



DEPARTMENT OF COMPUTER
ENGINEERING AND IT



AMIRKABIR UNIVERSITY
OF TECHNOLOGY

تمرین دوم درس یادگیری ماشین

دکتر ناظر فرد

پاییز ۹۸

- تمامی مستندات خود شامل گزارش و کدهای خود را در یک فایل فشرده با فرمت زیپ ذخیره کرده و با عنوان studentId_HW.zip بارگزاری نمایید.
- مهلت انجام این تمرین تا ساعت ۲۳:۵۵ روز ۳۰ آبان می باشد و به هیچ وجه تمدید نمی شود.

تمرینات تشریحی: (۳۰٪)

بخش اول - در این بخش به بررسی درخت تصمیم^۱ و جنگل تصادفی^۲ پرداخته می شود. لطفاً به سوالات زیر پاسخ دهید. (۸٪)

- ۱- هرس^۳ درخت تصمیم چه تاثیری بر بیش برآزش^۴ دارد؟ این هرس چه زمانی باید انجام شود؟ توضیح دهید.
- ۲- با توجه به جدول زیر به سوالات زیر پاسخ دهید.

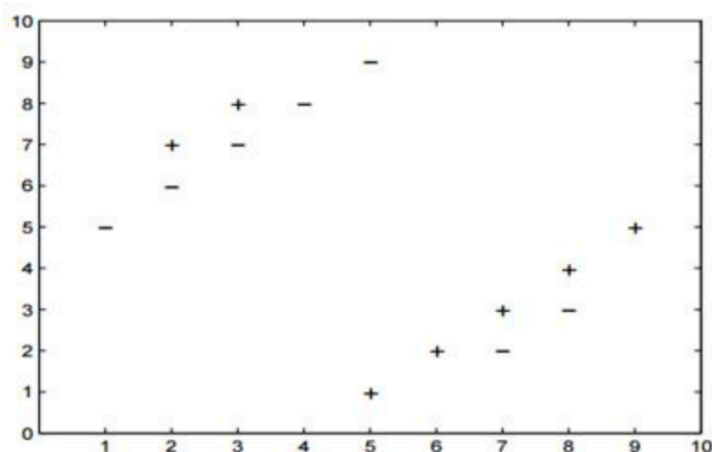
کلاس	ویژگی C	ویژگی B	ویژگی A
X	۱	۱	۱
X	۰	۱	۱
Y	۱	۰	۰
Y	۰	۰	۱

- الف - در ابتدا با توجه به ویژگی های موجود در جدول بیان نمایید که کدام ویژگی برای دسته بندی مناسب می باشد.
 - ب - سپس با استفاده از آنتروپی^۵ و بهره اطلاعات^۶ بهترین ویژگی را پیدا کنید.
 - پ - درخت تصمیم را با توجه به ویژگی های بخش قبل رسم نمایید.
- ۳- در رابطه با جنگل تصادفی و نقاط قوت آن در مقایسه با درخت تصمیم توضیح دهید. منبع را نیز ذکر نمایید (منبع مورد استفاده لازم است که مقاله باشد).

¹ Decision tree
² Random forest
³ Pruning
⁴ Overfit
⁵ Entropy
⁶ Information gain

بخش دوم- در این بخش به بررسی الگوریتم دسته‌بندی K نزدیک‌ترین همسایه^۸ پرداخته می‌شود. (۸٪)

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



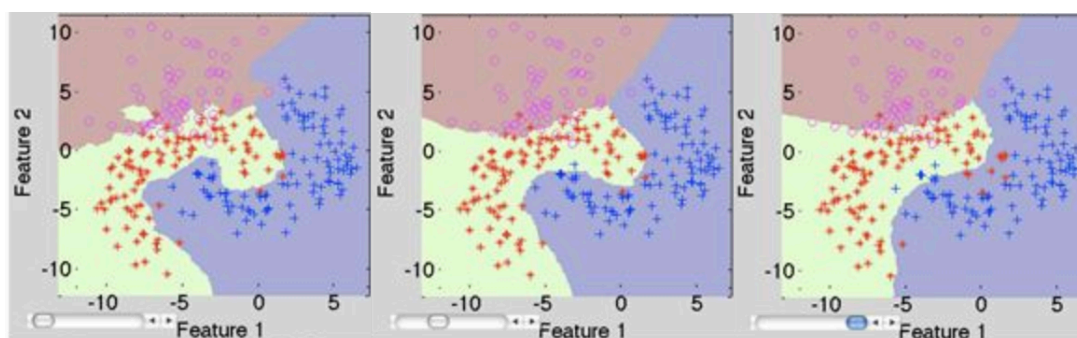
الف- بهترین مقدار K را برای الگوریتم K نزدیک‌ترین همسایه زمانی که از روش LOOCV^۹ استفاده شود را پیدا کنید.

ب- دقت ۱۰ این الگوریتم را برای مقدار K به دست آمده در قسمت الف محاسبه کنید.

۲- افزایش مقدار K در الگوریتم K نزدیک‌ترین همسایه چه تاثیری بر برازش الگوریتم دارد. توضیح دهید.

۳- برای کاهش تاثیر نویز بر دقت الگوریتم چه تغییری در الگوریتم K نزدیک‌ترین همسایه می‌توان اعمال کرد؟ توضیح دهید.

۴- در شکل زیر الگوریتم K نزدیک‌ترین همسایه با مقادیر مختلف K بر روی یک مجموعه داده تست و ناحیه‌بندی شده است. با توجه به مرزهای ناحیه‌بندی در این اشکال، مقادیر K را نسبت به یکدیگر مقایسه نمایید.



⁷ Classification

⁸ K-Nearest Neighbors (KNN-

⁹ Leave One Out Cross Validation

¹⁰ Accuracy

بخش سوم- آشنایی با ابزار وکا (۱۴٪)

۱- ابزار وکا را دانلود نمایید و با استفاده از مجموعه داده iris عملیات زیر را انجام دهید.

الف- ابتدا مجموعه داده iris را از میان مجموعه داده‌های خود وکا بارگزاری نمایید.

ب- سپس از میان ویژگی‌های مختلف این مجموعه داده، sepalwidth و petalwidth را برای دسته‌بندی انتخاب نمایید.

ج- این مجموعه داده را با استفاده از این دو ویژگی نمایش دهید.

د- الگوریتم K نزدیک‌ترین همسایه را در بخش دسته‌بندی انتخاب نمایید و به ازای مقادیر ۱ تا ۵ برای K، دقت الگوریتم را بررسی و گزارش نمایید. لازم به ذکر است که این الگوریتم با عنوان IBK نام‌گذاری شده است.

ه- بار دیگر عملیات الف تا د را برای حالتی که بر روی مجموعه داده ۱۵ درصد نویز اضافه شده است، انجام دهید (اضافه کردن نویز در این ابزار در بخش پیش‌پردازش می‌باشد).

و- دقت الگوریتم به‌دست آمده برای هر دو حالت به ازای مقادیر مختلف K را با یکدیگر مقایسه نمایید و تاثیر پارامتر K در مقابله با نویز را توضیح دهید.

۲- مجموعه داده vote را بارگزاری نموده و به سوالات زیر پاسخ دهید.

الف- با استفاده از الگوریتم درخت تصمیم، این مجموعه داده را دسته‌بندی نمایید و ماتریس پربشانی^{۱۱} آن را نشان دهید.

ب- بار دیگر قسمت الف را در حالتی که از هرس درخت استفاده نمی‌شود انجام دهید.

ج- در هر دو حالت درخت حاصل از دسته‌بندی را نمایش دهید.

د- نتایج حاصل از این دو بخش را با یکدیگر مقایسه نمایید و بیان کنید استفاده از هرس در درخت تصمیم چه سودی دارد.

در انجام تمرینات زیر به نکات زیر توجه نمایید:

- تمرین بدون گزارش فاقد ارزش می‌باشد و نمره‌ای به آن تعلق نمی‌یابد.
- کامنت گذاری کدها در حد لازم و کافی الزامی می‌باشد.
- گذاشتن عنوان^{۱۲} برای نمودارها و برچسب‌گذاری^{۱۳} محورهای نمودار الزامی می‌باشد.
- در پیاده‌سازی بخش‌های مختلف، امکان استفاده از کتابخانه‌های آماده مرتبط با الگوریتم‌های یادگیری ماشین را به طور کلی ندارید. موارد مجاز در صورت سوال بخش‌ها ذکر شده است.
- برای نمایش نمودارها و عملیات ماتریسی می‌توانید از کتابخانه‌های numpy و matplotlib استفاده کنید. همچنین برای خواندن داده‌ها به عنوان ورودی می‌توانید از pandas استفاده کنید.
- برای بهبود سرعت برنامه توصیه می‌شود که تا حد ممکن از عملیات ماتریسی استفاده شود.
- برای محاسبه معیارهای ارزیابی مانند دقت و ماتریس پیریشانی می‌توانید از کتابخانه استفاده نمایید. **در صورت پیاده‌سازی آن‌ها نمره تشویقی تعلق خواهد گرفت. (+۵)**
- تا حد ممکن سعی کنید اصول لازم برای گزارش مهندسی را رعایت نمایید. **در این راستا برای گزارش‌هایی با ظاهر عالی نمره تشویقی در نظر گرفته شده است. (+۵)**
- مطابق قوانین دانشگاه هرگونه کپی‌برداری ممنوع می‌باشد و در صورت مشاهده نمره هر دو طرف صفر در نظر گرفته می‌شود.

۱- هدف از این سوال آشنایی با الگوریتم K نزدیک‌ترین همسایه برای مسایل دسته‌بندی می‌باشد. مجموعه داده مشخص شده^{۱۴} را دانلود نمایید و از آن برای پاسخ دادن به سوالات زیر استفاده نمایید. (۳۰٪)

الف - الگوریتم K نزدیک‌ترین همسایه را با استفاده از معیار فاصله اقلیدسی^{۱۵} پیاده‌سازی نمایید. در این پیاده‌سازی از تمامی ویژگی‌های این مجموعه داده که شامل ۱۳ ویژگی است، استفاده نمایید. این مجموعه داده شامل ۳۰۳ شخص می‌باشد که به دو کلاس بیمار و سالم تقسیم شده‌اند. مقادیر ۱ تا ۷، ۱۰ و ۱۵ را برای این الگوریتم استفاده نمایید و بهترین مقدار K را بدست آورید. لازم به ذکر است که در این بخش داده‌ها را به دو بخش یادگیری^{۱۶} و آزمون^{۱۷} با نسبت ۲ به ۱ تقسیم کنید. در نهایت دقت الگوریتم و ماتریس پیریشانی را برای مجموعه داده آموزش و آزمون گزارش نمایید.

ب - قسمت الف را فقط برای بهترین مقدار K ولی بدون نرمال‌سازی ویژگی‌ها انجام دهید و نتیجه به‌دست‌آمده را با قسمت قبلی مقایسه نمایید. به صورت کلی استفاده از نرمال‌سازی چه تاثیری بر روی دسته‌بندی می‌تواند داشته باشد؟ توضیح دهید.

¹² Title

¹³ Labeling

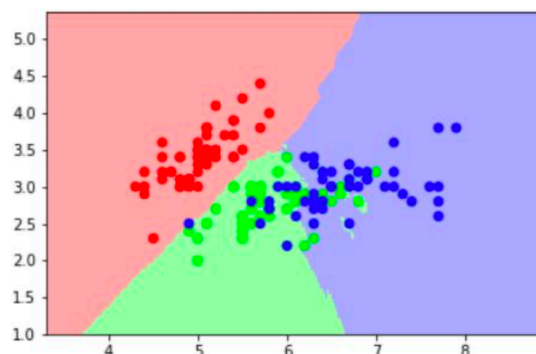
¹⁴ <https://archive.ics.uci.edu/ml/machine-learning-databases/heart-disease/processed.cleveland.data>

¹⁵ Euclidean distance

¹⁶ Train

¹⁷ Test

ج- نمودار ویژگی‌های سن^{۱۸} بر بیشترین ضربان^{۱۹} برای کلاس‌های مختلف رسم نمایید و مرزهای ناحیه‌ها که الگوریتم K نزدیک‌ترین همسایه مشخص می‌نماید را با رنگ‌های مختلف مانند شکل زیر نمایش دهید. در این دسته‌بندی فقط از این دو ویژگی استفاده نمایید (برای این تقسیم‌بندی از مجموعه داده یادگیری استفاده نمایید که با نسبت ۱ به ۲ تقسیم شده‌اند)^{۲۰}.



د- قسمت الف را بار دیگر با معیارهای فاصله زیر انجام دهید.

1. Manhattan
2. Chebyshev

ه- در این قسمت مجموعه داده یادگیری را با استفاده از روش K-fold cross validation دسته‌بندی نمایید و با انتخاب دلخواه K بهترین پارامترهای ممکن (اندازه همسایه‌ها و معیار فاصله) را برای این مجموعه داده بیابید. برای مدل به‌دست‌آمده ماتریس پیرایشی و دقت داده‌های تست را گزارش نمایید.

۲- مجموعه داده mnist را دانلود^{۲۱} و یا اینکه از طریق دستورات زیر آن را بارگزاری نمایید. (۲۵٪)

```
from sklearn import datasets
mnist = datasets.load_digits()
```

الف- این مجموعه داده را به سه دسته آموزش، ارزیابی^{۲۲} و آزمون با نسبت ۳، ۱، ۱ تقسیم نمایید. پس از تقسیم‌بندی از الگوریتم K نزدیک‌ترین همسایه برای دسته‌بندی استفاده نمایید. از مجموعه داده ارزیابی برای یافتن بهترین مقادیر K و بهترین معیار فاصله استفاده نمایید. پس از یافتن بهترین مدل دقت آموزش، ارزیابی و آزمون و ماتریس پیرایشی را برای داده‌های آزمون گزارش کنید. در نهایت ۱۵ نمونه از داده‌های آزمون را به همراه برچسب صحیح و پیش‌بینی شده مانند شکل زیر نمایش دهید.

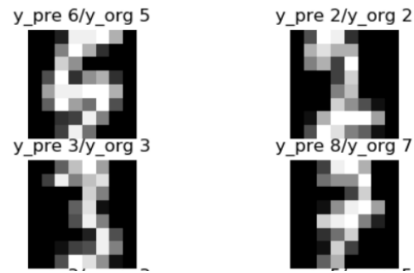
¹⁸ Age

¹⁹ Maximum heart rate

²⁰ https://scikit-learn.org/stable/auto_examples/neighbors/plot_classification.html

²¹ <http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>

²² Evaluation



ب- بار دیگر عملیات قسمت الف را با استفاده از کتابخانه آماده انجام دهید. توصیه می‌شود در این بخش از کتابخانه sklearn استفاده شود.

۳- در این تمرین از الگوریتم K نزدیک‌ترین همسایه برای رگرسیون استفاده نمایید. مجموعه داده regression که در کنار این فایل قرار داده شده است را به دو بخش آموزش و آزمون با نسبت ۲ به ۱ تقسیم نمایید. سپس بهترین مقدار K را با آزمون خطا به دست آورید. خطای MSE را برای این مدل برای هر دو مجموعه داده آزمون و آموزش گزارش کنید. مجموعه داده آموزش و آزمون را با رنگ‌های متفاوت در نموداری نمایش دهید (کمترین خطای MSE برای مجموعه داده آزمون شامل نمره تشویقی خواهد بود). (۵+۱۵٪)

موفق باشید