



پردیس دانشکده های فنی دانشگاه تهران

دانشکده برق و کامپیوتر



هوش مصنوعی

پروژه چهارم - یادگیری ماشین

مهلت تحویل: ۱۹ آذر

طراحان: مهسا اسکندری - آنیلا قندهاری - دیار محمدی

هدف پروژه:

هدف این پروژه آشنایی با روشهای یادگیری ماشین¹ به کمک کتابخانه SciKit-Learn است. این پروژه چهار در فاز تعریف شده است. در فاز صفر، به بررسی مجموعه داده ها و تجزیه و تحلیل داده های اکتشافی² میپردازید. در فاز اول و پیش پردازش³ را خواهید داشت. در فاز دوم، با استفاده از چند تخمین زن⁴ تعریف شده در کتابخانه SciKit-Learn مدلهایی را پیاده سازی و بهینه سازی خواهید کرد. نهایتاً در فاز سوم، با استفاده از مدلهای بهینه فاز دوم، به پیاده سازی چند روش یادگیری گروهی⁵ و تحلیل نتایج حاصل میپردازید.

- در طول این متن با علامت [Number] مواجه می شوید که به معنای شماره ی لینک مربوط به آن موضوع می باشد. لینک ها در انتها آمده اند.
- در طی این پروژه نیاز است که اثر تصمیمات خود را در نتیجه ی نهایی بررسی کنید. لذا بهتر است به نحوی کد بنویسید که اعمال تغییر در هر مرحله ی آن و بررسی نتیجه آسان باشد.
- در این پروژه تحلیل شما از تک تک مراحل کار، بیان مزایا و معایب هر روش و به طور کلی توانایی تشخیص استفاده از ابزار مناسب در هر مرحله از این پروژه از اهمیت بسیاری برخوردار است و شامل 80 درصد نمره پروژه خواهد بود. بنابراین در این پروژه کد شما تنها بخش محدودی از نمره را شامل می شود و بخش عمده نمره به توضیحات شما در حین ارائه و همچنین تحلیل های شما از هر بخش و ارائه نمودارهای مناسب در گزارش کار اختصاص دارد.

¹ Machine Learning

² Exploratory Data Analysis

³ Preprocessing

⁴ Regressor

⁵ Ensemble Learning

معرفی مجموعه داده⁶:

مجموعه داده‌ای که در اختیار شما قرار دارد [House Prices](#) است که یک مجموعه داده بارگذاری شده در سایت [Kaggle](#) است. این مجموعه داده شامل ویژگی⁷ های خانه های مختلف شامل خیابان، محله، خدمات، تعداد اتاق ها، نوع ساختمان، سال ساخت و ... می باشد و داده هدف⁸، قیمت خانه است.

این مجموعه داده شامل ۴ فایل است که به ترتیب توضیح داده خواهند شد.

train.csv:

دارای ۸۳ ستون است که یکی از آن‌ها قیمت فروش هر خانه است و مابقی ویژگی‌های دیگر هر خانه‌ی فروخته شده است. این داده‌ها برای فرآیند یادگیری استفاده می‌شوند.

test.csv:

دارای ۸۲ ستون است و کاملاً مشابه train.csv است با این تفاوت که ستون قیمت فروش خانه‌ها در آن وجود ندارد.

data_description.txt:

توضیحات کامل هر ستون مانند نام ستون، مفهوم آن و مقادیر آن است. توصیه می‌کنیم برای آشنایی با داده‌ها این فایل را حتماً ملاحظه فرمایید.

sample_submission.csv:

یک نمونه فایل خروجی برای شرکت در مسابقه‌ی مربوط به مجموعه داده در سایت Kaggle است.

⁶ Dataset

⁷ Feature

⁸ Target

فاز صفر: Visualization and EDA

در این فاز داده های خام ورودی را بررسی خواهید کرد. در آمار، EDA رویکردی برای تجزیه و تحلیل مجموعه داده ها برای جمع بندی خصوصیات اصلی آنها است که اغلب با روش های Visualization انجام می شود. چه از مدل های آماری استفاده بشود چه نشود در درجه اول EDA برای دیدن آنچه داده ها می توانند فراتر از وظیفه مدلسازی به ما بگویند، انجام می شود.

در هر مرحله نتایجی را که بدست می آورید در گزارشتان ذکر کنید

1. ساختار کلی داده ها را با متدهای describe و info از پکیج Pandas بررسی کنید.
2. درصد داده های از دست رفته ی هر ویژگی⁹ را بدست آورید.
3. نمودار وابستگی ویژگی ها به هم را رسم کنید (می توانید از برای این کار می توانید از Heatmap استفاده کنید). به نظرتان کدام ویژگی ها در تخمین قیمت خانه ها می توانند مفید باشند؟ حداقل چهار ویژگی را نام ببرید.
4. از قیمت خانه ها لگاریتم گرفته و دوباره نمودار وابستگی ویژگی ها را رسم کنید. آیا تفاوت وجود دارد؟ چرا؟
5. آیا سوال ۳ برای بررسی تمام ویژگی ها کافی است؟ چه ایرادی وجود دارد؟
6. ارتباط ویژگی هایی که در سوال سه نام بردید را با قیمت خانه ها دقیق تر بررسی کنید. می توانید از نمودارهای Hexbin و Scatter استفاده کنید.
7. نگاهی به فایل data_description.txt انداخته و سه ویژگی دسته ای که به نظرتان در پیش بینی قیمت مفید هستند را انتخاب کنید سپس با نمودارهای مناسب آن ها را نمایش دهید.
8. شما می توانید هر بررسی دیگری به هر طریقی که می خواهید بر روی داده ها داشته باشید تا دید مناسب تری نسبت به ویژگی ها داشته باشید و در فازهای بعدی بهتر عمل کنید.

⁹ Missing Values

فاز اول: Preprocessing

پیش پردازش داده ها [1] مرحله مهمی در فرآیند داده کاوی است. روش های جمع آوری اطلاعات اغلب به راحتی کنترل نمی شوند، و در نتیجه مقادیر خارج از محدوده (به عنوان مثال، تعداد فرزندان: 100)، ترکیب داده های "غیر ممکن" (به عنوان مثال، جنسیت: مرد، بارداری: بله)، مقادیر از دست رفته و غیره ممکن است در مجموعه داده ها وجود داشته باشند و می توانند نتایج گمراه کننده ای ایجاد کنند. بنابراین، قبل از هر کاری باید اینگونه داده ها را اصلاح کنیم. غالباً، پیش پردازش داده ها مهمترین مرحله پروژه یادگیری ماشین است، پس بهتر است کاملاً روی این موضوع مسلط باشید.

برای هر کدام از موارد زیر ممکن است روش های مختلفی وجود داشته باشد، تفاوت روش ها را در گزارش توضیح داده و دلیل استفاده از روشی که انتخاب کردید را ذکر کنید.

1. دو روش برای برطرف کردن مشکل داده های گمشده، ۱- حذف کل ستون و ۲- پر کردن مقادیر گمشده با آماره ها (مثلاً میانگین) هستند. دیگر روش های موجود برای هندل کردن داده های گمشده را در صورت وجود مختصراً توضیح دهید و این روش ها را مقایسه کنید.
2. با توجه به نتایج سوال دوم فاز صفر، کدام ویژگی ها بیشترین میزان داده ی گمشده را دارند؟ آیا باید ویژگی های با درصد بالای داده ی گمشده را از مجموعه داده ها حذف کرد؟ با روش های مناسب داده های گم شده را هندل کنید.
3. در ویژگی های عددی Normalizing یا Standardizing داده های عددی به چه منظور انجام می شود؟ به نظر شما آیا نیاز است که داده ها خود را در این پروژه نرمالیزه کنیم؟
4. برای این که مدل ما بتواند با داده های دسته ای کار کند چه روش هایی وجود دارد؟ آیا همه ی داده های دسته ای را باید با یک روش هندل کنیم؟ [2]
5. خوب است این جا فکر کنید که آیا باید همه ی ستون ها را نگه داریم؟ کدام ستون ها را می توانیم حذف کنیم؟ چرا؟ (در پاسخ گویی به این پرسش سوال سوم از فاز صفر نیز می تواند مفید باشد.)
6. نیاز است که داده ها را به دو دسته ی train و test تقسیم کنیم. یک روش این است که P درصد اول داده ها را برای یادگیری و مابقی را برای ارزیابی اختصاص دهیم. روی مقدار مناسب P بحث کنید. آیا روش های دیگری وجود دارد؟ نام ببرید. آیا نیاز است تقسیم داده ها به صورت تصادفی باشد؟ چرا؟

فاز دوم: Model Training, Evaluation and Hyperparameter Tuning

در این فاز از پروژه، سه مدل بر پایه K-Nearest-Neighbours و Decision Tree و Linear Regression با کمک کتابخانه SciKit-Learn پیاده‌سازی می‌کنید. سپس باید به بهینه‌سازی هایپرپارامترهای هر مدل بپردازید. تعداد همسایه‌ها برای الگوریتم KNN و عمق ماکزیمم برای الگوریتم Decision Tree کفایت میکنند.

1. دقت هر مدل را بر اساس معیارهای MAE و RMSE روی داده های train و test اندازه‌گیری کنید.
2. در مدل‌های که نیاز به بهینه کردن هایپر پارامتر است نمودار تغییرات معیارهای بالا را به ازای مقادیر مختلف هایپر پارامتر رسم کنید. (در صورتی که تعداد هایپر پارامترها بالا برود برای سهولت کار می‌توانید از ابزار Grid Search در کتابخانه SciKit-Learn استفاده کنید.)
3. در مورد underfitting و overfitting تحقیق کنید و بررسی کنید که آیا در مدل‌های شما underfitting یا overfitting اتفاق افتاده است؟
4. تاثیر پیش‌پردازش‌هایی که روی داده‌ها انجام دادید را به طور کامل بررسی کنید. (مثلا تاثیر روش‌های مختلف هندل کردن مقادیر گمشده را روی معیارهای نهایی ببینید.)

فاز سوم: Ensemble Methods

یادگیری گروهی به این معناست که از تجميع نتایج حاصل از تعدادی مدل، پیش‌بینی نهایی را انجام دهیم. در این فاز به پیاده‌سازی و تحلیل چند روش یادگیری گروهی می‌پردازیم.

1. Random Forest روشی است که تعدادی Decision Tree با ویژگی‌های مختلف و دیتاهای مختلف را در کنار هم قرار داده و هر کدام جداگانه یادگیری را انجام می‌دهند و در نهایت میان آن‌ها به طریقی رای‌گیری می‌شود. با کمک کتابخانه SciKit-learn این مدل را پیاده‌سازی کنید. تاثیر حداقل دوتا از hyperparameter ها را بر مدل بررسی کرده و معیارهای دقت را محاسبه کنید.
2. یکی دیگر از روشهای یادگیری گروهی Voting Regression است. در این روش، پیش‌بینی نهایی، میانگینی از خروجی همه مدل‌ها است. با استفاده از سه مدل بهینه شده در فاز اول، یک مدل بر این پایه پیاده‌سازی کنید و معیارهای دقت را اندازه‌گیری کنید.
3. در مورد Voting Regression علت شکست یا برتری آن‌را نسبت به روش‌های یادگیری فردی بررسی کنید. (برای مثال می‌توانید تفاوت جواب‌های هر کدام از سه مدل فاز قبل را نسبت به جواب‌های دیگر مدل‌ها بررسی کنید.)

[Towardsdatascience.com: Overfitting vs Underfitting a Complete Example](https://towardsdatascience.com/overfitting-vs-underfitting-a-complete-example-400000000000)

نکات پایانی:

- مقدار مطلوب برای MSE و RMSE حداکثر به ترتیب \$20000 و \$30000 است. توجه داشته باشید که کسب این دقت، تنها توسط یکی از مدل‌هایی که در پروژه پیاده سازی کردید کافی است و سایر مدل‌ها باید صرفاً دقت معقولی داشته باشند.
- همچنین به شدت توصیه میکنیم برای مقایسه روش‌های مورد استفاده و نتایج خودتان با دیگران، نوتبوک‌های موجود در رابطه با این مجموعه داده را ببینید و حتماً در این [مسابقه](#) Kaggle شرکت کنید و نتایج به دست آمده برای فایل test.csv را در مسابقه آپلود کنید. (توجه داشته باشید که این مجموعه داده، به صورت کامل بر روی وب قرار دارد و تعدادی افراد با سوءاستفاده از این موضوع اقدام به شرکت در مسابقه و گرفتن امتیازهای عالی و میزان خطای نزدیک به صفر می‌کنند، نتایج خود را با چنین افرادی مقایسه نکنید. کسب خطای 0.20 و کمتر از آن در این مسابقه بسیار عالی است.)
- در تمامی بخش‌های پروژه، استفاده از کتابخانه SciKit-Learn مجاز است. دقت کنید که هدف پروژه تحلیل نتایج است بنابراین از ابزارهای تحلیل داده بطور مثال نمودارها استفاده کنید و توضیحات مربوط به هر بخش از پروژه را بطور خلاصه و در عین حال مفید در گزارش خود ذکر کنید.
- توجه کنید که باید از **Jupyter Notebook** استفاده کنید و نتایج و گزارش خود را در یک فایل فشرده با عنوان **AI_CA4_<#SID>.zip** تحویل دهید. محتویات پوشه باید شامل تمامی پیاده سازی‌ها در یک فایل Notebook به همراه خروجی html آن ارائه دهید.
- در صورتی که سوالی در مورد پروژه داشتید بهتر است در فروم درس یا گروه تلگرام مطرح کنید تا بقیه از آن استفاده کنند؛ در غیر این صورت توسط ایمیل با طراحان در ارتباط باشید.
- هدف از تمرین، یادگیری شماسست، مشورت ایرادی ندارد اما تمرین را باید خودتان انجام دهید.