

بسمه تعالی



یادگیری ماشین

فروردین ۱۴۰۱

سلام بر تمام دانشجویان عزیز، چند نکته مهم:

۱. حجم گزارش به هیچ عنوان معیار نمره‌دهی نیست، در حد نیاز توضیح دهید.
۲. نکته‌ی مهم در گزارش نویسی روشن بودن پاسخها می‌باشد، اگر فرضی برای حل سوال استفاده می‌کنید حتما آن را ذکر کنید، اگر جواب نهایی عددی است به صورت واضح آن را بیان کنید.
۳. برای سوالات شبیه سازی، فقط از دیتاست داده شده استفاده کنید.
۴. شکل ها به طور واضح و در فرمت درست گزارش شوند. عکس‌ها را به صورت واضح و همراه با زیرنویس در گزارش خود بیاورید.
۵. هرگونه شباهت در گزارش و کد مربوط به شبیه سازی، به منزله تقلب می‌باشد و کل نمره تمرین صفر می‌شود.
۶. برای هر کد که در فایل نهایی ضمیمه می‌کنید، گزارش بنویسید. کدهای ضمیمه شده بدون گزارش مربوطه نمره‌ای نخواهند داشت. (این گزارش‌ها تنها معیار تفکیک کد شما و کدهای موجود در منابع مختلف مانند اینترنت خواهند بود).
۷. در صورت داشتن سوال، از طریق ایمیل mesbahamirhossein@gmail.com، سوال خود را مطرح کنید.

۱- سوال اول (۲۰ نمره)

توزیع نرمال $p(x) \sim N(\mu, \sigma^2)$ و تابع پنجره پارزن $\varphi(x) \sim N(0, 1)$ را در نظر بگیرید. نشان دهید که تخمین پنجره پارزن

$$P(x) = \frac{1}{nh_n} \sum_{i=1}^n \varphi\left(\frac{x - x_i}{h_n}\right)$$

برای h_n های کوچک دارای ویژگی های زیر است:

- $\tilde{p}_n(x) \sim N(\mu, h_n^2 + \sigma^2)$
- $p_n(x) - \tilde{p}_n(x) \cong \frac{1}{2} \left(\frac{h_n}{\sigma}\right)^2 \left[1 - \left(\frac{x-\mu}{\sigma}\right)^2\right] p(x)$
- $var[p_n(x)] \cong \frac{1}{2nh_n\sqrt{\pi}} p(x)$

۲- سوال دوم (۱۰ نمره)

متریک فاصله اقلیدسی را در d بعد در نظر بگیرید:

$$D(a, b) = \sqrt{\sum_{k=1}^d (a_k - b_k)^2}$$

فرض کنید عناصر هر بعد را در یک مقدار حقیقی غیر صفر ضرب میکنیم. یعنی $k = 1, 2, \dots, d$ داریم:

$$x_k = \alpha_k x_k$$

ثابت کنید پس از ضرب نیز این متریک فاصله همچنان یک فاصله استاندارد است یعنی ویژگی‌های گفته شده برای یک فاصله استاندارد را دارا می‌باشد. در مورد تاثیر این امر بر طبقه بند knn بحث کنید.

۳- سوال سوم (۲۰ نمره)

یک مسئله طبقه بندی با روش knn را در نظر بگیرید. مجموعه داده دو کلاسه D را نیز به صورت $D = \{x^q, \omega_i^q\}$ داریم. این داده ها نتایج یک نظر سنجی بوده و دیتاپوینت ها به صورت پرچسب خورده، مستقل از هم هستند و فرض می کنیم تعداد داده های دو کلاس یکسان است. برای هر سمپل تست نزدیک ترین k دیتاپوینت را به صورت $\{x_i\}_{i=1, \dots, k}$ نمایش می دهیم. هر دوی $p(x|1)$ و $p(x|2)$ توزیع یکنواخت بر روی یک کره به شعاع واحد دارند و مرکز دو ابر کره نیز از هم ۱۰ واحد فاصله دارند.

الف) نشان دهید اگر k فرد باشد متوسط احتمال خطا از رابطه زیر به دست می آید:

$$p_Q(e) = \frac{1}{2^Q} \sum_{j=0}^{\frac{k-1}{2}} \binom{Q}{j}$$

ب) با توجه به بخش قبل نشان دهید که در این حالت خطای طبقه بند نزدیک ترین همسایه کمتر از حالت $k \geq 2$ است و دلیل مشاهده این موضوع را توضیح دهید.

ج) نشان دهید: $\lim_{Q \rightarrow \infty} p_Q(e) = 0$

۴- سوال چهارم (۱۵ نمره)

بخش اول - ۵ نمره:

به سوالات زیر پاسخ دهید.

- مفهوم bias-variance trade off را با توجه به اندازه h_n در روش پارزن k_n در روش knn توضیح دهید.

- تفاوت روش های پارامتریک و نان پارامتریک را توضیح دهید.

- مشکلات روش های kernel based چیست.

- تفاوت مفهوم حجم در روش پارزن و knn را بررسی کنید.

بخش دوم (پیاده سازی) - ۱۰ نمره:

موارد خواسته شده در قسمت های مختلف را انجام داده و نتایج به دست آمده را تحلیل کنید.

الف) ۱۰۰۰ دیتاپوینت رندوم ۵ بعدی را تولید کرده و فاصله این ۱۰۰۰ دیتا پوینت را از هم حساب کرده و نمودار هیستوگرام فاصله ها را رسم نمایید. این کار را برای ابعاد ۲۰۰۰، ۵۰۰۰، ۱۰۰۰۰ و ۱۰۰۰۰۰۰ تکرار کنید.

ب) برای ۱۰۰۰ دیتا پوینت رندوم ۱۰۰۰۰ بعدی فاصله دیتاپوینت ها را از هم حساب کرده و نمودار هیستوگرام فاصله ها را رسم کنید. این کار را برای دیتاپوینت ها به تعداد ۲۰۰۰، ۵۰۰۰، ۱۰۰۰۰ و ۱۰۰۰۰۰۰ تکرار کنید.

ج) نتایج تحلیل خود را بیان کنید. با توجه به این نتایج عملکرد الگوریتم knn را چگونه ارزیابی می کنید. همچنین ارتباط این نتایج با curse of dimensionality را بیان کنید.

۵ - سوال پنجم (۱۵ نمره) - پیاده سازی

در این سوال میخواهیم به پیاده سازی روش تخمین نان پارامتری پارزن بپردازیم. لازم به ذکر است که الگوریتم خواسته شده در این سوال را باید بدون استفاده از کتابخانه های آماده موجود پیاده سازی کنید.

برای شروع ابتدا دیتاست [ted talks](#) را دانلود کنید.

الف) ستون **duration** این دیتاست را استخراج کرده و توزیع دیتای این ستون را با استفاده از روش پنجره پارزن با کرنل گوسی به دست آورده و نتیجه را نمایش دهید. اندازه پنجره را برابر با ۱۰ در نظر بگیرید.

ب) تاثیر اندازه پنجره را با ۳ مقدار ۲۰، ۵۰ و ۱۰۰ مختلف بررسی کنید.

ج) با استفاده از کتابخانه های آماده توزیع ستون **duration** را رسم کنید. با افزایش مقدار n روند تغییر و همگرا شدن به توزیع اصلی را روی یک نمودار نشان دهید. مقدار n را در بازه ۲۵۰ نمونه تا کل دیتا با **step** برابر با ۲۵۰ بررسی کرده و همگرایی برای n های مختلف را حتما روی یک نمودار نشان دهید.

د) نتیجه قسمت الف را با نتیجه توابع کتابخانه های آماده مقایسه کنید.

۶- سوال ششم (۲۰ نمره) - پیاده سازی

در این تمرین از دیتاست [fashion MNIST](#) استفاده خواهیم کرد. توصیه میشود با توجه به پردازش سنگین مورد نیاز این بخش از google colab برای این بخش استفاده کنید. ابتدا داده‌ها را به ۳ بخش train, validation و test جدا کنید. سپس یک مدل mlp برای طبقه‌بندی داده‌ها طراحی کنید.

الف) نمودار دقت و خطا (loss) را برای داده‌های train و validation و همچنین این مقادیر را برای داده‌های تست گزارش کنید.

ب) تاثیر تعداد لایه‌های مخفی و نرخ یادگیری و solver های مختلف را روی عملکرد مدل بررسی کنید. برای پارامترهای مختلف میتوانید گزینه‌های زیر را در نظر بگیرید.

- تعداد لایه‌های مخفی ۱، ۲، ۳ و ۴

- اندازه‌ی لایه‌های مخفی ۵۰، ۱۰۰، ۱۵۰ نورون

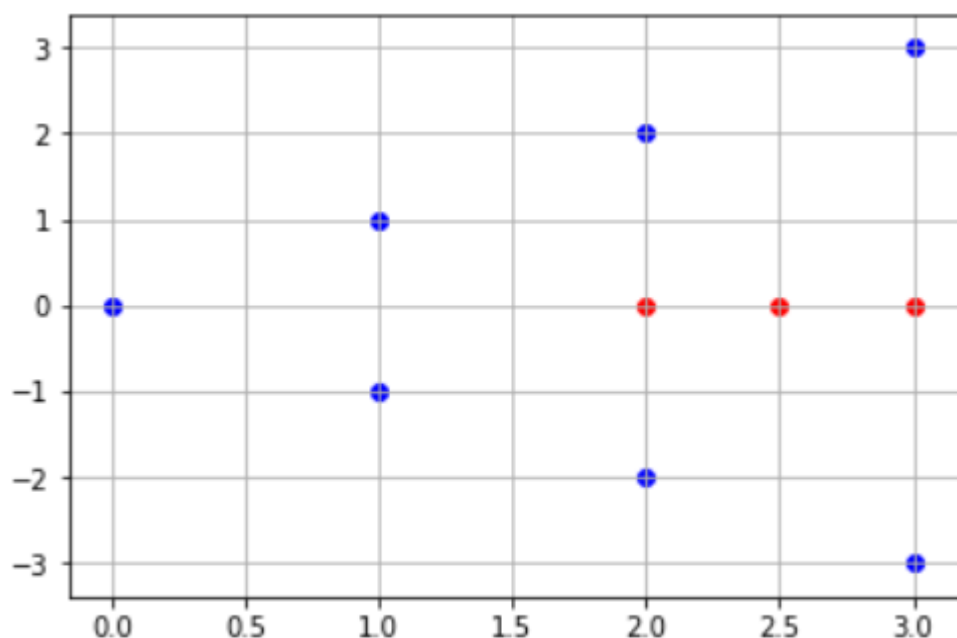
- solver های sgd, ADAM و rmsprop

- نرخ یادگیری ۰,۰۰۱، ۰,۰۱، ۰,۱، ۰,۵ و ۰,۹

ج) بهترین پارامترهایی که برای این مدل به دست آورده‌اید را در یک جدول گزارش کنید.

۷- سوال هفتم (۱۵ نمره)

تصویر زیر را در نظر بگیرید.



شکل ۱

برای جداسازی داده‌های دو کلاس یک شبکه پرسپترون طراحی کرده و پیاده سازی کنید. توجه داشته باشید که در این سوال مجاز به استفاده از کتابخانه و توابع آماده نمی‌باشید.

۸ - سوال هشتم (۲۰ نمره)

یک رستوران بر این است که بررسی نماید با توجه به عوامل موثر، افرادی که به رستوران مراجعه می کنند در صورتی که تمام میزها پر باشد، برای خالی شدن میز صبر می کنند یا نه؟

داده های ثبت شده از ۱۲ مراجعه کننده، جنبه های مختلف و اینکه صبر می کنند یا نه را در جدول ۱ مشاهده می فرمایید.

| 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 | $y_1 = \text{Yes}$ |
| x_2 | Yes | No | No | Yes | Full | \$ | No | No | Thai | 30-60 | $y_2 = \text{No}$ |
| x_3 | No | Yes | No | No | Some | \$ | No | No | Burger | 0-10 | $y_3 = \text{Yes}$ |
| x_4 | Yes | No | Yes | Yes | Full | \$ | Yes | No | Thai | 10-30 | $y_4 = \text{Yes}$ |
| x_5 | Yes | No | Yes | No | Full | \$\$\$ | No | Yes | French | >60 | $y_5 = \text{No}$ |
| x_6 | No | Yes | No | Yes | Some | \$\$ | Yes | Yes | Italian | 0-10 | $y_6 = \text{Yes}$ |
| x_7 | No | Yes | No | No | None | \$ | Yes | No | Burger | 0-10 | $y_7 = \text{No}$ |
| x_8 | No | No | No | Yes | Some | \$\$ | Yes | Yes | Thai | 0-10 | $y_8 = \text{Yes}$ |
| x_9 | No | Yes | Yes | No | Full | \$ | Yes | No | Burger | >60 | $y_9 = \text{No}$ |
| x_{10} | Yes | Yes | Yes | Yes | Full | \$\$\$ | No | Yes | Italian | 10-30 | $y_{10} = \text{No}$ |
| x_{11} | No | No | No | No | None | \$ | No | No | Thai | 0-10 | $y_{11} = \text{No}$ |
| x_{12} | Yes | Yes | Yes | Yes | Full | \$ | No | No | Burger | 30-60 | $y_{12} = \text{Yes}$ |

جدول ۱ - داده های ثبت شده از ۱۲ مراجعه کننده

توضیح فیچرهای مختلف نیز به شرح زیر است:

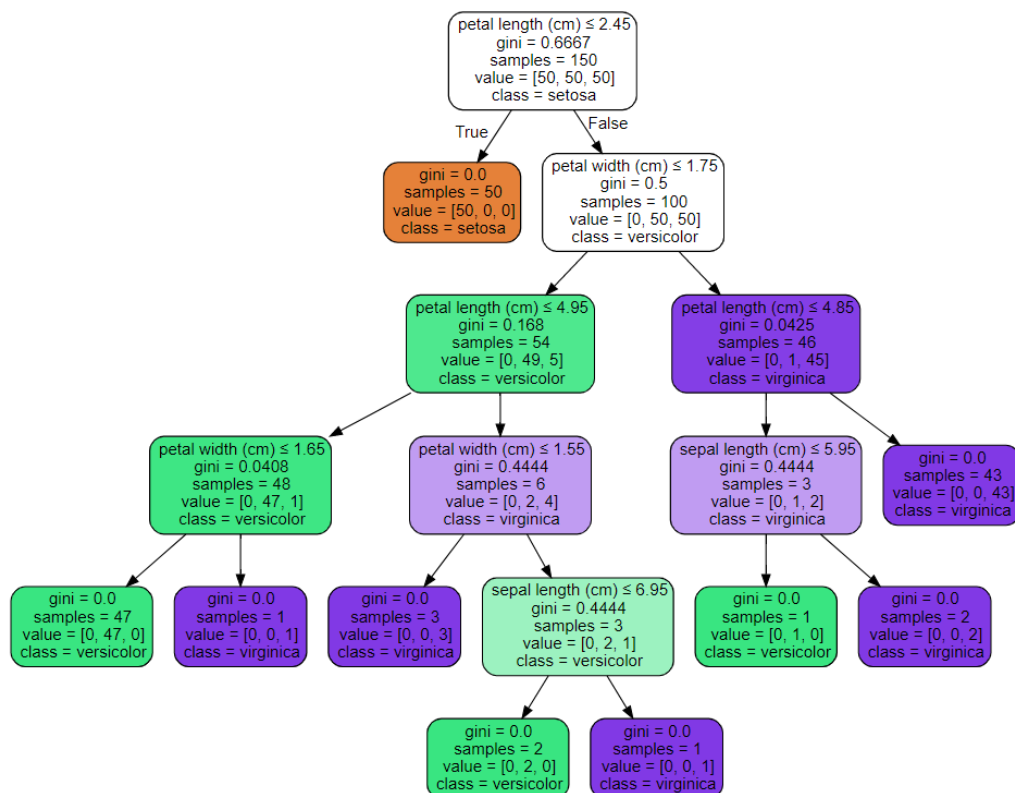
| | |
|-----|---|
| 1. | Alternate: whether there is a suitable alternative restaurant nearby. |
| 2. | Bar: whether the restaurant has a comfortable bar area to wait in. |
| 3. | Fri/Sat: true on Fridays and Saturdays. |
| 4. | Hungry: whether we are hungry. |
| 5. | Patrons: how many people are in 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). |

- مرحله (یا لایه) اول درخت تصمیم را با استفاده از معیار آنتروپی به صورت دستی حل کنید.

- طبقه بند درخت تصمیم را بدون استفاده از پکیج آماده و با در نظر گرفتن معیار آنتروپی کرده پیاده سازی کرده و نتایج پیاده سازی را گزارش کنید.

- طبقه بند درخت تصمیم را با استفاده از پکیج‌های آماده و با در نظر گرفتن معیار آنتروپی پیاده سازی کرده و آن را رسم کنید. شکل ۲ نشان‌دهنده نمونه از تصویر رسم شده برای درخت تصمیم می‌باشد.

- تفسیر خود از این درخت تصمیم را شرح دهید.



شکل ۲ - درخت تصمیم رسم شده توسط پکیج scikit learn