



# مبانی یادگیری ماشین

نیم سال دوم ۰۱-۰۲

مدرس: دکتر حامد ملک

دانشکده‌ی مهندسی کامپیوتر

موعده تحویل: ۱۴۰۱/۱۲/۲۰

تمرین سری دوم

## مسائل تئوری

### مسئله‌ی ۱.

چرا در رگرسیون لجستیک از تابع هزینه آنتروپی متقاطع<sup>۱</sup> به جای خطای میانگین مربعات<sup>۲</sup> استفاده می‌شود؟

### مسئله‌ی ۲.

تا کنون مسائل طبقه‌بندی دودویی<sup>۳</sup> را آموخته‌اید و می‌دانید که کلاس نهایی داده بر اساس احتمال بیشتر یک کلاس انتخاب می‌شود. اکنون در این سوال از شما خواسته می‌شود در رابطه با مسئله طبقه‌بندی چندکلاسه<sup>۴</sup> تحقیق کنید و در رابطه با تکنیک یک در مقابل همه<sup>۵</sup> بنویسید.

### مسئله‌ی ۳.

یک دیتاست چند متغیره را در نظر بگیرید، فرض کنید مشاهده کرده‌ایم که یکی از ضرایب محاسبه‌شده در عملیات رگرسیون خطی مقدار خیلی بزرگ منفی نسبت به باقی متغیرها پیدا کرده است کدام یک از گزاره‌های زیر صحیح است؟ توضیح دهید.

- این ویژگی تاثیر زیادی روی مدل دارد و باید حفظ شود.
- این ویژگی تاثیر زیادی روی مدل ندارد و باید در نظر گرفته نشود.
- نمی‌توان بدون در دست داشتن اطلاعات بیشتر درمورد این ویژگی نظر داد.

### مسئله‌ی ۴.

با ارائه دلیل صحیح یا غلط بودن هر یک از گزاره‌های زیر را ثابت کنید.

- اگر بایاس زیاد است اضافه کردن تعداد داده‌های آموزش کمک زیادی به کم کردن بایاس نمی‌کند.

<sup>1</sup>Cross Entropy

<sup>2</sup>Mean squared error (MSE)

<sup>3</sup>Binary Classification

<sup>4</sup>Multiclass Classification

<sup>5</sup>One-vs-rest

- کم کردن خطای مدل روی داده های آموزش منجر به کاهش خطای مدل روی داده های تست می شود.
- افزایش پیچیدگی مدل رگرسیون همواره منجر به کاهش خطای مدل روی داده ی آموزش و افزایش خطای مدل روی داده ی تست می شود.

## مسائل کدی

### پروژه ی ۱.

مجموعه داده `Car_prices_classification.csv` شامل اطلاعات ۱۱۸ هزار ماشین است که مجموعاً به پنج دسته قیمت تقسیم شده اند. در این مسئله مراحل زیر را انجام دهید:

- الف) عملیات پیش پردازش<sup>۶</sup> را با توجه به هدف مسئله انجام دهید.
- ب) داده ها را با نسبت مناسب به داده آموزشی و داده تست تقسیم کنید. (دلیل انتخاب درصد نسبت داده آموزشی و تست را گزارش کنید.)
- پ) مدل رگرسیون لجستیک را با توجه به مباحث تدریس شده در کلاس درس از پایه پیاده سازی کنید (چرخ را دوباره اختراع کنید: ) و درصد دقت مدل را روی داده تست گزارش کنید.
- ت) مدل رگرسیون لجستیک را به کمک کتابخانه `scikit-learn` پیاده سازی کنید و مجدداً درصد دقت مدل را روی داده تست گزارش کنید.
- ث) در پایان نمودار کاهش خطا را ترسیم نمایید. همچنین نمودار افزایش دقت<sup>۷</sup> و دقت نهایی را گزارش کنید.

### پروژه ی ۲.

الگوریتم پرسپترون<sup>۸</sup> را برای یک شبکه عصبی با تعداد  $m$  نورون پیاده سازی کنید. این شبکه را با استفاده از مجموعه داده Iris<sup>۹</sup> تعلیم دهید. این داده به صورت کلی دارای سه کلاس است. مسئله را با یکسان کردن برچسب کلاس برای نمونه های دو کلاس `versicolor` و `setosa` به یک مسئله دو کلاس تبدیل کنید. از سوی دیگر برای تسهیل ترسیم، تنها دو ویژگی طول کاسبرگ و عرض کاسبرگ از چهار ویژگی را در نظر بگیرید و نمونه ها را رسم کنید. برای هر کلاس یک رنگ لحاظ کنید، مرز تصمیم نورون تعلیم دیده را به همراه بردار وزن به دست آمده را در کنار نمونه ها ترسیم کنید. نمودار تغییرات خطا را رسم کرده و مقدار خطای نهایی را ارائه کنید.

## نکات تمرین

- در این تمرین هدف این است که تا حد امکان از کتابخانه های `Numpy`، `Pandas` و `Matplotlib` استفاده کنید. عدم استفاده از این کتابخانه ها و استفاده از `pure python` منجر به کسر نمره خواهد شد.
- در صورت هرگونه تقلب نمره صفر برای شما لحاظ می گردد.

<sup>۶</sup>Preprocessing

<sup>۷</sup>Accuracy

<sup>۸</sup>Perceptron

<sup>۹</sup>قابل دسترسی با استفاده از `sklearn.datasets.load_iris`

— استفاده از زبان غیر از پایتون مجاز نیست.

— تمرین تحویل حضوری خواهد داشت.

— پیوند عمومی نوت‌بوک خود را در سامانه کوئرا بارگذاری نمایید.