



پردیس دانشکده‌های فنی

بسمه تعالی
دانشکده مهندسی برق و کامپیوتر
تمرین سری چهارم درس یادگیری ماشین



دانشگاه تهران

سلام بر تمام دانشجویان عزیز، چند نکته مهم:

۱. حجم گزارش به هیچ عنوان معیار نمره‌دهی نیست، در حد نیاز توضیح دهید.
۲. در این تمرین شما حداکثر ۱۳۰ نمره می‌توانید کسب کنید بنابراین از پاسخ دادن به سوالات با مجموع بارم بیشتر از ۱۳۰ نمره خودداری کنید.
۳. حداقل به سه سوال شبیه‌سازی باید پاسخ دهید.
۴. برای سوالات شبیه‌سازی، فقط از دیتاست داده شده استفاده از کنید. شکل‌ها، به طور واضح و در فرمت درست گزارش شود.
۵. در سوال ۱۰ از دیتاست **FACES** و در سوال ۹ تا ۱۳ از دیتاست **TinyMNIST** استفاده کنید. در صورت زمانبر بودن پردازش کد در مورد سوالات ۹، ۱۲ و ۱۳ می‌توانید از **TinyMNIST_loader** که پیوست شده نیز استفاده نمایید. اما دقت بفرمایید که در مورد سوالات ۱۰ و ۱۱ مجاز به کاهش تعداد ویژگی‌ها/تعداد داده‌ها نیستید.
۶. هرگونه شباهت در گزارش و کد مربوط به شبیه‌سازی، به منزله **تقلب** می‌باشد و کل نمره تمرین **صفر** می‌شود.
۷. در صورت داشتن سوال، از طریق ایمیل afsaneh.h.ebrahimi@gmail.com سوال خود را مطرح کنید.

۱. فرض کنید که $p_x(x|w_i)$ بیانگر چگالیهای احتمالی با میانگین μ_i و واریانس Σ_i برای $i = 1, 2$ باشد که لزوماً هم نرمال نیست. $y = w^t x$ بیانگر **projection** باشد و چگالیهای تک بعدی $p(y|w_i)$ دارای میانگین μ_i و واریانس δ_i^2 هستند. نشان دهید که تابع هزینه‌ی $J(w)$

$$J(w) = \frac{(\mu_1 - \mu_2)^2}{\delta_1^2 + \delta_2^2}$$

به وسیله‌ی $w = (\Sigma_1 + \Sigma_2)^{-1}(\mu_1 - \mu_2)$ بیشینه می‌گردد. (۲۰ نمره)

۲. در مسالهی طبقه‌بندی C کلاسه، ماتریس پراکندگی درون کلاسی و بین کلاسی به ترتیب به صورت زیر تعریف میشود:

$$S_w = \sum_{k=1}^C \sum_{x^q \in w_i} (x^q - \mu_k)(x^q - \mu_k)^T$$

$$S_B = \sum_{k=1}^C N_k (\mu_k - \mu)(\mu_k - \mu)^T$$

- ا. نشان دهید که $rank(S_B) \leq C - 1$ و در چه شرایطی $rank(S_B) = C - 1$ ؟ (۵ نمره)
- ب. درباره‌ی حداکثر تعداد مقادیر ویژه‌ی ناصفر ماتریس جدایی‌پذیری $S_B^{-1} S_w$ بحث نمایید. (۵ نمره)

ت. نشان دهید که $S_T = S_W + S_B$ (۱۰ نمره)

۳. عبارت

$$J = \frac{1}{n_1 n_2} \sum_{y_i \in Y_1} \sum_{y_j \in Y_2} (y_i - y_j)^2$$

پراکندگی کل درون گروهی (within group scatter) را اندازه میگیرد نشان دهید که این عبارت را میتوان به صورت زیر هم نوشت: (۲۰ نمره)

$$J = (m_1 - m_2)^2 + \frac{1}{n_1} s_1^2 + \frac{1}{n_2} s_2^2$$

۴. فرض کنید که پهنجری پارزن به صورت زیر تعریف شده باشد:

$$\varphi(x) = \begin{cases} e^{-x} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

و همچنین $p(x) \sim U(0, a)$ باشد. نشان دهید که میتوان میانگین پهنجری پارزن را با n نمونه به صورت زیر تخمین زد: (۲۰ نمره)

$$\bar{p}_n(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}} \right) & 0 \leq x \leq a \\ \frac{1}{a} \left(e^{\frac{a}{h_n}} - 1 \right) e^{-\frac{x}{h_n}} & a \leq x \end{cases}$$

۵. یک رستوران بر این است که بررسی نماید با توجه به عوامل موثر، افرادی که به رستوران مراجعه میکنند در صورتی که

تمام میزها پر باشد آیا برای خالی شدن میز صبر میکند یا نه؟

داده‌های ثبت شده از ۱۲ مراجعه کننده، جنبه‌های مختلف و اینکه صبر میکنند/نمیکند را در جدول ۱ مشاهده

میفرمایید.

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
x ₁	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0-10	y ₁ = Yes
x ₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30-60	y ₂ = No
x ₃	No	Yes	No	No	Some	\$	No	No	Burger	0-10	y ₃ = Yes
x ₄	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10-30	y ₄ = Yes
x ₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	y ₅ = No
x ₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0-10	y ₆ = Yes
x ₇	No	Yes	No	No	None	\$	Yes	No	Burger	0-10	y ₇ = No
x ₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0-10	y ₈ = Yes
x ₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	y ₉ = No
x ₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10-30	y ₁₀ = No
x ₁₁	No	No	No	No	None	\$	No	No	Thai	0-10	y ₁₁ = No
x ₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30-60	y ₁₂ = Yes

جدول ۱ داده‌های ثبت شده از ۱۲ مراجعه کننده

توضیح جنبه‌های مختلف که در Input attributes آمده است به شرح زیر است:

1.	Alternate: whether there is a suitable alternative restaurant nearby.
2.	Bar: whether the restaurant has a comfortable bar area to wait in.
3.	Fri/Sat: true on Fridays and Saturdays.
4.	Hungry: whether we are hungry.
5.	Patrons: how many people are in the restaurant (values are None, Some, and Full).
6.	Price: the restaurant's price range (\$, \$\$, \$\$\$).
7.	Raining: whether it is raining outside.
8.	Reservation: whether we made a reservation.
9.	Type: the kind of restaurant (French, Italian, Thai or Burger).
10.	WaitEstimate: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60).

حال برای این مساله decision tree را رسم نمایید و با توجه به داده‌های موجود قضاوت نمایید که آیا مراجعه کننده در صورت پر بودن میزها صبر میکند یا خیر. (۲۰ نمره)

۶.

در مورد روش Adaboost مطالعه بفرمایید.

به عنوان مثال می‌توانید به لینک زیر مراجعه نمایید:

http://www.cs.man.ac.uk/~nikolaon/~nikolaon_files/Introduction_to_AdaBoost.pdf

۵. به طور کلی در مورد این روش و نقاط قوت/ضعف آن نسبت به decision tree توضیح مختصری بدهید. (۵ نمره)

ب. فرض کنید $h_t(x)$ یک طبقه بند ضعیف باشد که در مرحله t مشاهده شده است و α_t بیانگر وزنهای باشد. طبقه بند نهایی به شکل زیر است:

$$H(x) = \text{sign}(f(x)), \quad f(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

نشان دهید که برای خطای داده‌های آموزشی طبقه بند نهایی رابطه زیر برقرار است: (y_i) لیبل حقیقی برای x_i است. (۱۵ نمره)

$$\frac{1}{m} \sum_{i=1}^m (H(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-f(x_i)y_i)$$

۷.

- ا. برای نمونه‌های $X = \{-7, -5, -4, -3, -2, 0, 2, 3, 4, 5, 7\}$ در یک مساله یک بعدی، تخمین پنجره‌ی پارزن $P_j(x)$ را برای یک پنجره‌ی مستطیلی بیابید. ضمناً $h_j = \frac{1}{\sqrt{j}}$ است. پنجره را برای $j=1, 4, 11$ در یک شکل رسم نمایید. (۱۰ نمره)
- ب. در رابطه $h_j = \frac{h}{\sqrt{j}}$ ، h را ثابت بگیرید و درباره شکل $P_j(x)$ به ازای h های مختلف بحث کنید. (۱۰ نمره)

۸. یکی از تکنیکهای متداول در فشرده‌سازی تصاویر PCA میباشد که تعداد Principle Components (PCs) در کیفیت تصویر و نرخ فشرده سازی تصویر تاثیرگذار است. حال ما برآنیم که با استفاده از PCA تصاویر را فشرده کنیم و به فضایی با تعداد ویژگیهای کمتر منتقل شویم تا عملیات تشخیص تصویر بهتر انجام شود. در این سوال شما از شامل ۲۱۳ تصویر که دارای ۶ حالت (Happy, Fear, Angry, Disgust, Surprise and sad) استفاده میکنید. برای بارگذاری تصاویر از image_loader.py که با تصاویر پیوست شده است استفاده نمایید.
- ا. مقادیر ویژه از PCA را به ترتیب کاهشی رسم نمایید و بیان نمایید که چگونه میتوان تعداد کامپوننت مناسب را در فرآیند فشرده سازی تشخیص داد. (۱۰ نمره)
- ب. ۴ مقدار ویژه اول و ۴ مقدار ویژه نهایی (eigenfaces) را نشان دهید و تحلیل کنید. (۵ نمره)
- ا. حال با یک طبقه بند k-NN با مقادیر $k = 1, 2$ طبقه بندی را انجام دهید و CCR را گزارش کنید. (۱۰ نمره)

۹. الگوریتم forward selection را پیاده‌سازی نمایید. شما میتوانید از Naïve bayes optimal classifier در الگوریتم خود به عنوان طبقه‌بند استفاده کنید.

ا. CCR را بر حسب تعداد ویژگی‌های انتخاب شده در یک نمودار رسم نمایید. (۱۰ نمره)

ب. تعداد بهینه‌ی ویژگی‌ها را برای بهترین عملکرد طبقه‌بند بیان نمایید. (۱۰ نمره)

۱۰. حال لیبل کلاس‌ها را در نظر گرفته و ماتریس پراکندگی درون کلاسی و بین کلاسی را محاسبه نمایید تا روش LDA را پیاده‌سازی کنیم.

أ. از LDA کمک بگیرید و مقادیر ویژه را مرتب نمایید و مقادیر ویژه ماتریس جداپذیری را در قالب نزولی رسم نمایید. (۱۰ نمره)

ب. در یک نمودار مقدار $trace(S_W^{-1}S_B)$ (separability measure) نسبت به تعداد ویژگی‌ها رسم نمایید و در مورد تاثیر تعداد ویژگی‌ها بر آن بحث کنید. (۱۰ نمره)

ت. تعداد بهینه‌ی ویژگی‌ها را بر اساس دو بخش پیشین بیان نمایید و حال با انتقال به زیرفضای جدید که فقط شامل ویژگی‌های بهینه است، طبقه‌بند Naïve Bayes optimal classifier با تخمین پارامتری گوسی را پیاده‌سازی کنید و مقدار CCR را گزارش کنید. (۵ نمره)

۱۱. حال لیبل کلاس‌ها را در نظر گرفته و مقادیر و بردار ویژه‌های ماتریس کواریانس را حساب نمایید تا در این مساله روش PCA را اعمال کنیم.

أ. مقادیر ویژه‌ی ماتریس کواریانس را بر حسب شماره ویژگی رسم نمایید. (۱۰ نمره)

ب. تعداد بهینه‌ی ویژگی‌ها را بر اساس بخش پیشین بیان نمایید و حال با انتقال به زیرفضای جدید که فقط شامل ویژگی‌های بهینه است، طبقه‌بند Naïve Bayes optimal classifier با تخمین پارامتری گوسی را پیاده‌سازی کنید و مقدار CCR را گزارش کنید. (۱۰ نمره)

۱۲. برای داده‌ی مذکور، ابتدا یک طبقه‌بند بهینه‌ی بیز با روش پنجره پارزن برای تخمین pdf با حالتهای زیر طراحی کنید و مقدار CCR را گزارش نمایید:

أ. دو پنجره‌ی مستطیلی و گوسی را بررسی نمایید. (۱۰ نمره)

ب. تاثیر اندازه‌ی پنجره را بررسی کنید. (برای ۳ مقدار مختلف) (۱۰ نمره)

۱۳. سوال قبل را با استفاده از روش K نزدیکترین همسایه برای تخمین pdf دوباره تکرار کنید.

أ. برای سه مقدار مختلف k الگوریتم را تکرار کنید و نتایج را گزارش و تحلیل نمایید. (۱۰ نمره)

ب. در چه صورتی نتایج این سوال و سوال پیش میتوانند به یکدیگر همگرا شوند؟ توضیح دهید.. (۱۰ نمره)

۱۴. در این سوال قصد داریم، با استفاده از ICA، ۳ منبع تولیدکننده موسیقی که به طور همزمان کار میکنند را تشخیص دهیم. جهت اینکار، ابتدا سه سیگنال با طولهای برابر و به صورت سینوسی، پالسی و دندان اره‌ای تولید کنید. سپس، این سه سیگنال را با یکدیگر و با یک ضریب از نویز گوسی جمع کنید. در نهایت، پس از تولید mixing matrix ، با استفاده از ICA، این سه منبع را تشخیص دهید. (۲۰ نمره)