

به نام خدا



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)



DEPARTMENT OF COMPUTER
ENGINEERING AND IT

دانشگاه صنعتی امیرکبیر

دانشکده‌ی مهندسی کامپیوتر

تمرین سری دوم یادگیری ماشین

دکتر احسان ناظر فرد

طراح سوال:

محمد رضا امامی ناصری

سید اردلان قریشی

پاییز ۱۳۹۹

توضیحات مهم:

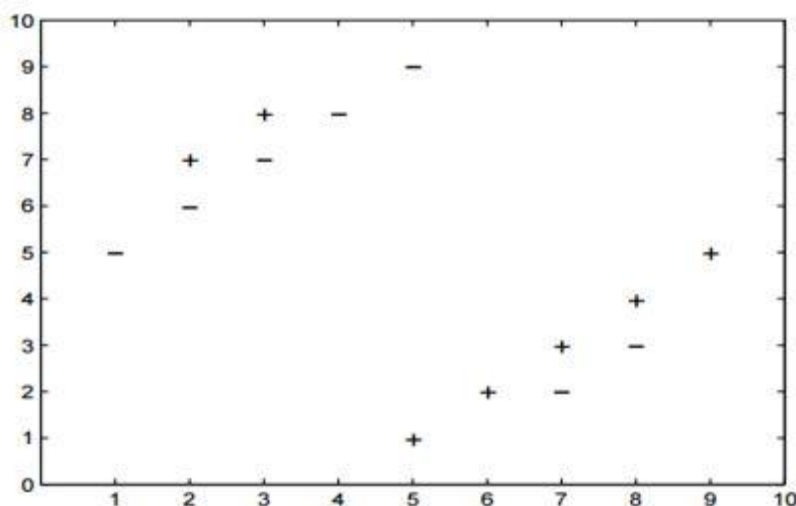
- تمامی مستندات خود شامل گزارش و کدهای خود را در یک فایل فشرده با فرمت zip ذخیره کرده و با عنوان #StudentId_HW2.zip بارگذاری نمایید (به عنوان مثال 99131000_HW2.zip).
- مهلت انجام تمرین تا ساعت ۲۳:۵۵ روز چهارشنبه مورخ ۲۶ آذر می باشد و به هیچ وجه تمدید نمی شود.
- تمرین بدون گزارش فاقد ارزش می باشد و نمره ای به آن تعلق نمی یابد.
- تا حد ممکن سعی کنید اصول لازم برای گزارش مهندسی را رعایت نمایید (به بهترین گزارش نمره تشویقی تعلق می گیرد).
- مطابق قوانین دانشگاه هرگونه کپی برداری ممنوع می باشد و در صورت مشاهده نمره ای هر دو طرف صفر در نظر گرفته می شود.
- شما مجاز هستید برای تمامی تمرین ها ۷ روز در کل و با سقف حداکثر ۳ روز برای هر تمرین، تاخیر بدون کسر نمره داشته باشید. به ازای هر روز تاخیر بیشتر، ۱۰٪ از نمره ای تمرین مربوطه کسر می شود.
- در صورت داشتن هرگونه ابهام می توانید از طریق ایمیل زیر سوال خود را مطرح نمایید:

MLAUTFALL99@gmail.com

سوالات تشریحی

۱- برای یافتن بهترین مقدار پارامتر K در الگوریتم K -نزدیک‌ترین همسایه^۱، چه راهکاری را پیشنهاد می‌کنید؟^۳

۲- با توجه به شکل زیر به سوالات پاسخ دهید. ۸٪



الف) بهترین مقدار K برای الگوریتم K -نزدیک‌ترین همسایه زمانی که از روش LOOCV^۲ استفاده می‌شود را محاسبه کنید. دقت^۳ الگوریتم را به ازای این K گزارش نمایید.

ب) مشکل انتخاب مقدار K خیلی بزرگ و خیلی کوچک چیست؟ توضیح دهید.

ج) نقطه (2,1) با توجه به K به دست آمده به کدام کلاس تعلق می‌یابد؟

۳- هر کدام از الگوریتم‌های K -نزدیک‌ترین همسایه و درخت تصمیم^۴ را از نظر پارامتریک و غیر پارامتریک^۵ بودن بررسی کنید. ۲٪

۴- هر کدام از الگوریتم‌های K -نزدیک‌ترین همسایه و درخت تصمیم را از Generative و Discriminative بودن بررسی کنید. ۲٪

¹ K-Nearest Neighbors algorithm (KNN)

² Leave One Out Cross Validation

³ accuracy

⁴ decision tree

⁵ parametric and nonparametric

۵- به چه الگوریتم‌هایی تنبل^۶ گفته می‌شود؟ K-نزدیک‌ترین همسایه تنبل است یا خیر؟ توضیح دهید. ۲٪

۶- هرس^۷ درخت تصمیم چه تاثیری بر بیش‌برازش^۸ دارد؟ این هرس چه زمانی باید انجام شود؟ توضیح دهید. ۴٪

۷- در جدول ۱ مجموعه داده‌ای نمایش داده شده است که در آن افراد با توجه به ویژگی‌هایی مثل سن، درآمد و... اقدام به خرید یا عدم خرید یک کالا کرده‌اند. هدف ما تخمین^۹ این است که آیا فرد مورد نظر قصد خرید کالا را دارد یا خیر. ۱۰٪

age	income	student	credit	Buy
youth	high	no	fair	-
youth	high	no	excellent	-
middle	high	no	fair	+
senior	medium	no	fair	+
senior	low	yes	fair	+
senior	low	yes	excellent	-
middle	low	yes	excellent	+
youth	medium	no	fair	-
youth	low	yes	fair	+
senior	medium	yes	fair	+
youth	medium	yes	excellent	+
middle	medium	no	excellent	+
middle	high	yes	fair	+
senior	medium	no	excellent	-

جدول ۱

الف- با توجه به ویژگی آن‌تروپی^{۱۰} و بهره اطلاعات^{۱۱} درخت تصمیم بهینه^{۱۲} را برای این مجموعه داده بیابید.

ب- داده‌های زیر را با توجه به درخت تصمیم به دست آمده، دسته‌بندی^{۱۳} کنید.

$X_1 = (\text{age} = \text{youth}, \text{income} = \text{high}, \text{student} = \text{yes}, \text{credit} = \text{fair})$

$X_2 = (\text{age} = \text{senior}, \text{income} = \text{low}, \text{student} = \text{no}, \text{credit} = \text{excellent})$

$X_3 = (\text{age} = \text{middle-aged}, \text{income} = \text{medium}, \text{student} = \text{no}, \text{credit} = \text{fair})$

⁶ lazy

⁷ pruning

⁸ overfitting

⁹ predict

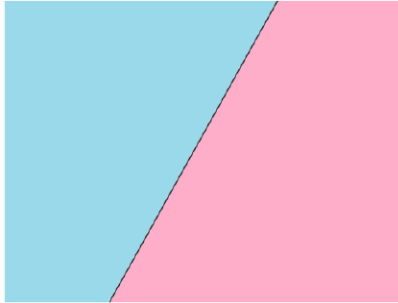
¹⁰ entropy

¹¹ information gain

¹² optimal

¹³ classify

۸- در دسته‌بندی^{۱۴} داده‌ها با درخت تصمیم می‌توانیم با افزایش عمقِ درخت^{۱۵} توابع بسیار پیچیده‌ای بسازیم. آیا برای مسائل جداپذیر خطی^{۱۶} مانند شکل زیر می‌توان دسته‌بندی را با عمق محدود درخت انجام داد؟ ۴٪



۹- الگوریتم جنگل تصادفی^{۱۷} را مختصر توضیح دهید. چرا این الگوریتم با محدود نگه داشتن عمق درخت می‌تواند داده‌ها را جدا کند؟ ۵٪

¹⁴ classification

¹⁵ tree-depth

¹⁶ linearly separable

¹⁷ random forest

هدف از این بخش آشنایی با ابزار شناخته‌شده‌ی وکا است. برای انجام تمرینات این بخش می‌توانید از آموزش‌های موجود در اینترنت کمک بگیرید.

(به عنوان نمونه: <https://www.tutorialspoint.com/weka/index.htm>)

۱- مجموعه‌داده‌ی پیوست شده به نام breast-cancer.arff را بارگذاری کرده و با استفاده از درخت تصمیم (J48) داده‌ها را دسته‌بندی کنید. (۷۰ درصد داده‌ها را به عنوان داده‌های آموزش^{۱۸} و ۳۰ درصد داده‌ها را به عنوان داده‌ی آزمون^{۱۹} در نظر بگیرید). ۴٪

خروجی مورد نظر: شکل درخت تصمیم نهایی، ماتریس درهم ریختگی^{۲۰}، مقادیر Precision, FN, FP, TP, Recall, F1measure و تحلیل لازم.

۲- پارامتر unpruned چه چیزی را کنترل می‌کند؟ این مقدار را از false به true تغییر داده و دوباره موارد قسمت ۱ را انجام دهید. درخت آموزش داده شده در این قسمت چه تفاوتی با قسمت ۱ دارد؟ توضیح دهید. ۴٪

خروجی مورد نظر: شکل درخت تصمیم نهایی، ماتریس درهم ریختگی، مقادیر Precision, FN, FP, TP, Recall, F1measure و تحلیل لازم.

۳- حال موارد قسمت ۱ و ۲ را بار دیگر، این بار با اعمال ۱۵ درصد نویز^{۲۱} به ریشه‌ی درخت‌های دو قسمت قبل تکرار کنید و نتایج حاصل را با نتایج مراحل قبلی مقایسه کرده و بررسی کنید که هرس کردن درخت چه تاثیری در برخورد با نویز دارد. ۲٪

خروجی مورد نظر: ماتریس درهم ریختگی، مقادیر Precision, FN, FP, TP, Recall, F1measure و تحلیل لازم.

¹⁸ train
¹⁹ test

²⁰ confusion matrix
²¹ noise

سوالات پیاده‌سازی

توضیحات مهم:

- در روند اجرا انتخاب مقادیر برای تقسیم داده‌ها به مجموعه آموزش، ارزیابی و... به عهده دانشجو می‌باشد.
- حتما پارامترهای انتخاب شده برای برنامه خود و هرگونه شرایطی که در نظر گرفته‌اید را در گزارش خود بیاورید.
- برای بهبود سرعت برنامه توصیه می‌شود از عملیات ماتریسی استفاده کنید.
- در هر مرحله، نتایج خود را تحلیل کنید.
- کدهای خود را برای خوانایی بیشتر **کامنت گذاری** کنید.
- در تمامی سوال‌ها تنها مجاز به استفاده از کتابخانه‌های `numpy`، `matplotlib` و `pandas` می‌باشید.
- در پیاده‌سازی بخش‌های مختلف، امکان استفاده از کتابخانه‌های آماده مرتبط با الگوریتم‌های یادگیری ماشین را به طور کلی ندارید. **موارد مجاز در صورت سوال ذکر شده است.**
- گذاشتن عنوان برای نمودارها و برچسب گذاری محورهای نمودار الزامی می‌باشد.

در این بخش برای دسته‌بندی از مجموعه داده‌های مرتبط با بیماری سرطان پستان^{۲۲} استفاده شده است. دلیل این امر اول مواجه شدن عزیزان دو طراح این سری تمرین با این بیماری شوم و دوم نشان دادن یکی از هزاران کاربرد الگوریتم‌های دسته‌بندی در دنیای واقعی بوده است. با توجه به اینکه این بیماری در صورت تشخیص سریع می‌تواند درمان بسیار بهتری داشته باشد، قسمتی از پایان‌نامه‌ی مادرم که خود در دوران دست و پنجه نرم کردن با این بیماری اقدام به نوشتن و دفاع از آن کرده است را به امید اطلاع‌رسانی هرچند کوچک آورده‌ام. سید اردلان قریشی

سرطان پستان اولین سرطان شایع و دومین علت مرگ ناشی از سرطان در بین زنان ۳۵ تا ۵۵ ساله ایران بوده و سالانه حدود ۷ هزار مورد جدید به بیماران قبلی اضافه می‌شود. از هر ۸ زن ایرانی در فاصله سنی ۴۵ تا ۵۵ سال، یک نفر شانس ابتلا به این سرطان را دارد، در ایران سن ابتلا در زنان در سال‌های اخیر کاهش یافته است، به عبارتی در دنیا مادر بزرگ‌ها به این بیماری مبتلا می‌شوند و در ایران مادران!

این بیماری تأثیرات شدید روانی ایجاد می‌کند و تشخیص و درمان آن، تجربه‌ای همراه با استرس و اضطراب می‌باشد. با اینکه جراحی رایج‌ترین درمان سرطان می‌باشد، اما این بیماران پس از جراحی از یک سو با مشکلاتی از جانب سرطان و جراحی و از سوی دیگر با درمان‌هایی مثل شیمی‌درمانی و عوارض جانبی ناخوشایندی مثل ریزش مو، حالت تهوع، اِدم لنفاوی^{۲۳} و سایر مشکلات مواجه می‌شوند. درمان‌های طولانی، توانایی زنان را در برقراری نقش اجتماعی به عنوان زن خانه‌دار یا شاغل، مورد تهدید قرار می‌دهند. سطح بالای استرس، تأثیر منفی طولانی مدت بر خود باوری زنان دارد که در نهایت تأثیر نامطلوبی بر عملکرد خانوادگی و نقش زناشویی گذاشته و همچنین منجر به پایین آمدن سطح کیفی زندگی آنان می‌شود. (محمدی گل، ۱۳۹۵).

توصیه می‌شود که معاینه‌ی شخصی و چک کردن علایم از سن ۲۰ سالگی به صورت ماهانه توسط هر فرد انجام شود. علاوه بر آن از سن ۲۰ تا ۳۹ سالگی هر ۳ سال این معاینه توسط پزشک نیز انجام گردد و از ۴۰ سالگی به بعد به طور سالانه معاینه توسط پزشک انجام شود. (همایی فاطمه، ۱۳۹۰).

²² breast cancer

²³ lymphedema

۱- یک تابع^{۲۴} بنویسید که با گرفتن ورودی^{۲۵}های مجموعه داده، K و معیار فاصله^{۲۶}، الگوریتم KNN را اجرا کند. از آن تابع برای دسته‌بندی مجموعه داده‌ی mammographic_masses.data استفاده کنید. ۲۵٪

- برای ارزیابی^{۲۷} استفاده از 10-fold cross validation الزامی است.
- در این مجموعه داده برخی نمونه^{۲۸}ها دارای ویژگی‌هایی با مقادیر نامعلوم^{۲۹} می‌باشند. به ابتکار خود روشی را برای حل این مشکل بیابید و روش خود را در گزارش شرح دهید.
- پیاده‌سازی ماتریس درهم ریختگی و k-fold cross validation الزامی نیست و می‌توانید از کتابخانه‌ی آماده برای این دو مورد استفاد کنید. پیاده‌سازی این موارد توسط خود دانشجو نمره تشویقی دارد. (در گزارش ذکر شود).
- پیش پردازش^{۳۰} لازم را بر روی مجموعه داده انجام داده و در گزارش خود بیاورید.

الف) الگوریتم KNN را به ازای مقادیر مختلف ۱،۳،۵،۷،۱۵،۳۰ برای K و با فاصله اقلیدسی^{۳۱} اجرا کرده و تاثیر مقادیر مختلف K را تحلیل کنید.

خروجی مورد نظر: دقت الگوریتم و ماتریس درهم ریختگی

ب) به ازای بهترین مقدار K که در قسمت الف یافته‌اید و با فاصله‌های اقلیدسی، منهتن^{۳۲} و کسینوسی^{۳۳}، الگوریتم را اجرا کنید.

خروجی مورد نظر: دقت الگوریتم و ماتریس درهم ریختگی

ج) با استفاده از کتابخانه‌های آماده و به ازای مقادیر مختلف ۱،۳،۵،۷،۱۵،۳۰ برای K و فاصله اقلیدسی، داده‌ها را دسته‌بندی کرده و پیاده‌سازی خود را از نظر دقت و سرعت با این کتابخانه مقایسه کنید.

خروجی مورد نظر: دقت الگوریتم، ماتریس درهم ریختگی و زمان اجرای هر الگوریتم

۲- در این قسمت از الگوریتم KNN برای رگرسیون استفاده نمایید. مجموعه داده‌ی regression.xlsx را با نسبت ۷۰ به ۳۰ تقسیم کرده و سپس بهترین مقدار برای K را با آزمون و خطا بیابید. خطای MSE را برای این مدل برای هر دو مجموعه‌ی داده آزمون و آموزش گزارش کنید. ۱۵٪

خروجی مورد نظر: خطای MSE برای هر دو مجموعه آزمون و آموزش

²⁴ function
²⁵ input
²⁶ distance criterion
²⁷ evaluation
²⁸ instance

²⁹ missing value
³⁰ preprocessing
³¹ Euclidean distance
³² Manhattan distance
³³ cosine distance/similarity

۳- درخت تصمیم بهینه را با استفاده از کتابخانه‌های آماده برای مجموعه داده‌ی breast-cancer-wisconsin- train.data آموزش داده و مجموعه داده‌ی breast-cancer-wisconsin-test.data را دسته‌بندی کنید. ۱۰٪

- در این مجموعه داده برخی نمونه‌ها دارای ویژگی‌هایی با مقادیر نامعلوم می‌باشند. به ابتکار خود روشی را برای حل این مشکل بیابید و روش خود را در گزارش شرح دهید.
- پیش پردازش‌های لازم را بر روی مجموعه داده انجام داده و در گزارش خود بیاورید.

خروجی مورد نظر: دقت الگوریتم برای مجموعه داده آزمون و ماتریس درهم ریختگی

با آرزوی موفقیت!