



بنام خدا
دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



درس سیستم‌های هوشمند

تمرین شماره 2

آبان ۱۴۰۱

فهرست سوالات

- سوال ۱ - درخت تصمیم (تحلیلی) 3
- الف) طراحی طبقه‌بند 3
- ب) آزمون طبقه‌بند 3
- ج) رویکرد حریصانه الگوریتم ID3 4
- د) افزایش قوام طبقه‌بند 4
- سوال ۲ - پیاده‌سازی الگوریتم درخت تصمیم 5
- الف) پیاده‌سازی مدل درخت 5
- ب) بهبود بخشی الگوریتم درخت تصمیم 5
- ج) استفاده از جنگل تصادفی 5
- سوال ۳ - یادگیری بر اساس معیار 6
- طبقه‌بند k همسایه نزدیک 6
- الف) طراحی طبقه‌بند 6
- ب) محاسبه توزیع احتمال تعلق به هر کلاس 6
- یادگیری بر اساس معیار 7
- الف) بررسی کارکرد روش یادگیری 7
- ب) ترسیم دادگان انتقال یافته در فضای جدید 7
- ج) مقایسه عملکرد طبقه‌بند 8
- د) ضریب همبستگی 8
- ه) GMM 9
- نکات تحویل: 9

سوال ۱ - درخت تصمیم (تحلیلی)

در این سوال بررسی می‌کنیم که چگونه دو نوع جاندار را، براساس رنگ، تعداد پا، قد و محل زندگی، می‌توان از هم تشخیص داد. ستون جاندار، برچسبی را که می‌خواهیم پیش بینی کنیم، نشان می‌دهد.

جدول 1-1: اطلاعات مورد نیاز برای طبقه‌بندی دو نوع جاندار (داده‌های آموزش)

شماره	رنگ	تعداد پا	قد	محل زندگی	جاندار
۱	قهوه‌ای	2	بلند	خشکی	A
۲	قهوه‌ای	3	کوتاه	خشکی	B
۳	سبز	2	بلند	آب	B
۴	سبز	3	بلند	آب	B
۵	قهوه‌ای	2	کوتاه	آب	A
۶	قهوه‌ای	2	بلند	آب	A
۷	قهوه‌ای	2	کوتاه	خشکی	B
۸	سبز	2	کوتاه	آب	A
۹	سبز	3	بلند	آب	B
۱۰	قهوه‌ای	2	بلند	خشکی	A

الف) طراحی طبقه‌بند

با استفاده از جدول ۱-۱، یک طبقه‌بند درخت تصمیم^۱ برای تشخیص نوع جاندار (A یا B)، بر مبنای بهره‌ی اطلاعات^۲ و با الگوریتم ID3 را آموزش دهید.

ب) آزمون طبقه‌بند

با استفاده از طبقه‌بند قسمت الف، نوع جاندار هر یک از نمونه‌های زیر را (جدول ۱-۲) مشخص کرده و عملکرد مدل را به کمک ماتریس آشفتگی^۳ بررسی کنید.

^۱ Decision Tree

^۲ Information Gain

^۳ Confusion Matrix

	A	B
A	2	1
B	1	2

جدول 1-2: اطلاعات مورد نیاز برای طبقه‌بندی دو نوع جاندار (داده‌های آزمون)

شماره	رنگ	تعداد پا	قد	محل زندگی	جاندار
۱	قهوه‌ای	3	بلند	خشکی	B
۲	سبز	2	بلند	خشکی	A
۳	سبز	2	کوتاه	خشکی	A
۴	قهوه‌ای	2	کوتاه	آب	B
۵	قهوه‌ای	2	بلند	خشکی	A

ج) رویکرد حریصانه^۴ الگوریتم ID3

تمام شرایط لازم برای اینکه جاندار از نوع A باشد و یا از نوع B باشد را در نظر بگیرید. (برای مثال اگر تعداد پاها، ۳ باشد، جاندار از نوع B است). در هر کدام از این شروط، حداکثر از ۴ ویژگی استفاده شده است. آیا می‌توانید، درخت تصمیم جدیدی طراحی کنید که فقط با استفاده از ۲ ویژگی بتواند نوع جاندار را تشخیص دهد و هم چنان باعث صفر شدن خطا در مجموعه آموزشی شود؟ (به این معنا که درخت تصمیم جدید هم چنان برای تمامی داده‌های آموزش صدق کند) جواب خود را توجیه کنید.

د) افزایش قوام طبقه‌بند

چرا طبقه‌بندهای درخت تصمیم در برابر بیش‌برازش^۵ مقاوم^۶ نیستند؟ دو روش برای جلوگیری از این مشکل ارائه دهید.

^۴ Greedy approach

^۵ Overfitting

^۶ Robust

سوال ۲ - پیاده‌سازی الگوریتم درخت تصمیم

در این بخش، با استفاده از داده‌های پیوست شده هدف آن است که الگوریتم درخت تصمیم را بدون استفاده از کتابخانه‌های آماده، پیاده‌سازی کنیم.

الف) پیاده‌سازی مدل درخت

در ابتدا، به داده‌های مورد نیاز دسترسی پیدا می‌کنیم:

```
Import pandas as pd
Train = pd.read_csv("titanic-train.csv")
Test = pd.read_csv("titanic-test.csv")
```

هدف آن است که بتوانیم با استفاده از ویژگی‌های افرادی که سوار بر کشتی تایتانیک بودند، نجات یافتن یا نیافتن هر کدام را پیش‌بینی کنیم. به نظر شما، چه ویژگی‌هایی می‌توانند در این امر تأثیرگذار باشند؟ (مشخصاً، اسم فرد بی‌تأثیر خواهد بود، اما سن او می‌تواند اطلاع مفیدی باشد). با توجه به این دید، در ابتدا، پیش‌پردازش‌های لازم را بر روی این داده‌ها انجام دهید و هر کدام را توضیح دهید.

(** برای پیش‌پردازش‌های لازم، می‌توانیم از این لینک، کمک بگیرید)

عمق درخت را (تعداد ویژگی‌های استفاده شده برای پیش‌بینی) در ابتدا ۳ در نظر بگیرید.

(**دقت شود که در مراحل گوناگون، می‌توان از یک ویژگی چند مرتبه استفاده کرد)

معیار انتخاب ویژگی برتر را به دلخواه، از میان معیارهای معرفی شده در کلاس انتخاب نمایید. در هر مرحله، ماتریس آشفتگی^۷ را نمایش دهید و همچنین دقت مدل خود را گزارش کنید. با افزایش عمق درخت، دقت طبقه‌بند چه تغییری می‌کند؟ توضیح دهید.

ب) بهبود بخشی الگوریتم درخت تصمیم

ایرادات الگوریتم درخت تصمیم را بیان نمایید و توضیح دهید چگونه می‌توان عملکرد این الگوریتم را با استفاده از روش‌های Bagging و جنگل تصادفی^۸ بهبود دهیم.

ج) استفاده از جنگل تصادفی

حال، با استفاده از الگوریتم جنگل تصادفی، سعی بر آن است که دقت طبقه‌بند افزایش یابد. در این راستا، تعدادی درخت تصمیم (حداقل ۳ درخت) که هر کدام بر اساس ویژگی‌های تصادفی بر روی تعدادی داده‌های تصادفی پیاده می‌شوند. با استفاده از رأی اکثریت^۹ بایستی پیش‌بینی مورد نظر را اعلام کنیم. حال دقت و ماتریس آشفتگی را گزارش دهید.

^۷ Confusion Matrix

^۸ Random Forest

^۹ Majority Voting

سوال ۳ - یادگیری بر اساس معیار

در این سوال قرار است که تاثیر روش های یادگیری بر اساس معیار را بر طبقه بند k همسایه نزدیک بررسی کنیم. در ابتدا به کمک دستور زیر، دیتاست “wine” را به کمک کتابخانه Scikit-learn بخوانید:

```
from sklearn.datasets import load_wine
data = load_wine()
```

** (در تمامی مراحل سوال نیاز به استانداردسازی 10 دادگان نمی باشد) **

طبقه بند k همسایه نزدیک

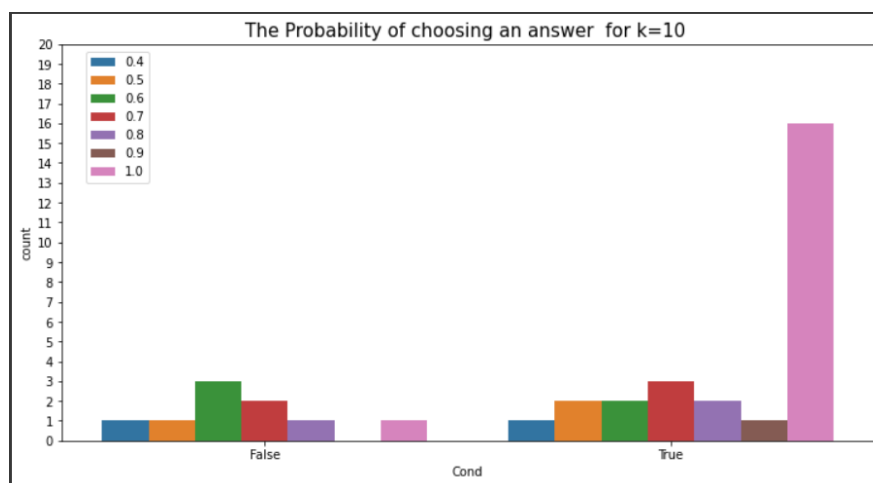
** (در این قسمت مجاز به استفاده از کتابخانه Scikit-learn نیستید و باید تمامی پیاده سازی را به کمک کتابخانه NumPy انجام دهید) **

الف) طراحی طبقه بند

در ابتدا 20 درصد دادگان را به دادگان آزمون و بقیه را به دادگان آموزش اختصاص دهید. به ازای تعداد همسایه های 10، 5، 1 و 20 دقت طبقه بند و ماترس آشفتگی را بر روی دادگان آزمون گزارش دهید.

ب) محاسبه توزیع احتمال تعلق به هر کلاس

برای تعداد همسایه های ذکر شده در بخش قبل، احتمال تعلق دادگان آزمون را به هر کلاس مشابه شکل آمده در زیر رسم کنید. با تغییر تعداد همسایه ها، نحوه تغییر در توزیع احتمال ها را بررسی کنید. برای کدام مقدار همسایه فکر می کنید مدل بهتر عمل کرده است؟ معیار خود را برای این انتخاب توضیح دهید.



شکل ۳-۱: نمودار توزیع احتمالی تعلق به هر کلاس

¹⁰ Normalization

یادگیری بر اساس معیار

*** (در این قسمت می‌توانید از کتابخانه‌های آماده Scikit-learn و metric-learning استفاده کنید)***

به کمک دو یادگیری بر اساس معیار $LMNN^{11}$ و $LFDA^{12}$:

الف) بررسی کارکرد روش یادگیری

هدف تعریف هر یک از یادگیری‌ها، قیود¹³ تعریف شده در فرآیند یادگیری و علت وجود هر یک از قید ها را توضیح دهید.

ب) ترسیم دادگان انتقال یافته در فضای جدید

1. در هر یک از دو یادگیری بر اساس معیار، پارامتری بنام k وجود دارد. بنظر شما کارکرد این پارامتر چیست و چه تفاوتی با پارامتر k در طبقه بند k همسایه نزدیک دارد ؟

2. در این قسمت می‌خواهیم تاثیر یادگیری بر اساس معیار در افراز داده ها در 2 بعد ببینیم. از آنجایی داده اولیه ما دارای 13 بعد می‌باشد، نمایش آن در فضای دو بعدی امکان پذیر نیست. بدین منظور دو راه کار در پیش دارید:

➤ به کمک کتابخانه metric-learning (راه کار پیشنهادی) : پارامتری را در فراخوان توابع مربوطه پیدا کنید که به کمک آن بتوان دادگان را به فضای با بعد پایین تر انتقال داد.

➤ به کمک روش های کاهش بعد در کتابخانه Scikit-learn: می‌توانید بعد از انتقال دادگان به فضای جدید به کمک روش های کاهش بعد مانند PCA^{14} ، دادگان را به فضای با بعد پایین تر انتقال داد.

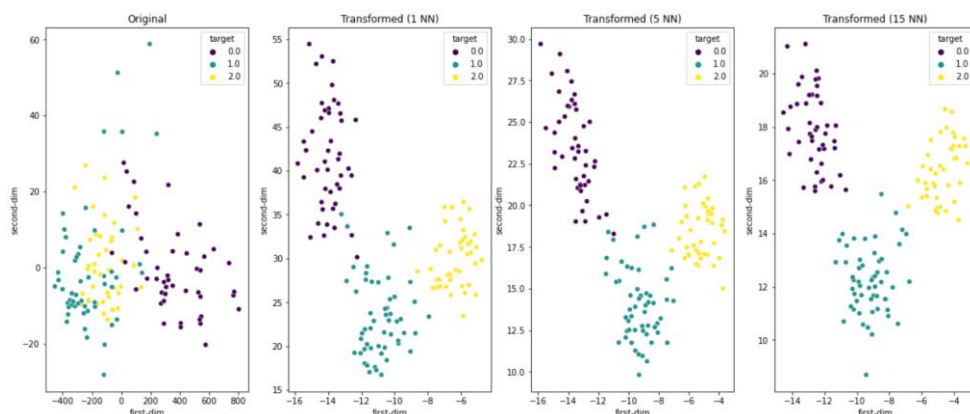
بعد از کاهش بعد دادگان اصلی و انتقال یافته در فضای جدید، به ازای 3 مقدار مختلف برای پارامتر k (1,5,15) نحوه افراز دادگان هر کلاس با کلاس خود و کلاس های دیگر را رسم و تحلیل کنید. برای کدام مقدار k ، دادگان در فضای جدید قابلیت تفکیک پذیری بیشتری دارند ؟ چرا ؟ *** (نمودار مشابه در شکل زیر آورده شده است).***

¹¹ Largest Margin Nearest Neighbor

¹² Local Fisher Discriminant analysis

¹³ Constraints

¹⁴ Principal analysis component



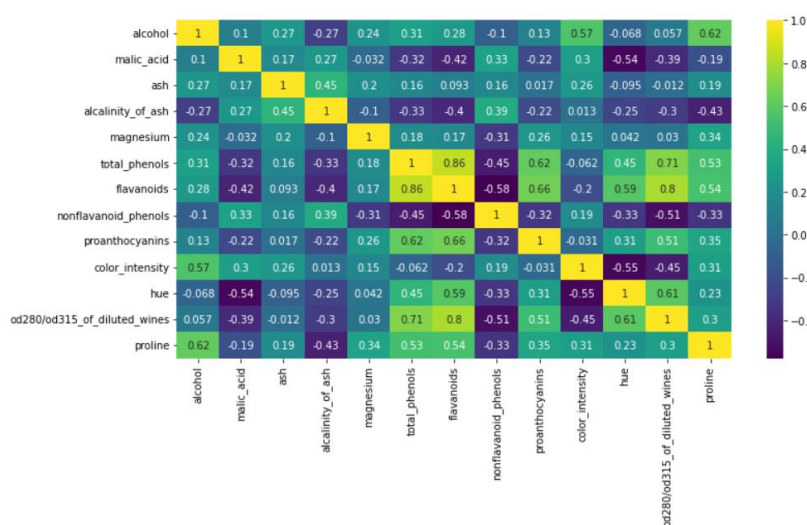
شکل ۳-۲: افراز دادگان اصلی و انتقال یافته به ازای مقادیر مختلف k

ج) مقایسه عملکرد طبقه بند

برای بهترین مقدار k بدست آمده از قسمت قبل، دقت و ماتریس آشفته‌گی طبقه بند را این بار برای دادگان انتقال یافته در فضای جدید به ازای تعداد همسایه مشابه بخش الف در قسمت طبقه بند k همسایه بدست آورده و مقایسه کنید.

د) ضریب همبستگی

یکی از اطلاعات مفیدی که میتوان همواره از دادگان استخراج کرد، همبستگی بین ستون‌های ویژگی^{۱۵} می‌باشد. بدین صورت که میتوانیم ضریب همبستگی بین هر دو ستون ویژگی از دیتاست خود را داشته باشیم. این اطلاعات از این جهت سودمند است که می‌توانیم تاثیر متقابل ستون‌های ویژگی را در فرآیند یادگیری بیشتر درک کنیم. در کتابخانه Pandas، میتوانید به کمک دستور `corr()` همبستگی دو به دو بین ستون‌های ویژگی را در قالب یک آرایه دو بعدی بدست آورید و سپس رسم کنید که نمونه‌ای از آن در پایین آورده شده است :



شکل ۳-۳: نمایش ماتریس همبستگی یک دیتاسیت دلخواه

¹⁵ Feature set

برای هر کدام از دو روش یادگیری بر اساس معیار، ماتریس های همبستگی را بررسی کنید. ستون های ویژگی در فضای انتقال یافته برای هر کدام از روش ها چه ویژگی متمایزی دارند. ستون های ویژگی بدست آمده در روش LMNN، چه اطلاعات مهمی را در فضای جدید آشکار می کنند ؟

ه) GMML

در [این مقاله](#)، روش GMML به عنوان روشی جدید در یادگیری بر اساس معیار معرفی شده است. مدل سازی مساله را توضیح دهید. روش ارائه شده چه تفاوتی با روش LMNN مطرح شده در ابتدا دارد ؟

نکات تحویل:

- مهلت تحویل این تمرین 6 آذر میباشد.
- انجام این تمرین به صورت یک نفره است.
- برای انجام این تمرین تنها مجاز به استفاده از زبان برنامه نویسی پایتون هستید.
- در صورت وجود تقلب نمره تمامی افراد شرکت کننده در آن -۱۰۰ لحاظ میشود.
- لطفا پاسخ تمرین خود را (به همراه کد/گزارش سوال کامپیوتری) به صورت زیر در صفحه درس آپلود نمایید:

HW [HW number] _ [Last name] _ [Student number].zip

- در صورت وجود هر گونه ابهام یا مشکل میتوانید از طریق ایمیل با مسئولان حل تمرین در تماس باشید:

مسئول تمرین سوال ۱: عاطفه ملاحقر (ut.ac.ir@ati.mollabagher)

مسئول تمرین سوال ۲: دریا افزلی (darya.afzali@ut.ac.ir)

مسئول تمرین سوال ۳: شایان واصف (sh.vassef@ut.ac.ir)