

# به نام خدا

یادگیری ماشین



مهلت ارسال: ۲۰ آبان ۱۳۹۶

تکلیف نخست - پاییز ۱۳۹۶

قبل از اینکه شروع به حل تمرین‌ها کنید، حتماً یک بار فایل **Policies.pdf** (موجود در وبگاه) که حاوی نکات مهم در تحویل تکالیف هست را مطالعه فرمایید.  
حل سوالاتی که امتیازی هستند، اجباری نبوده و نمره‌ی اضافی خواهد داشت.  
طراح هر یک از سوالات در زیر مشخص شده است. در صورتی که سوالی دارید می‌توانید سوالاتان را در گروه مطرح کنید یا از دستیار درس بپرسید.

۱- با توجه به شکل زیر به سوالات پاسخ دهید:

به نام خداوند بخشنده بختايشگر

A A A

1 ἐν ἀρχῇ ἐποίησεν ὁ θεὸς τὸν οὐρανὸν καὶ τὴν γῆν

קִיצוּ שְׁכוּרִים וּבְכוּ וְהִלְלוּ כָּל־שְׁתֵּי יַיִן עַל־עֹסִים בְּי 5  
בְּרַת מְפִיכֵם:

شکل ۱

- (الف) همان‌طور که می‌دانید تعریف‌های گوناگونی برای یادگیری ماشین ارائه شده است. با توجه به گستردگی تعاریف، تعریفی از یادگیری ماشین ارائه کنید که بتواند به بهترین صورت آنالیز داده‌های شکل ۱ را به عنوان کاربردی از یادگیری ماشین توصیف کند.
- (ب) برای شکل ۱ مساله‌ای طراحی کنید و با استفاده از یادگیری ماشین برای آن راه‌حلی پیشنهاد نمایید.
- (ج) مساله تشخیص خط را در نظر بگیرید، یادگیری ماشین چه کمکی در بدست آوردن راه‌حل می‌تواند انجام دهد؟
- (د) آیا می‌توانیم ادعا کنیم در حل مساله تشخیص خط به یادگیری ماشین نیازمندیم؟

۲- عموماً بین پیچیدگی فضای فرضیه، مقدار خطای generalization و تعداد داده‌های آموزشی رابطه‌ای متناظر وجود دارد بیان کنید تغییرات هریک چه اثری بر دیگری دارد.

۳- روش یادگیری Naive Bayes Classifier بر پایه این فرض ساده (Naive) عمل می‌کند که:

ویژگی‌های مختلف مستقل هستند.

Naive Bayes classifier با استفاده از رابطه‌ای که بین tokenها (برای نمونه کلمات) وجود دارد به منظور spam filtering استفاده می‌شود، به این شکل که بتواند با استفاده از تکنیک دسته‌بندی bayes با محاسبه احتمال spam بودن یا نبودن یک email اقدام به جداسازی آن‌ها نماید. حال فرض کنید پیام‌های مشکوک آلوده با کلمه "replica" باشد. بیشتر افرادی که ایمیل دریافت می‌کنند می‌دانند وجود چنین کلمه‌ای نشان دهنده احتمال بالای spam بودن آن است. برای نمونه می‌تواند یک پیشنهاد جهت فروش کپی جعلی یک ساعت با برند شناخته شده باشد. یک spam detection software عموماً از وجود بسیاری از جزئیات اطلاعاتی ندارد از این رو اقدام به محاسبه احتمال می‌نماید.

$Pr(S)$ : احتمال spam بودن هر پیام.

$Pr(W|S)$ : احتمال اینکه کلمه "replica" در پیام‌های spam دیده شود.

$Pr(H)$ : احتمال اینکه هر پیامی که دیده شده است spam نباشد.

$Pr(W|H)$ : احتمال اینکه کلمه "replica" در پیام‌های سالم دیده شود.

الف) فرض کنید  $Pr(S|W)$  احتمال spam بودن پیام به در میان پیام‌های حاوی کلمه "replica" باشد. با در نظر گرفتن توابع احتمالاتی مورد نیاز به صورت پارامتری احتمال  $Pr(S|W)$  را محاسبه نمایید.

ب) حداقل ۴ دلیل ذکر کنید که چرا در نظر گرفتن فرض  $idd^1$  باعث مختل شدن عملکرد spam filter می‌شود. و در هر کدام نحوه ایجاد اختلال spam را توضیح دهید.

ج) براساس پیاده‌سازی که جهت طراحی spam filter انجام می‌دهید، Spam filter شما ممکن است به spam حساسیت‌های مختلفی نشان دهد یا حتی اصلاً حساس نشود. کسانی که spam را تولید می‌کنند سعی می‌کنند به طریقی این حساسیت را کاهش دهند تا بتوانند از این موانع عبور کنند. به عنوان یک فروشنده جنس قلایبی! سعی کنید ایده‌هایی جهت عبور از این spam filter ارائه دهید.

<sup>1</sup> Independent and identically distributed random variables

۴- فرض کنید در یک غار تاریک با ۳ خروجی غیر قابل تشخیص و پنهان در دیوارها گیر افتاده اید. یکی از درها به مسیری ۳ ساعته برای خروج به سمت بیرون ختم می‌شود. دو در دیگر نیز دارای مسیری است که شما را در طی ۱ و ۲ ساعت دوباره به غار برمی‌گرداند. شما راهی برای تشخیص تمایز درها ندارید. تخمین بزنید چقدر طول می‌کشد به بیرون برسید.

۵- فرض کنید داده های  $X_1, \dots, X_n$  داده هایی iid از یک توزیع یکنواخت حول دایره ای به شعاع  $\theta$  در فضای  $R^2$  باشند به طوریکه  $X_i \in R^2$  و

$$P(x|\theta) = \begin{cases} \frac{1}{\pi\theta^2} & \text{if } |x| \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

با استفاده از MLE،  $\theta$  را تخمین بزنید.

۶- (امتیازی) در این سوال با برنامه نویسی در محیط Python آشنا خواهید شد. برای شروع کدهای زیر را در محیط برنامه کپی نمایید سپس گام به گام مراحل خواسته شده در زیر را طی کنید. داده‌هایی که مورد استفاده قرار می‌گیرد داده‌های "Fisher iris" می‌باشد این داده‌ها شامل ۴ اندازه حقیقی به عنوان ویژگی‌های سه نوع گل زنبق (iris) می‌باشد.

```
import numpy as np
import matplotlib.pyplot as plt

iris = np.genfromtxt("data/iris.txt", delimiter=None) # load the text file
Y = iris[:, -1] # target value is the last column
X = iris[:, 0:-1] # features are the other columns
```

الف) با استفاده از برنامه نویسی تعداد ویژگی‌ها و تعداد داده‌ها را بدست آورید.

ب) برای هر ویژگی هیستوگرام مقادیر داده‌ها را رسم کنید.

ج) برای هر ویژگی میانگین آن‌ها را محاسبه کنید

د) برای داده‌های هر ویژگی مقادیر واریانس و انحراف معیار را بدست آورید.

ه) داده‌ها را با محاسبه تفاضل میانگین از داده‌ها بخش بر انحراف معیار نرمال نمایید.

۷- روش Bayes یکی از راه‌های بدست آوردن مرز تصمیم در دسته‌بندی‌هاست<sup>۲</sup>. در این حالت ما مرز تصمیم بهینه را برای دو کلاس  $C_1$  و  $C_2$  با استفاده از نسبت likelihood به دست می‌آوریم.

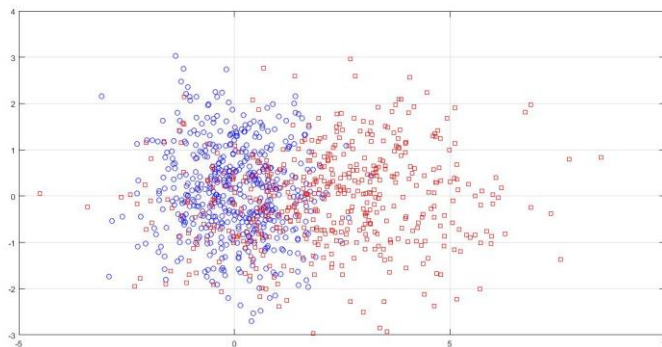
داده‌های  $Class1$  و  $Class2$  مربوط به شکل زیر را در نظر بگیرید هر یک شامل ۵۰۰ داده دویعدی است. قصد داریم این داده‌ها را به وسیله یک مرز از هم جدا سازیم:

فرض کنیم می‌دانیم، توزیع داده‌ها در کلاس  $C_1$  بصورت گوسی با بردار میانگین  $[0,0]$  و واریانس ۱ و برای کلاس  $C_2$  با بردار میانگین  $[0,2]$  و واریانس ۴ می‌باشد. و threshold ما ۱ باشد.

الف) ابتدا داده‌ها را در محیط متلب وارد کنید و شکل هر کدام را ترسیم کنید.

ب) با استفاده از تکنیک تخمین با استفاده از بیشینه‌ی درست‌نمایی سعی کنید مرزی برای جداسازی این دو کلاس را به صورت دستی و با کد متلب بیابید.

ج) به نظر تان مرز بدست آمده بهینه است.



شکل ۲- پراکندگی داده‌ها با هر دو کلاس

<sup>۲</sup> Bayesian Decision Boundary