

به نام او ...

تمرین دوم درس داده کاوی

در این تمرین دو مجموعه داده برای حل دو مسئله‌ی رگرسیون در اختیار شما قرار داده می‌شود. هر دو بخش تمرین به صورت رقابتی بوده و نتایج بهتر نمره‌ی بیشتری در پی خواهند داشت.

بخش اول (۵۰ نمره)

مجموعه داده‌ی این قسمت از نمره‌های تعدادی دانشجو در ۷ درس مختلف تشکیل شده است. هدف آموزش مدلی با استفاده از ۶ نمره‌ی اول برای پیش‌بینی نمره‌ی هفتم می‌باشد. مراحل زیر را برای این بخش طی نمایید.

۱- ابتدا تعریف کوتاهی برای عبارات زیر ارائه کنید.

- K-fold Cross-Validation
- MSE, MAE, RMSE
- Covariance, Correlation
- Regression toward the mean
- L1 norm, Lasso Regression
- L2 norm, Ridge Regression

۲- نمودار هر یک از ۶ ویژگی نسبت به ویژگی هدف را رسم کرده و با در نظر گرفتن مفهوم Correlation در مورد هر یک از نمودارها به صورت شهودی بحث کنید.

۳- به جای مقادیر نامشخص مجموعه داده، مقدار ۰ قرار داده شده است. راهکاری برای پر کردن مقادیر نامشخص ارائه دهید. هرچه قدر این راه حل هوشمندانه‌تر باشد نمره‌ی بیشتری تعلق خواهد گرفت.

۴- از روش‌های زیر و یک روش دلخواه دیگر برای آموزش مدل و پیش‌بینی ویژگی هدف استفاده کنید (نیازی به پیاده‌سازی روش‌ها به طور مستقیم نیست و می‌توانید از کتابخانه‌هایی مانند Scikit استفاده کنید). خطای هر کدام را با 10-fold Cross-Validation با در نظر گرفتن معیار RMSE گزارش کنید. توجه داشته باشید که تغییر هاپر پارامترهای هر روش می‌تواند تاثیر قابل توجهی در خطا داشته باشد.

- Linear Regression
- Lasso
- Gradient Boosting

برای روش دلخواه دیگر می‌توانید از رگرسیون با استفاده از توابع چندجمله‌ای، رگرسیون با توابع سینوسی و کسینوسی، یا رگرسیون با بسط‌های نمایی استفاده کنید. همچنین می‌توانید از Regularization برای تعیین تنظیم پیچیدگی مدل استفاده نمایید. مقادیر ویژگی هفتم را برای مجموعه داده‌ی آزمایشی پیش‌بینی کرده و به

همراه ۶ ویژگی دیگر در یک فایل csv به اسم P1_submission.csv ذخیره کنید. کد نهایی باید با خواندن فایل‌های train.csv و test.csv از پوشه ی ۱، فایل خروجی را تولید کند. از هر روش دلخواهی برای پیش‌بینی مقادیر ویژگی هدف می‌توانید استفاده کنید. در نظر داشته باشید که خطای کم‌تر معادل با نمره‌ی بیش‌تر خواهد بود. بنابراین از ایده‌های خود برای کاهش خطا استفاده کنید. در نظر داشته‌باشید که می‌توانید زیرمجموعه‌ای از ویژگی‌ها را برای آموزش مدل استفاده کرده یا ویژگی جدیدی از روی ویژگی‌های موجود بسازید.

بخش دوم (۵۰ نمره)

در این قسمت یکی از مسابقه‌های سایت Kaggle در مورد پیش‌بینی قیمت خانه در نظر گرفته شده است. از این [لینک](#) می‌توانید به صفحه‌ی مسابقه دسترسی داشته‌باشید. محدودیت خاصی برای این قسمت اعمال نمی‌شود. البته انتظار می‌رود حداقل‌های زیر در حل مسئله بررسی شود.

- حل مشکل وجود مقادیر نامشخص
- استفاده از 10-fold cross-validation با معیار RMSE برای گزارش خطای مدل‌های مختلف
- استفاده از حداقل ۳ مدل مختلف برای یافتن بهترین روش
- بررسی زیرمجموعه‌های مختلف از ویژگی‌ها برای آموزش مدل

کد نهایی باید با خواندن فایل‌های train.csv و test.csv از پوشه‌ی ۲، مدل انتخاب شده‌ی نهایی را آموزش داده و پیش‌بینی قیمت خانه‌های مجموعه داده‌ی test.csv را در فایل P2_submission.csv با توجه به ساختار فایل sample_submission.csv موجود در پوشه‌ی ۲ تولید کند.

گزارش:

گزارش بایستی در قالب فایل PDF باشد. لطفاً فایل Word نفرستید. در گزارش تحلیل خود را در رابطه با تمام کارهایی که انجام داده‌اید بیان نمایید. بخش مهمی از کار مربوط به فرآیند رسیدن به نتیجه‌ی نهایی می‌باشد. بنابراین گزارش بخش قابل توجهی از نمره را به خود اختصاص می‌دهد.

فایل گزارش خود را به شکل «Report2_StdNum.pdf» نامگذاری کنید. (مانند Report2_9131081.pdf)

کد:

کد اجرایی می‌تواند در محیط R یا Python تهیه شود. عدم وجود کد معادل نمره‌ی صفر خواهد بود.

فایل کد خود را به شکل «DM2_P1_StdNum» برای بخش اول و «DM2_P2_StdNum» برای بخش دوم نامگذاری کنید. علاوه بر این مدل نهایی مورد استفاده برای پیش بینی ویژگی هدف را در کدهای جداگانه به نام‌های «DM2_P1_Submit_StdNum» و «DM2_P1_Submit_StdNum» قرار دهید.

بارگذاری:

تمام فایل های مورد نظر را در قالب یک فایل فشرده در سایت درس بارگذاری نمایید.

فایل فشرده را به شکل «DM2_StdNum» نامگذاری کنید. (مانند DM2_9131081)

مهلت ارسال تمرین ساعت ۲۳:۵۵ دقیقه‌ی روز دوشنبه مورخ ۱۶ اسفند می‌باشد.

به ازای هر روز تاخیر در ارسال تمرین، برای دو روز اول ۵ و روز های بعد ۱۰ درصد از نمره‌ی آن از دست خواهد رفت.

هر گونه سوال در مورد تمرین را می‌توانید از طریق ایمیل AUT.DM2017@gmail.com بپرسید.

موفق باشید