

بخش اول: پرسش های تشریحی

- سوال (۱) درستی یا نادرستی عبارات زیر را با ذکر دلیل مشخص کنید.
- الف) در یک الگوریتم یادگیری، کاهش خطای آموزش، منجر به کاهش خطای آزمون می شود.
- ب) با کاهش تعداد داده های آموزش، الگوریتم بیشتر مستعد بیش برازش می شود.
- ج) افزایش پیچیدگی مدل در رگرسیون همیشه سبب کاهش خطای آموزش و افزایش خطای آزمون می شود.
- د) از میان معیارهای MSE و RMSE و MAE معیار MAE در برابر داده های پرت و نویزی عملکرد بهتری دارد.
- سوال (۲) فرض کنید که از رگرسیون چندجمله ای استفاده می کنید و با رسم منحنی ها متوجه می شوید که اختلاف زیادی بین خطای آموزشی و خطای اعتبارسنجی وجود دارد. توضیح دهید که چه اتفاقی رخ داده است و سه راه حل برای حل این مشکل ارائه دهید.
- سوال (۳) مجموعه داده آموزشی شامل n داده به فرم (x_i, y_i) در اختیار داریم که x_i دارای d بعد است. تابع هزینه SSE به صورت زیر محاسبه می شود:

$$J(w) = \sum_{i=1}^n (y_i - w^T x_i)^2$$

- الف) نشان دهید که در رگرسیون خطی با تابع هزینه SSE ، w بهینه به این صورت است:
- $$\hat{w} = (X^T X)^{-1} X^T y$$
- ب) دو مشکلی که استفاده مستقیم از این رابطه دارد را ذکر کنید و برای هر کدام راه حلی ارائه دهید.
- ج) اگر یک جمله منظم ساز تُرم ۲ به صورت $\|w\|^2$ به رابطه loss اضافه کنیم فرم بسته w بهینه را به دست آورید.
- د) توضیح دهید اضافه کردن جمله منظم ساز چه مزیت هایی نسبت به حالت عادی دارد.

سوال (۴) با در نظر گرفتن الگوریتم های زیر برای رگرسیون خطی، به سوالات پاسخ دهید:

- ۱- معادله ی نرمال (Normal Equation)
- ۲- روش های مبتنی بر گرادیان نزولی شامل:
 - Batch GD: گرادیان ها بر اساس کل داده آموزش محاسبه می شوند.
 - Stochastic GD: گرادیان ها بر اساس یک نمونه گیری از داده محاسبه می شوند.
 - Mini-Batch GD: گرادیان ها با استفاده از مجموعه های تصادفی کوچکی از نمونه ها به دست می آیند.

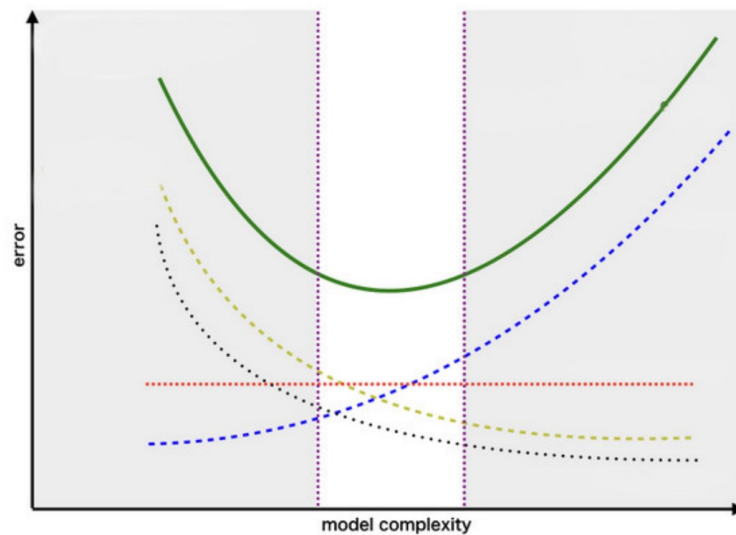
الف) اگر یک مجموعه داده آموزش با میلیون‌ها ویژگی آذر اختیار داشته باشید، از کدام الگوریتم رگرسیون خطی می‌توان استفاده کرد؟ دلایل انتخاب خود را شرح دهید.

ب) فرض کنید که ویژگی‌های موجود در مجموعه داده آموزش در مقیاس‌های بسیار متفاوتی هستند. این ویژگی روی هر کدام از الگوریتم‌های مطرح شده، چه تاثیری می‌گذارد و راه حل شما چیست؟

سوال ۵) فرض کنید یک مدل با واریانس بالا داریم و برای بهبود عملکرد مدل تصمیم گرفته‌ایم که تعداد داده‌های آموزشی را افزایش دهیم. به نظر شما چه نتیجه‌ای حاصل خواهد شد؟ سوال قبل را در مورد مدلی که بایاس بالایی دارد نیز پاسخ دهید.

سوال ۶) در شکل داده شده موارد زیر را مشخص کنید.

خطای آموزش، واریانس، مربع بایاس^۴، خطای ناشی از نویز^۵، خطای تست، محدوده $overfitting$ ، محدوده $underfitting$



بخش دوم: پیاده‌سازی

سوال اول:

مجموعه داده `dataset1.csv` را که در فولدر `dataset1` قرار دارد در نظر بگیرید.

الف) داده‌ها را رسم کنید.

ب) شافل کردن داده‌ها به چه منظور انجام می‌شود؟ آیا مجموعه داده این سوال نیازی به این اقدام دارد؟ در صورت لزوم این مورد را بر روی مجموعه داده اعمال کنید.

ج) با استفاده از روش گرادیان نزولی نمودارهایی با درجه‌های ۵ و ۸ و ۱۰ و با تعداد تکرارهای ۵۰۰۰ و ۱۰۰۰۰ بر روی داده‌ها برازش دهید و نمودارهای حاصل را رسم کنید. این عمل را برای سه معیار خطای MSE و $RMSE$ و MAE تکرار کنید و نتایج را مقایسه کنید. همچنین در صورت مشاهده بیش‌برازش آن را گزارش کنید.

² feature

³ scale

⁴ Bias squared

⁵ irreducible error

- د) برای هر یک از موارد قسمت قبل نمودار خطای آموزش و آزمون و نمودار اندازه قدم را رسم کنید. محور افقی نشان‌دهنده تکرارها و محور عمودی مقدار خطا و اندازه قدم را نشان دهد.
- ه) قسمت قبل را با معادله‌ی نرمال و بدون در نظر گرفتن ضریب λ تکرار کنید و نتایج را مقایسه کنید.
- و) به ازای درجه ۸، با روش معادله‌ی نرمال و با در نظر گرفتن سه مقدار مختلف برای ضریب λ نموداری بر روی داده‌ها برازش دهید و مقدار خطای RMSE را برای داده‌های آموزش و آزمون رسم کنید. تاثیر ضریب λ را بررسی نمایید.

سوال دوم:

مجموعه داده CSM_dataset که در فولدر dataset2 قرار دارد را در نظر بگیرید. این مجموعه داده مربوط به فیلم‌ها و امتیاز IMDB آنها می‌باشد که از سه منبع YouTube و Twitter و IMDB جمع‌آوری شده‌است. در این سوال می‌خواهیم مقدار متغیر Ratings را پیش‌بینی کنیم.

الف) در برخی از ستون‌های مجموعه داده مقادیر گمشده وجود دارد. با استفاده از روش‌های مناسب، مقادیر گمشده را پر کنید و به طور مختصر توضیح دهید که چگونه این کار را انجام دادید. همچنین اگر داده‌ها نیاز به پیش‌پردازش دیگری دارند، در این قسمت انجام دهید و در گزارش خود ذکر کنید.

ب) نمودار همبستگی بین ویژگی‌ها را رسم کنید و در گزارش بیاورید.

ج) با توجه به نمودار رسم شده آیا می‌توان یک یا چند ویژگی را حذف کرد؟ دلیل خود را ذکر کنید.

د) با استفاده از روش گرادیان نزولی یک نمودار بر روی داده‌ها برازش دهید. (پارامترها را به گونه‌ای انتخاب کنید که به بهترین خروجی دست یابید) این عمل را یک بار با استفاده از کل ویژگی‌ها و یک بار با استفاده از ویژگی‌های منتخب تکرار کنید. نمودار خطای آموزش و آزمون و نمودار طول گام را رسم کنید.

* سوال سوم (امتیازی)

داده‌های مربوط به ویژگی‌های شخصی و هزینه‌ی پزشکی افراد در فولدر dataset3 قرار دارد. هر داده دارای ۶ ویژگی ورودی و یک خروجی می‌باشد که هزینه‌ی پزشکی فرد را نشان می‌دهد. همانطور که در داده‌های سوال دیده می‌شود، ویژگی‌های جنسیت، منطقه و سیگاری بودن از نوع categorical هستند. برای تبدیل این ویژگی‌ها به ویژگی‌های عددی در این سوال می‌خواهیم از دو روش integer encoding و one hot encoding (OHE) استفاده کنیم.

الف) توضیح دهید که integer encoding و one hot encoding (OHE) چگونه انجام می‌گیرد.

ب) ویژگی‌های جنسیت و سیگاری بودن را با استفاده از integer encoding و ویژگی منطقه را با استفاده از OHE به مقدار عددی تبدیل کنید. (برای جنسیت: female=0, male=1 - برای سیگاری بودن: no=0, yes=1)

ج) به نظر شما چرا برای ویژگی منطقه از OHE استفاده می‌کنیم و از integer encoding استفاده نمی‌کنیم؟

د) رگرسیون خطی تعمیم یافته را بدون جمله‌ی منظم ساز پیاده سازی کرده، نتایج را روی داده‌های تست گزارش کنید. تابع basis برای ویژگی سن به صورت زیر باشد:

⁶ iterations

⁷ Missing values

$$\phi_{age}(x_{age}) = x_{age}^2$$

و برای سایر ویژگی‌ها از تابع همانی استفاده کنید:

$$\phi(x) = x$$

ه) در این قسمت ابتدا تعداد داده‌های آموزش را ۱۰ در نظر بگیرید و با گام‌های ۱۰ تایی، تا ۱۰۰۰ افزایش دهید. تغییرات خطای تست و آموزش را با افزایش داده آموزش بررسی کنید. برای مقایسه پذیر بودن خطا، از MSE برای تابع هزینه استفاده کنید. نمودار این تغییرات را در گزارش بیاورید.

و) روش stochastic gradient descent را بر روی این مجموعه داده پیاده سازی کنید و مقدار دقت و خطا را گزارش کنید.

پایدار باشید

نکات مربوط به تحویل تمرین

- کدهای خود را ترجیحا به زبان پایتون و در محیط jupyter پیاده‌سازی کنید.
- نظم و خوانایی در نوشتن گزارش و کدها از اهمیت بالایی برخوردار است. کدهای خود را تا حد امکان کامنت‌گذاری کنید.
- در پیاده‌سازی بخش‌های مختلف امکان استفاده از کتابخانه‌های آماده مربوط به الگوریتم‌های یادگیری ماشین را به طور کلی ندارید. مگر در مواردی که در صورت سوال ذکر شده باشد.
- برای خواندن داده‌ها می‌توانید از کتابخانه pandas و برای نمایش نمودارها و عملیات ماتریسی می‌توانید از کتابخانه‌های numpy و matplotlib استفاده کنید. برای محاسبه معیارهای ارزیابی مانند دقت، ماتریس درهم‌ریختگی و تقسیم داده‌ها به مجموعه‌های آموزش و آزمون نیز استفاده از کتابخانه آماده مجاز است.
- در صورتی که داده‌ها را به دسته‌های آموزشی، تست (و یا validation) تقسیم می‌کنید، درصد هر کدام را در گزارش ذکر کنید.
- فایل‌های کد و گزارش خود را در قالب یک فایل فشرده با فرمت HW01_StdNumber.zip که StdNumber شماره دانشجویی شماست، در سامانه بارگذاری کنید.
- سوالات ستاره دار(*) دارای نمره اضافی بر تمرین است.
- مطابق قوانین دانشگاه هر گونه کپی‌برداری ممنوع است و در صورت مشاهده، نمره هر دو طرف صفر می‌شود.
- در صورت وجود هر گونه سوال یا ابهامی با ایمیل درس در تماس باشید:

ml.ce.aut@gmail.com