بردار پشتیبان با کرنل‌های متفاوت که برای داده‌های یکسان آموزش دیده‌اند، می‌تواند معیاری برای میزان کارایی مدل باشد.

ج) ماشین‌های بردار پشتیبان همواره در برابر بیش‌برازش مقاوم می‌باشند.

د) وجود داده‌های پرت و نویز بر روی ماشین‌های بردار پشتیبان بی‌تاثیر است.

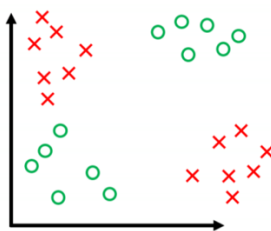
سوال ۲) برای داده‌های زیر الگوریتم آدابوست را اجرا کرده و در تکرار اول مرز تصمیم h_1 بدست آمده است. وزن α_1 که توسط

الگوریتم آدابوست به h_1 اختصاص داده می‌شود چقدر است؟ (وزن‌های اولیه برابر $\frac{1}{8}$ می‌باشند).



شکل ۱

سوال ۳) امکان استفاده از ماشین‌های بردار پشتیبان برای مسئله زیر را مورد بررسی قرار دهید. در صورت امکان راهکار خود را بطور کامل با ذکر مقادیر و جزییات شرح دهید.



شکل ۲

^۱ Support Vector Machine (SVM)

سوال ۴) سوالات زیر را در مورد روش های ترکیبی پاسخ دهید.

الف) آیا الگوریتم آدابوست در برابر داده ی نویز حساس است؟ چرا؟

ب) الگوریتم آدابوست با استفاده از هر نوع دسته بند ضعیف و یا ترکیب چند دسته بند ضعیف در نهایت به خطای صفر می رسد.

- کدهای خود را به زبان پایتون و ترجیحا در محیط jupyter پیاده سازی کنید. می توانید تحلیل خودتان را به عنوان سلول های متنی در همان محیط ارائه کنید.
- نظم در نوشتن گزارش و کدها می تواند به کسب نمره ی بهتر به شما کمک کند. برنامه نوشته شده خوانا و کامنت گذاری مناسب داشته باشد.
- در پیاده سازی بخش های مختلف، امکان استفاده از کتابخانه های آماده مرتبط با الگوریتم های یادگیری ماشین را به طور کلی ندارید. موارد مجاز در صورت سوال بخش ها ذکر شده است.
- برای نمایش نمودارها و عملیات ماتریسی می توانید از کتابخانه های numpy و matplotlib استفاده کنید.
- همچنین برای خواندن داده ها به عنوان ورودی می توانید از pandas استفاده کنید.
- برای محاسبه معیارهای ارزیابی مانند دقت و ماتریس درهم ریختگی می توانید از کتابخانه آماده استفاده کنید.
- مطابق قوانین دانشگاه هرگونه کپی برداری ممنوع می باشد و در صورت مشاهده نمره هر دو طرف صفر در نظر گرفته می شود.
- در صورت داشتن سوال می توانید با ایمیل تدریس یاران درس در تماس باشید:

hse.khalilian08@gmail.com , hamid.dargahi0072@gmail.com

سوال های پیاده سازی

مسئله ۱) پیاده سازی ماشین بردار پشتیبان

دیتاست **lsvt-voice-rehabilitation** را از ادرس زیر دانلود کرده و به سوالات پاسخ دهید. برای ارزیابی از -fold-cross- ۱۰

validation استفاده کنید. در این تمرین مجاز به استفاده از کتابخانه آماده هستید.

<https://archive.ics.uci.edu/ml/datasets/LSVT+Voice+Rehabilitation>

الف) داده ها را به روش های زیر دسته بندی کنید سپس به سوالات که در زیر آمده است پاسخ دهید.

- کرنل خطی
- کرنل چند جمله ای (پارامترهای d, r)

- کرنل RBF (پارامتر گاما)

- سیگموئید (پارامتر r)

(ب) معیار دقت^۲ و F1-measure را برای هر یک از دسته‌بندی‌های بالا بدست آورید (برای هر یک از پارامترهای یاد شده

حداقل سه مقدار متفاوت در نظر بگیرید).

(ذ) تاثیر پارامترهای هر کرنل بر کارایی مدل‌ها را تحلیل کنید.

(د) بهترین مدلی را که یافتید، مشخص نمایید.

مسئله ۲) بررسی عملکرد مدل‌های ترکیبی

توضیح: داده‌های مورد استفاده در این تمرین مربوط به تصاویر اعداد دست نوشته انگلیسی است. از هر تصویر، ۱۶ ویژگی استخراج شده که این ویژگی‌ها به همراه برچسب تصویر متناظر، در اختیار شما قرار داده شده است. هدف، طراحی دسته‌بند مناسب برای این داده‌ها است. لازم به ذکر است که داده‌های یادگیری و تست به صورت مجزا و به ترتیب در فایل data_train.csv و data_test.csv قرار دارد.

(الف) با استفاده از روش Random Forest یک دسته‌بند مناسب طراحی کنید. دسته‌بندی را با استفاده از ترکیب نتیجه‌ی ۱۵ درخت تصمیم انجام دهید. هر درخت تصمیم را نیز با استفاده از ۳ ویژگی که به صورت تصادفی انتخاب می‌شود، آموزش دهید. برای آموزش درخت‌های تصمیم می‌توانید از کتابخانه‌ی scikit-learn استفاده کنید. بیشینه‌ی عمق درخت‌ها را نیز ۳ در نظر بگیرید. دقت دسته‌بندی و همچنین ماتریس درهم‌ریختگی را برای داده‌های تست در گزارش خود ارائه کنید.

(ب) با استفاده از روش AdaBoost یک دسته‌بند مناسب طراحی کنید. دسته‌بندی را با استفاده از ترکیب نتیجه‌ی ۱۰ درخت تصمیم انجام دهید. برای آموزش درخت‌های تصمیم همچون قبل می‌توانید از کتابخانه‌ی scikit-learn استفاده کنید. دقت دسته‌بندی را برای داده‌های تست در گزارش خود ارائه کنید.

(ج) با استفاده از روش AdaBoost قسمت قبل را با ۵، ۲۰، و ۵۰ دسته‌بند درخت تصمیم تکرار کرده و برای هر کدام دقت دسته‌بندی را برای داده‌های تست در گزارش خود ارائه کنید.

***د)** در چند سال اخیر روش‌های Gradient Boosting بسیار مورد توجه قرار گرفته‌اند. یکی از کتابخانه‌هایی که برای این روش‌ها منتشر شده است، کتابخانه‌ی XGBoost است. با بررسی این کتابخانه، یک دسته‌بند مناسب برای داده‌های این تمرین طراحی کنید. انتخاب پارامترهای بهینه در این قسمت به عهده شما است. پس از پایان یادگیری، دقت دسته‌بندی را بر روی داده‌های تست برای دسته‌بند پیشنهادی خود در گزارش ارائه کنید.

^۲ Accuracy