

به نام خدا

تمرین سری اول درس یادگیری ماشین

پاییز ۱۴۰۴ - دانشگاه گیلان

مدرس: یاسمن برشبان

سوال ۱: بایاس - واریانس (۱۰ نمره)

- الف) با افزایش تعداد داده های آموزشی، واریانس و بایاس مدل یاد گرفته شده چگونه تغییر میکند؟ (۵ نمره)
ب) در مساله رگرسیون، از تابع درجه ۵ برای تخمین استفاده شده است، اما خطای آموزش قابل توجهی داریم.
درجه تخمین را باید چگونه تغییر دهیم؟ (۵ نمره)

سوال ۲: دسته بندی خطی: (۱۵ نمره)

- الف) مشکل استفاده از تابع هزینه SSE در مساله دسته بندی چیست؟ (۷ نمره)
ب) برتری تابع هزینه روش دسته بندی پرسپترون نسبت به تابع هزینه ای که تنها تعداد نمونه های غلط دسته بندی شده را در نظر می گیرد چیست؟ (۸ نمره)

سوال ۳: LDA (۱۵ نمره)

فرض کنید داده های مساله، نقاط زیر در فضای دو بعدی باشند. راستای بهینه را با استفاده از الگوریتم LDA محاسبه کنید.

$$C_1 = \{(1, 1), (1, 2), (2, 1), (2, 4), (3, 1), (3, 3)\}$$

$$C_2 = \{(2, 2), (3, 4), (4, 2), (5, 1), (5, 4), (5, 5)\}$$

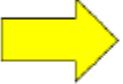
سوال ۴: رگرسیون (۶۰ نمره - پیاده سازی)

در این سوال می خواهیم از رگرسیون خطی برای پیش‌بینی هزینه پزشکی افراد بر اساس ویژگی های شخصی آن ها استفاده کنیم. داده ها در دو کلاس `test.csv` و `train.csv` ذخیره شده اند. هر داده دارای ۶ ویژگی ورودی (X) است و یک خروجی هزینه پزشکی (y) دارد.

لینک داده

همان طور که در داده های سوال دیده می شود، ویژگی های جنسیت، منطقه و سیگاری بودن از نوع دسته بندی شده هستند (Categorical) و عددی نیست. برای انکود کردن این ویژگی ها روش های متفاوتی وجود دارد. در این تمرین از دو روش استفاده می کنیم:

- به هر دسته، یک عدد یکتا (معمولًا بین ۰ تا $n-1$) اختصاص داده می شود.
 - برای n دسته در یک ویژگی، n متغیر از نوع باینری در نظر گرفته می شود و برای وردی که مقدار i را دارد، متغیر i ام می شود و دیگر متغیرها ۰ می شوند.
- همانند شکل زیر



Color	Red	Yellow	Green
Red	1	0	0
Red	1	0	0
Yellow	0	1	0
Green	0	0	1
Yellow	0	0	1

در داده های سوال، ویژگی ها را به صورت زیر تغییر دهید:
جنسیت: مقادیر male و female را به ترتیب تبدیل به ۰ و ۱ میکنیم

سیگاری بودن: مقادیر no و yes را به ترتیب تبدیل به ۰ و ۱ می کنیم.
منطقه: region در این ویژگی ۴ ستون وجود دارد. از روش OHE برای این ویژگی ها استفاده کنید. دلیلی که برای OHE region از integer encoding استفاده می شود این است که در روش OHE ترتیب بین مقادیر عددی ممکن است خطا ایجاد کند. در حالی که بین مناطق هیچ ترتیبی وجود ندارد اما در OHE این مشکل برطرف می شود.

قسمت الف) رگرسیون بدون منظم سازی:

رگرسیون خطی را بدون جمله منظم ساز پیاده سازی کرده، نتایج را روی داده های تست گزارش کنید. (۱۰ نمره)

ب) در ابتدا تعداد داده های تمرین 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 تابع هزینه از MSE استفاده کنید. نمودار این تغییرات را در گزارش بیاورید. (۱۰ نمره)

ج) روش های stochastic gradient descent و batch gradient descent را پیاده کنید و مقدار w و خطای خطا را در هر حالت به همراه گزارش کنید. (۲۰ نمره)

د) رگرسیون با منظم ساز: ۲۰ نمره

در این قسمت به تابع هزینه SSE جمله منظم ساز L2 را با ضریب α بیفزایید. با استفاده از ۵ fold cross validation بهترین مقدار پارامتر را برای α از بین اعضای مجموعه ۰.۰۰۰۱، ۰.۰۰۱،، ۱،، ۱۰۰۰۰ بیابید سپس نتیجه را بر روی داده های تست بدست آورید.

نمودار خطای خطا را بر حسب لگاریتم α رسم کنید. نتایج نهایی و مقدار خطای خطا را برای داده های تست و آموزش در گزارش بنویسید.

نکته:

- ۱- فایل گزارش را به همراه تمام کدها در یک فایل ZIP به اسم HW1 به همراه نام، نام خانوادگی و شماره دانشجویی آپلود کنید.
- ۲- در این تمرین می توانید از کتابخانه های آماده مانند NumPy، Matplotlib و scikit-learn برای پیاده سازی مدل استفاده کنید. هدف تمرین، درک رفتار الگوریتم و تحلیل نتایج مدل است، نه پیاده سازی از صفر.
- ۳- در صورت تشابه پاسخ دو دانشجو، نمره هر دو دانشجو از آن سوال صفر خواهد شد.

موفق باشید