



به نام خدا



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده مهندسی برق و کامپیوتر

سیستم‌های هوشمند

تمرین شماره 2

نیمه زمان پور

810198407

آذر ماه 1401

فهرست سوالات

سوال 1	4
الف)	4
ب)	9
ج)	10
د)	10
سوال 2	12
الف)	12
ب)	13
ج)	14
سوال 3	15
طبقه بند k همسایه نزدیک	15
الف)	15
ب)	16
یادگیری بر اساس معیار	20
الف)	20
ب)	22
ج)	23
د)	24
ه)	26
پیوست	27
سوال 2) توضیح کد	27

سوال 3) توضیح کد 28.

منابع: 29.



سوال 1

(الف)

ساختن درخت تصمیم را با انتخاب یک ویژگی از بین 4 ویژگی مورد برای نظر انتخاب می‌کنیم. و برای هر کدام information gain را حساب می‌کنیم.

$$Entropy = E(S) = - \sum p_i \log_2 p_i$$

$$IG(S, A) = Entropy(S) - \sum_{v \in values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

انتخاب بر اساس محل زندگی:

	$S = [+5, -5]$ $E = 1$	
	محل زندگی	
خشکی		آب
$S = [+2, -2]$ $E = 1$		$S = [+3, -3]$ $E = 1$
$IG(S, habitat) = 1 - 0.4 * 1 - 0.6 * 1 = 0$		

انتخاب بر اساس قد:

	$S = [+5, -5]$ $E = 1$	
	قد	
کوتاه		بلند
$S = [+2, -2]$ $E = 1$		$S = [+3, -3]$ $E = 1$
$IG(S, habitat) = 1 - 0.4 * 1 - 0.6 * 1 = 0$		

انتخاب بر اساس تعداد پا:

	$S = [+5, -5]$ $E = 1$	
	تعداد پا	
3		2
$S = [0, -3]$ $E = 0$		$S = [+5, -2]$ $E = 0.863$
$IG(S, habitat) = 1 - 0.3 * 0 - 0.7 * 0.863$ $= 0.396$		

انتخاب بر اساس رنگ

	$S = [+5, -5]$ $E = 1$	
	رنگ	
سبز		قهوه‌ای
$S = [1, -3]$ $E = 0.811$		$S = [+4, -2]$ $E = 0.811$
$IG(S, habitat) = 1 - 0.4 * 0.811 - 0.6 * 0.811$ $= 0.125$		

از بین مقادیر information gain برای انتخاب ویژگی ریشه درخت. 0، 0.396، 0.082 بیشترین را برمی‌گزینیم. که تعداد پا می‌باشد.

حال در این ریشه به سراغ شاخه سمت چپ می‌رویم. در این شاخه تمام نمونه‌ها B هستند. پس این node تبدیل به برگ می‌شود.

در شاخه سمت راست حالا به سراغ انتخاب بهترین attribute برای شاخه می‌رویم:

انتخاب بر اساس محل زندگی:

	$S = [+5, -2]$ $E = 0.863$	
--	-------------------------------	--

	محل زندگی	
خشکی		آب
$S = [2, -1]$ $E = 0.918$		$S = [+3, -1]$ $E = 0.811$
$IG(S, \text{habitat}) = 0.863 - \frac{3}{7} * 0.918 - \frac{4}{7} * 0.811 = 0.006$		

انتخاب بر اساس محل قد:

	$S = [+5, -2]$ $E = 0.863$	
	قد	
کوتاه		بلند
$S = [2, -1]$ $E = 0.918$		$S = [+3, -1]$ $E = 0.811$
$IG(S, \text{habitat}) = 0.863 - \frac{3}{7} * 0.918 - \frac{4}{7} * 0.811 = 0.006$		

انتخاب بر اساس رنگ:

	$S = [+5, -2]$ $E = 0.863$	
	رنگ	
سبز		قهوه‌ای
$S = [1, -1]$ $E = 1$		$S = [+4, -1]$ $E = 0.722$
$IG(S, \text{habitat}) = 0.863 - \frac{2}{7} * 1 - \frac{5}{7} * 0.722 = 0.062$		

مقدار information gain ویژگی رنگ از همه بیشتر شد. پس شاخه بعدی ویژگی رنگ است.

حالا در شاخه سمت چپ به سراغ attribute بعدی می‌رویم.

انتخاب بر اساس قد

	$S = [1, -1]$ $E = 1$	
	فد	
کوتاه		بلند
$S = [+1, 0]$ $E = 0$		$S = [0, -1]$ $E = 0$
$IG(S, habitat) = 1 - 0.5 * 0 - 0.5 * 0 = 1$		

انتخاب بر اساس محل زندگی

	$S = [1, -1]$ $E = 1$	
	محل زندگی	
خشکی		آب
$S = [0, 0]$ $E = 0$		$S = [+1, -1]$ $E = 1$
$IG(S, habitat) = 1 - 1 * 1 = 0$		

واضا ویژگی قد به راحتی داده ها را تقسیم کرده و 2 برگ به ما می دهد.

حالا به سراغ شاخه سمت راست رنگ می رویم.

انتخاب بر اساس قد

	$S = [4, -1]$ $E = 0.722$	
	فد	
کوتاه		بلند
$S = [+1, -1]$ $E = 1$		$S = [+3, 0]$ $E = 0$
$IG(S, habitat) = 0.722 * 0.6 * 0 - 0.4 * 1 = 0.322$		

انتخاب بر اساس محل زندگی

	$S = [4, -1]$ $E = 0.722$	
	محل زندگی	
خشکی		آب
$S = [2, -1]$ $E = 0.918$		$S = [+2, 0]$ $E = 0$
$IG(S, habitat) = 0.722 - 0.6 * 0.918 - 0.4 * 0 = 0.015$		

ویژگی قد را به عنوان شاخه بعدی انتخاب می‌نماییم. شاخه سمت راست تبدیل به برگ می‌شود. چون فقط یک نوع داده دارد. حالا آخرین attribute را نیز به عنوان شاخه برای شاخه سمت چپ ویژگی قد می‌گذاریم.

انتخاب بر اساس محل زندگی

	$S = [2, -1]$ $E = 0.918$	
	محل زندگی	
خشکی		آب
$S = [0, -1]$ $E = 1$		$S = [+1, 0]$ $E = 0$
$IG(S, habitat) = 0.918 - 0.5 * 0 - 0.5 * 0 = 0.918$		

دو شاخه آخر تبدیل به برگ می‌شود. و ساختن درخت کامل می‌شود.

		تعداد پا				
	3		2			
B			رنگ			
		سبز		قهوه‌ای		
	قد				قد	
کوتاه		بلند		کوتاه		بلند
A		B		محل زندگی		A
			خشکی		آب	
		B				A

(ب)

Table 1: داده های آزمون و نتایج مدل

شماره	رنگ	تعداد پا	قد	محل زندگی	جاندار	خروجی مدل
1	قهوه‌ای	3	بلند	خشکی	B	B
2	سبز	2	بلند	خشکی	A	B
3	سبز	2	کوتاه	خشکی	A	A
4	قهوه‌ای	2	کوتاه	آب	B	A
5	قهوه‌ای	2	بلند	خشکی	A	A

دقت مدل 60% شد.

Table 2: ماتریس آشفتگی^۱ مدل

Label Prediction	A B	
	A	B
A	2	1
B	1	1

به دلیل کم بودن داده های آموزش، مدل عملکرد خوبی ندارد. و دو جاندار جای هم پیش بینی می شوند. مدل توانست دقت داده های آموزش را صفر کند ولی دقت داده های آزمون پایین است. یعنی واریانس آن بالا است و مدل overfit شده است.

ج)

بطور کلی الگوریتم ID3 بصورت حریصانه^۲ عمل کرده. و هر بار بهترین attribute را انتخاب می کند. تا به برگ های درخت برسد. روش Reduced-Error Pruning روشی برای هرس کردن درخت است. تا اندازه آن کوچک شود. این روش در ابتدا دقت مجموعه train را کاهش داده ولی دقت مجموعه test را افزایش می دهد. چون می خواهیم دقت مجموعه آموزش همچنان صفر بماند. حذف هر کدام از برگ ها باعث ضرر می شود. و نمی توان فقط با استفاده از 2 ویژگی درخت را ساخت. راه دیگر این است که $nCr(2^4)$ ویژگی را امتحان کرده و ببینیم که هیچ کدام دقت مجموعه آموزش را صفر نمی کند.

د)

درخت های تصمیم به دلیل آن که در مراحل ساخت تا جایی پیش می روند که یا ویژگی باقی نمانده باشد یا برگ ها همگی از یک نوع باشند. مقاومت به نویز کمی دارند. و وجود نویز باعث می شود که درخت بزرگتر شده تا اثر آن را خنثی کند. که این مسئله باعث بزرگ شدن درخت و overfitting می شود. دلیل دیگر نیز کمبود داده آزمون است که باعث خطای صفر آموزش و خطای بالای آزمون می شود.

برای حل این مشکل دو راهکار وجود دارد:

¹ Confusion matrix
² greedy

1. از ابتدا نگذاریم درخت به حد زیادی بزرگ شود. حدی برای حداکثر عمق درخت در نظر بگیریم. (max depth) یا اگر تعداد نمونه های باقی مانده در یک شاخه از حدی کمتر شد نیز درخت را متوقف کنیم. و دیگر تقسیم بندی انجام ندهیم (min sample split)
2. ابتدا بگذاریم درخت بطور کامل رشد کند. بعد از پایین شروع به "هرس" کردن درخت می کنیم. یک روش معمول Reduced-Error Pruning است. در این روش هر کدام از شاخه ها را به همراه زیرشاخه های آن حذف می کنیم. و جای آن پرتکرار ترین Label را می گذاریم. اگر حذف آن باعث بهبود دقت مجموعه آزمون شد آن را تبدیل به برگ می کنیم. در غیر این صورت آن را نگه می داریم این کار را تا جایی انجام می دهیم که هرگونه تغییر برای مدل مضر باشد.

سوال 2

الف)

ویژگی های موثر فرد سن، کلاس بلیط (که موقعیت فرد بر کشتی را مشخص می کند)، تعداد همراه (شامل فرزند، والدین، همسر) و جنسیت می باشد.

دیتاست شامل 891 مسافر کشتی تایتانیک است. ازین تعداد 700 مسافر را به عنوان داده آموزش و 191 نفر را داده آزمون انتخاب می کنیم. برای داده ها ستون Survived را جدا می کنیم. بعد برای داده آزمون ستون Cabin را که شماره کابین افراد است. و اطلاعات خاصی به ما نمیدهد و همچنین در داده آموزش وجود ندارد را حذف می کنیم. در dataset تعدادی missing value وجود دارد. برای ستون های Age, Fare از میانگین بقیه ستون ها برای پر کردن نمونه های خالی استفاده می کنیم. و این دو متغیر پیوسته را از مقدار اطراف میانگین خود سه تکه کرده و گسسته می کنیم (افراد جوان میانسال و پیر براس ستون Age و بلیط های ارزان و متوسط و گران برای ستون Fare). برای بقیه ستون ها به دلیل گسسته بودن مقادیر، از مد داده های ستون استفاده می کنیم.

برای درخت تصمیم گیری، از معیار انترپوی استفاده می کنیم.

در $depth = 3$ دقت درخت برابر 81.68% و ماتریس آشفستگی بصورت زیر است:

Table 3: ماتریس آشفستگی درخت در عمق 3

Label Prediction	Label	
	1	0
1	114	6
0	29	42

در $depth = 4$ دقت درخت برابر 83.77% و ماتریس آشفستگی بصورت زیر است:

Table 4: ماتریس آشفستگی درخت در عمق 4

Label Prediction	Label	
	1	0
1	110	10
0	21	50

در $\text{depth} = 5$ دقت درخت برابر 81.68% و ماتریس آشفته‌گی بصورت زیر است:

Table 5: ماتریس آشفته‌گی درخت در عمق 5

Label Prediction	1	0
1	108	12
0	23	48

در $\text{depth} = 6$ دقت درخت برابر 82.20% و ماتریس آشفته‌گی بصورت زیر است:

Table 6: ماتریس آشفته‌گی در عمق 6

Label Prediction	1	0
1	109	11
0	23	48

در $\text{depth} = 7$ دقت درخت برابر 91.62% و ماتریس آشفته‌گی بصورت زیر است:

Table 7: ماتریس آشفته‌گی در عمق 7

Label Prediction	1	0
1	112	8
0	8	63

با افزایش عمق طبقه بند از عمق 3 تا 6 دقت مدل تقریباً یکسان و حدود 82% است. به این معنا که ویژگی‌های اضافه شده به درخت تاثیر زیادی بر دقت آن ندارند. اما در عمق 7 که حداکثر عمق ممکن است. دقت به 92% می‌رسد.

ویژگی‌ها به ترتیب از پایین به بالا Sex, Pclass, Fare, SibSp, Parch, Age, Embarked می‌باشند.

(ب)

مهم‌ترین ایراد درخت تصمیم مستعد بودن به overfitting است. که باعث می‌شود. واریانس مدل زیاد شده. و نسبت به نویز حساس باشد. (برای داده‌های نویزی برای این که دقت مجموعه آموزش همچنان کم

باقی بماند، درخت شاخه های ضائد اضافه می کند که عمق آن را زیاده از حد می کند) ایراد دیگر آن سخت بودن اسفاده در متغیر های پیوسته است. که زیاد به گسسته سازی دارد.

حالا برای حل مشکل اول باید به گونه ای از overfitting جلوگیری کنیم. و واریانس را کاهش دهیم. به جای یک درخت کامل که نویز بالایی دارد. از چند درخت باهم استفاده می کنیم. و با رای اکثریت^۱ میانگین آن ها را می گیریم. از دو روش bootstrap aggregation و feature selection استفاده می کنیم. یعنی در دیتاست قسمتی از داده ها را تکثیر و جایگزین همان تعداد داده بصورت رندم می کنیم. این کار bootstrap aggregation نام دارد. سپس بر روی این داده ها درخت های خود را آموزش می دهیم. با این تفاوت که هر بار تعدادی از ویژگی ها را حذف می کنیم. (مثلا p ویژگی از کل m ویژگی) و سپس آموزش می دهیم. با این روش ها تعدادی درخت به وجود می آید. که پیش بینی نتایج را از طریق رای اکثریت بین جواب درخت ها انجام می دهیم. (در روش بالا اگر $m = p$ بود به آن روش bagging می گویند. در غیر اینصورت جنگل تصادفی است).

(ج)

برای $p = 6$ و تعداد درخت 5، دقت جنگل تصادفی برابر 84.82% و ماتریس آشفستگی آن بصورت زیر است:

Table 8 : ماتریس آشفستگی جنگل تصادفی

Label Prediction	Label	
	1	0
1	106	14
0	15	56

¹ Majority voting

سوال 3

طبقه بند k همسایه نزدیک

(الف)

با پیاده سازی الگوریتم KNN نتایج زیر برای دیتاست آزمون بدست آمد

$K = 1$:

Accuracy = 83.33%

Table 9: ماتریس آشفتگی برای $K = 1$

Prediction \ Label	0	1	2
0	10	0	0
1	1	12	1
2	0	4	8

$K = 5$

Accuracy = 77.78%

Table 10: ماتریس آشفتگی برای $K = 5$

Prediction \ Label	0	1	2
0	10	0	0
1	1	11	2
2	2	3	7

$K = 10$

Accuracy = 72.22%

Table 11 : ماتریس آشفته‌گی برای $K = 10$

Prediction \ Label	0	1	2
0	10	0	0
1	1	11	2
2	2	3	7

$K = 20$

Accuracy = 72.22%

Table 12 : ماتریس آشفته‌گی برای $K = 20$

Prediction \ Label	0	1	2
0	10	0	0
1	1	10	3
2	1	5	6

بهترین عملکرد مدل در $K = 1$ بدست آمد.

(ب)

با محاسبه احتمال کلاس پیش‌بینی شده و رسم هیستوگرام آن نمودارهای زیر برای K های مختلف بدست می‌آید:

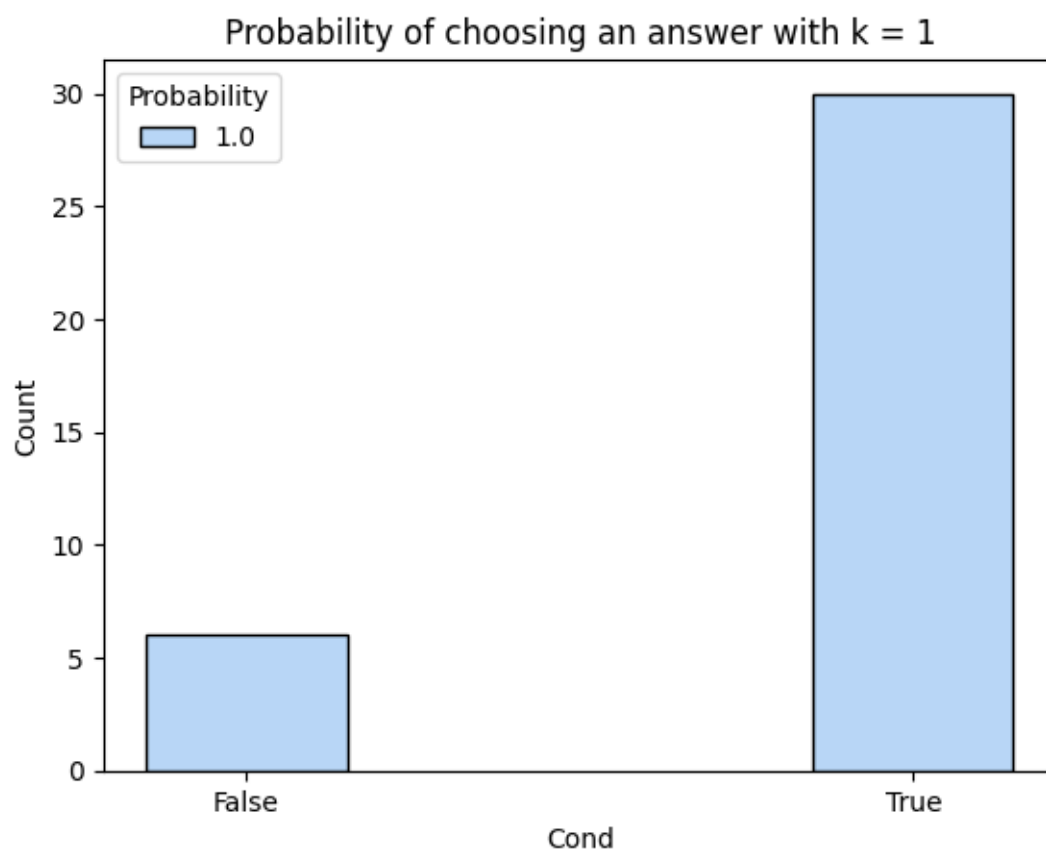


Figure 1: نمودار توزیع احتمال متعلق به هر کلاس برای $K = 1$

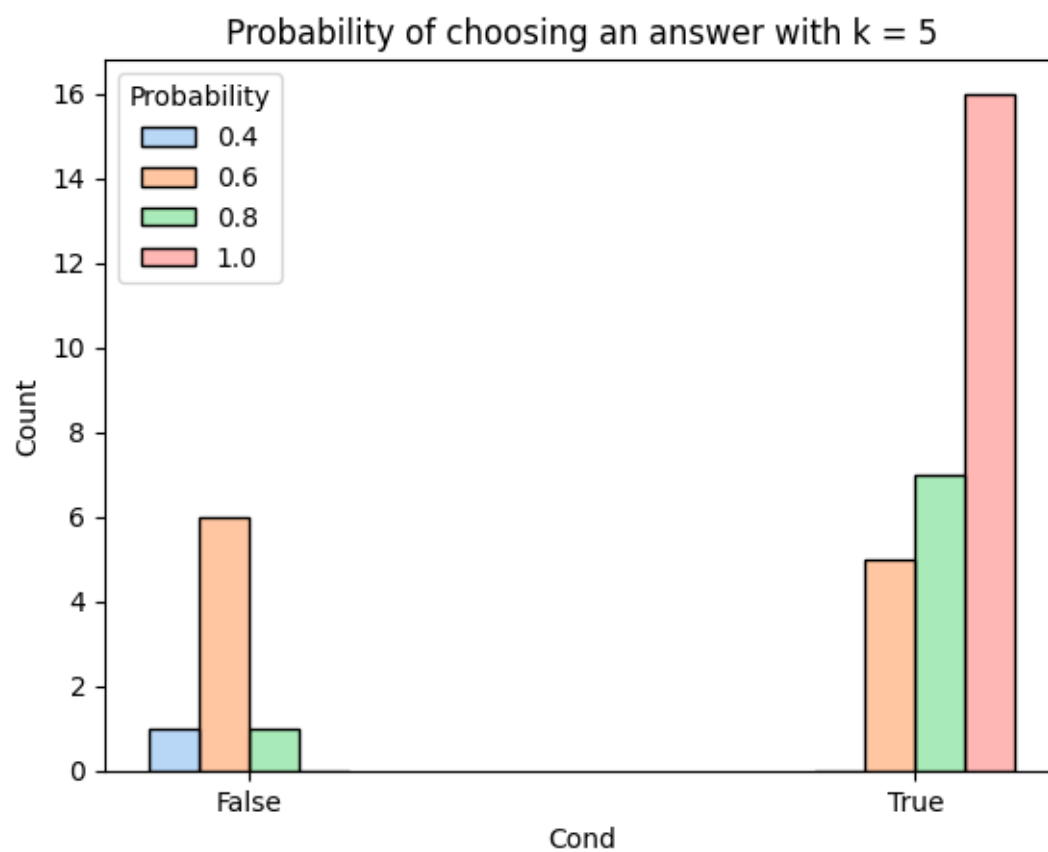


Figure 2: نمودار توزیع احتمال متعلق به هر کلاس برای $K = 5$

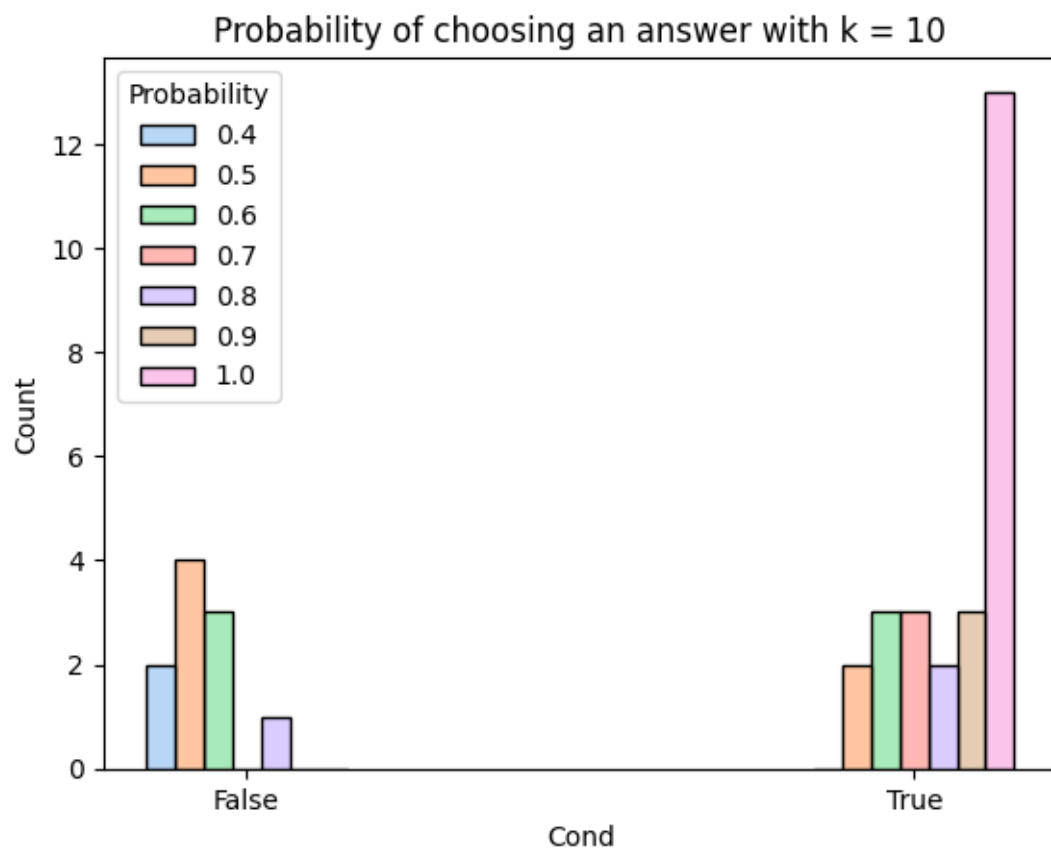


Figure 3: نمودار توزیع احتمال متعلق به هر کلاس برای $K = 10$

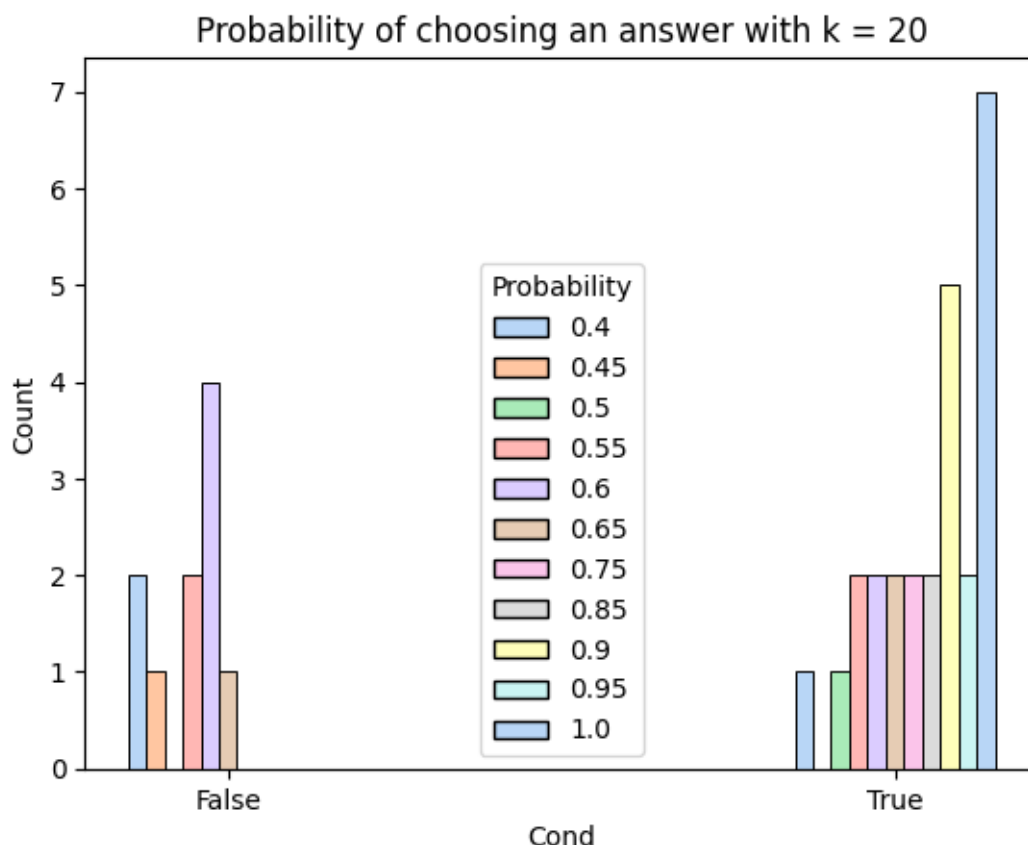


Figure 4: نمودار توزیع احتمال متعلق به هر کلاس برای $K = 20$

با اتکا به درصد درستی به تنهایی، در $K = 1$ مدل بهتر عمل کرده است. ولی با توجه به نمودارها می‌توان دید که در $K = 10$ احتمال درست بودن $p = 1$ به شکل چشم‌گیری از بقیه احتمال‌ها درست‌تر است. به این معنا که قطعیت آن بالا تر است. در حالی که در بقیه K ها احتمال $p = 1$ اختلاف چشم‌گیری ندارد.

یادگیری بر اساس معیار

(الف)

بطور کلی اهداف یادگیری بر اساس معیارهای LMNN, LFDA این است. که داده‌ها با یک تبدیل خطی به فضایی بروند که در آن فضا داده‌های مشابه به هم نزدیک‌تر و داده‌های غیر مشابه از هم دورتر شده باشند. در این معیارها ما از فاصله ماهالانوبیس^۱ به جای فاصله اقلیدسی استفاده می‌کنیم.

$$d(\vec{x}_i, \vec{x}_j) = (\vec{x}_i - \vec{x}_j)^T M (\vec{x}_i - \vec{x}_j)$$

¹ Mahalanobis distance

در معادله بالا M همان ماتریس تبدیل خطی ما است. که نیمه معین مثبت می باشد.

حالا ماتریس M را طوری تعریف می کنیم که فاصله بین داده های مشابه را مینیمم کند.

$$\min_M \sum_{i,j \in N_i} d(\vec{x}_i, \vec{x}_j)$$

و از طرف فاصله دیگر نمونه های غیر مشابه را که کمتر از یک واحد با نمونه های مشابه، با داده فاصله دارند را نیز حداکثر کند. یعنی یک margin بین داده های هم مشابه و همسایه با داده های غیر مشابه و همسایه ایجاد کند. (قسمت $[\cdot]_+ = \max(\cdot, 0)$ برای این است که داده های مشابه که همینطوری از داده های غیر مشابه به نمونه نزدیک تر هستند. دیگر در بهینه سازی نقش نداشته باشند).

$$\min_M \sum_{i,j \in N_i, l, y_l \neq y_i} [d(\vec{x}_i, \vec{x}_j) + 1 - d(\vec{x}_i, \vec{x}_l)]_+$$

معادله نهایی بصورت زیر است.

$$\min_M \sum_{i,j \in N_i} d(\vec{x}_i, \vec{x}_j) + \lambda \sum_{i,j,l} \xi_{ijl}$$

$$\forall i,j \in N_i, l, y_l \neq y_i$$

$$d(\vec{x}_i, \vec{x}_j) + 1 - d(\vec{x}_i, \vec{x}_l) \leq \xi_{ijl}$$

$$\xi_{ijl} \geq 0$$

M : semi definite positive

در واقع ξ_{ijl} همان نقش تابع max را بازی می کند. و برای این که زیر رادیکال منفی نشود نیز ماتریس M باید نیمه معین مثبت باشد.

در روش LFDA ابتدا تعریف می کنیم.

$$S^{(w)} = \frac{1}{2} \sum_{i,j} W_{i,j}^{(w)} (x_i - x_j)(x_i - x_j)^T$$

$$S^{(b)} = \frac{1}{2} \sum_{i,j} W_{i,j}^{(b)} (x_i - x_j)(x_i - x_j)^T$$

که $S^{(w)}$ ماتریس خوشه درون-کلاسی^۱ است. و $S^{(b)}$ ماتریس خوشه میان-کلاسی^۲ است. داریم:

$$A_{i,j} = \text{affinity of } (x_i, x_j)$$

¹ within-class scatter
² between-class scatter matrix

$$W_{i,j}^{(w)} = \begin{cases} \frac{A_{i,j}}{n_l} & \text{if } y_i = y_j = l \\ 0 & y_i \neq y_j \end{cases}$$

$$W_{i,j}^{(b)} = \begin{cases} A_{i,j} \left(\frac{1}{n} - \frac{1}{n_l} \right) & \text{if } y_i = y_j = l \\ \frac{1}{n} & y_i \neq y_j \end{cases}$$

$$T_{LFDA} = \operatorname{argmax}_{T \in R} [\operatorname{tr}(T^T S^{(w)} T)^{-1} T^T S^{(b)} T]$$

(ب)

1. در طبقه بند k همسایه، پارامتر k تعداد همسایه نزدیکی را نشان می‌دهد. که با استفاده از آن ها لیبل نمونه ها مشخص می‌شود. ولی در روش های یادگیری بر اساس معیار پارامتر K در واقع تعداد همسایه با لیبل مشابه است که قرار است نزدیک نگه داشته شوند. و بقیه همسایه ها با ایجاد یک margin دور شوند.
2. روش اول را انتخاب نموده و به کمک کتابخانه metric-learning بعد داده ها را کاهش دادیم.
3. به کمک دستور های fit, transform ابعاد داده‌ها کاهش یافت و به فضای 2 بعدی رفتند. علاوه بر آن تبدیل ماتریسی تغییر پایه نیز بر آن ها اعمال شد. برای K های مختلف نمودار های زیر بدست آمد:

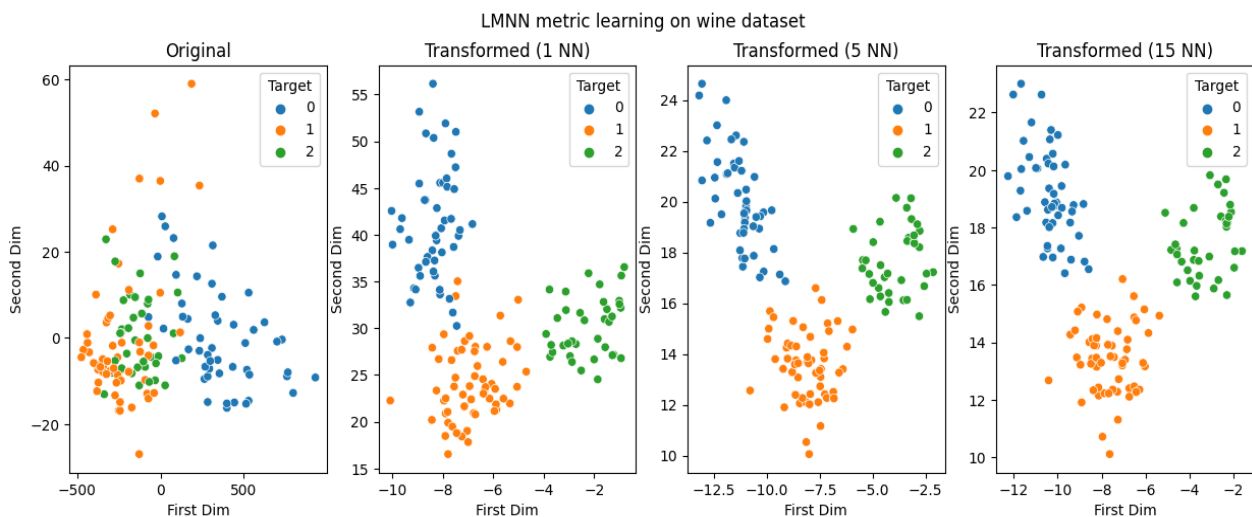


Figure 5 افزایش دادگان اصلی و انتقال یافته به ازی مقادیر مختلف K با روش LMNN

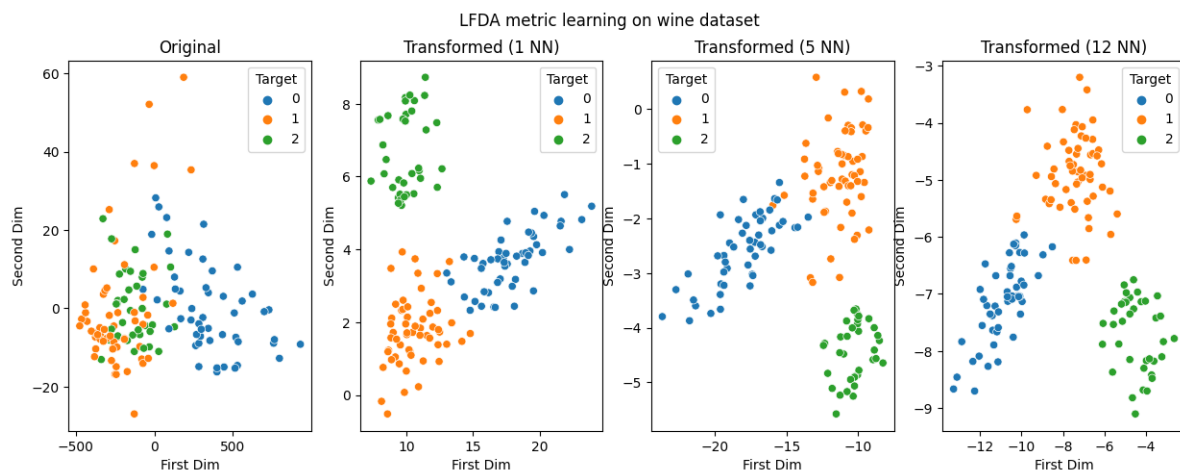


Figure 6 افزایش دادگان اصلی و انتقال یافته به ازای مقادیر مختلف K با روش LFDA

4. هر چقدر مقدار K بیشتر شود تفکیک پذیری بیشتری دارد. دلیل آن این است. که مقدار k همانطور که در بخش الف گفته شد. K عدد از همسایه های مشابه نزدیک را در نظر می گیرد. و ماتریسی میسازد که داده های مشابه همسایه را نزدیک تر و داده های غیر مشابه را با یک margin دور می سازد. طبیعتاً هر چه مقدار k بیشتر باشد. همسایه های بیشتری به هم نزدیک می شوند و غیر همسان ها دور تر و قابلیت تفکیک پذیری بالاتر می رود.

(ج)

در بخش قبل مقدار $K = 1$ بهترین عملکرد را داشت. با دادگان انتقال یافته در فضای جدید برای هر دو روش LFDA, LMNN داریم:

$$\text{Accuracy} = 94.44 \%$$

Table 13 ماتریس آشفتگی مدل در فضای انتقال یافته

Prediction \ Label	0	1	2
0	12	0	0
1	1	16	1
2	0	0	6

طبقه بند در فضای انتقال یافته جدید بسیار بهتر عمل کرده که دلیل آن تفکیک پذیری بهتر با استفاده از یادگیری معیار است.

(د)

با رسم ماتریس همبستگی دیتاست برای حالت بدون انتقال می‌بینیم که تعدادی از داده ها با هم همبستگی مثبت و تعدادی همبستگی مثبت دارند. با اعمال روش LMNN می‌بینیم که همبستگی داده های که مثبت بوده بیشتر شده و همبستگی داده هایی که منفی بوده اند بیشتر شده است. اما در ماتریس همبستگی روش LFDA می‌بینیم که تقریباً همبستگی بین داده ها صفر شده است. یعنی ستون های ویژگی از هم مستقل شده اند. که باعث می‌شود. تفکیک پذیری داده ها بیشتر شود.

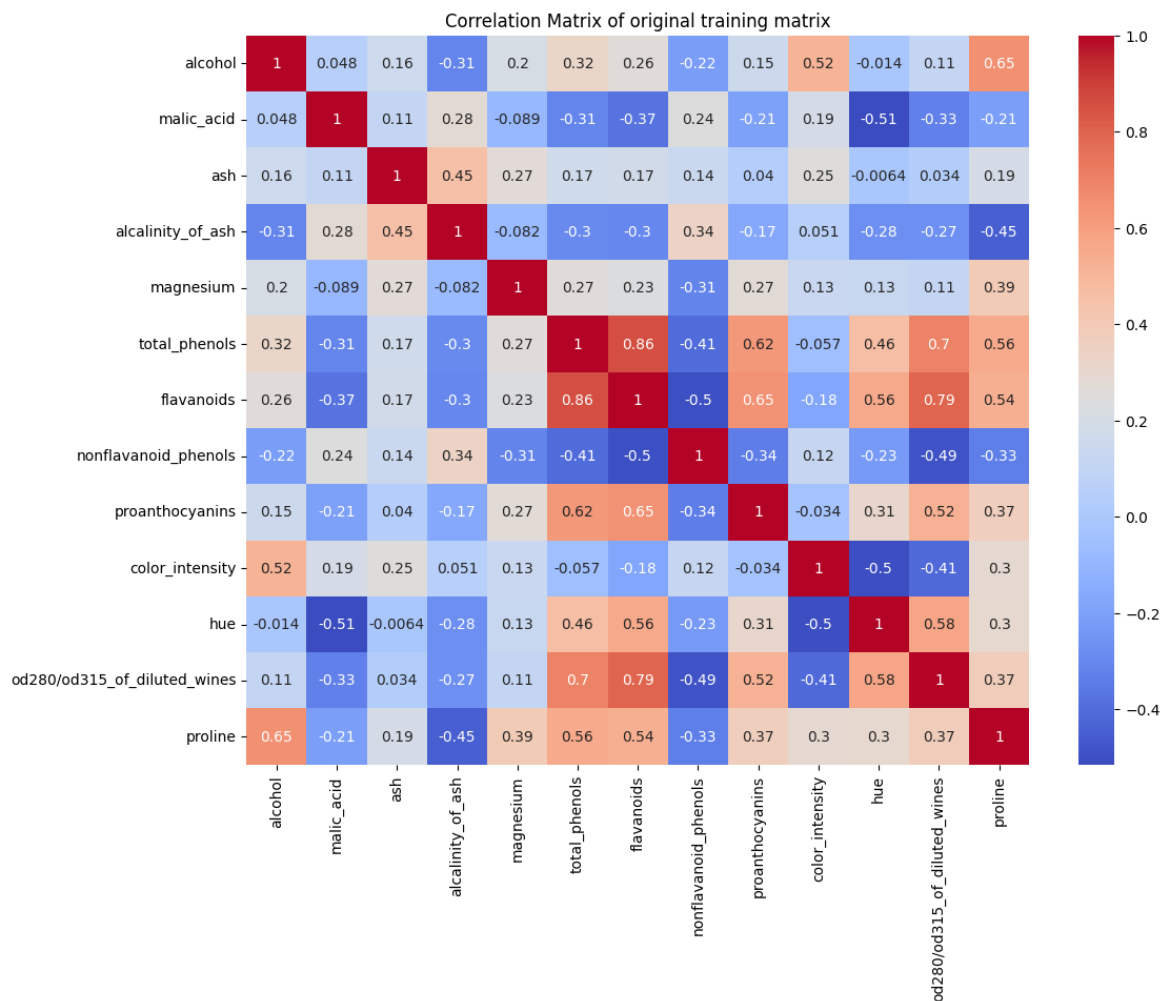


Figure 7 ماتریس همبستگی دیتاست اصلی

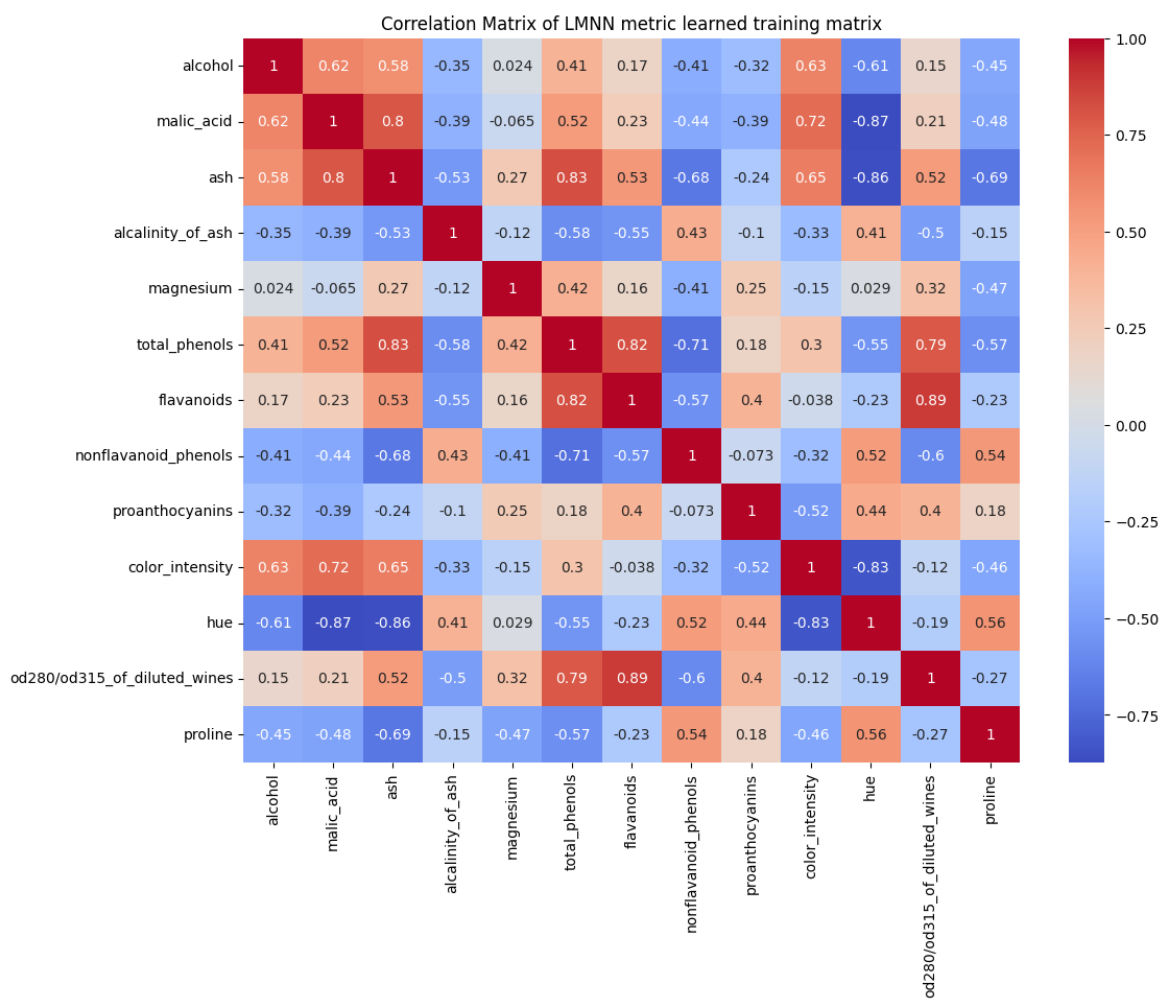


Figure 8 ماتریس همبستگی دیتاست انتقال یافته به روش LMNN

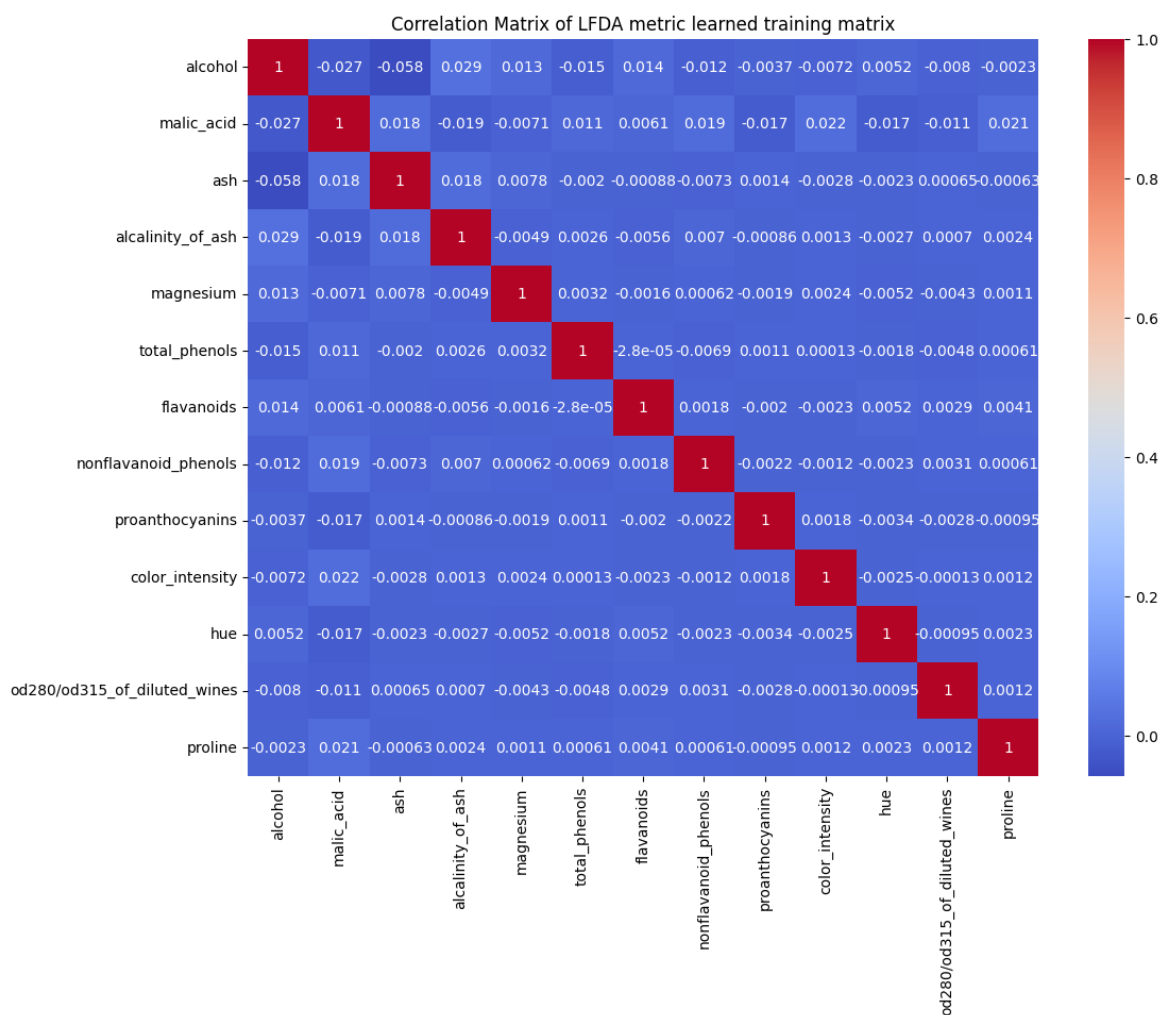


Figure 9 ماتریس همبستگی دیتاست انتقال یافته به روش LFDA

(۵)

در روش GMLL پایه روش همان دور کردن نقاط غیر مشابه و نزدیک کردن نقاط مشابه به هم است. همانند روش هایی مثل MMC این روش نیز یک بخش دارد که فاصله داده های نزدیک را کم می کند.

$$\min_M \sum_{(x_i, x_j) \in S} d_A(\vec{x}_i, \vec{x}_j)$$

و بخش دیگر که فاصله داده های غیر مشابه را کم می کند. با این تفاوت که از خاصیت ماتریس های مثبت معین استفاده می کند که $A > b \implies B^{-1} > A^{-1}$ یعنی همزمان هر چقدر که ماتریس A فاصله Mahalanobis را کم می کند. ماتریس A^{-1} فاصله Mahalanobis را زیاد می کند. که این دقیقاً خواسته ما از مسئله است.

$$\begin{aligned} \min_M \sum_{(x_i, x_j) \in S} &= d_A(\vec{x}_i, \vec{x}_j) + \sum_{(x_i, x_j) \in D} d_{A^{-1}}(\vec{x}_i, \vec{x}_j) \\ &= \sum_{(x_i, x_j) \in S} \text{tr}(A(x_i - x_j)(x_i - x_j)^T) + \sum_{(x_i, x_j) \in D} \text{tr}(A^{-1}(x_i - x_j)(x_i - x_j)^T) \end{aligned}$$

فرق روش LMNN و GMML این است که روش LMNN به شکل غیر متقارن با داده ها رفتار می کند. و هر کدام را جداگانه از تابع max عبور می دهد. ولی روش GMML فاصله Mahalanobis را بصورت یک حاصل جمع بر روی کل تابع هزینه اعمال می کند.

پیوست

سوال (2) توضیح کد

در قسمت اول کد درخت تصمیم ابتدا داده ها را وارد کرده. و بر اساس توضیحات سوال 2 پیش پردازش را انجام می دهیم. توابع entropy, information_gain, best_attribute را تعریف می کنیم. این توابع برای انتخاب بهترین ویژگی برای تقسیم درخت هستند. تابع best_attribute برای هر ستون دیتاست ورودی بهره اطلاعات را به کمک معیار آنتروپی حساب می کند. و بهترین ویژگی را برمی گرداند. جهت تعریف بهتر درخت، یک کلاس Node تعریف می کنیم. که شامل attribute های پدر، فرزند و زیر درخت و ... است.

حالا به سراغ تعریف تابع درخت تصمیم گیری می رویم. این درخت بصورت بازگشتی پیاده سازی شده. در ابتدا یک نود ریشه می سازیم این نود دارای بهترین ویژگی گرفته شده از تابع best_attribute است. برای هر مقدار از این ویژگی به کمک for یک نود فرزند ساخته و قسمتی از دیتاست که از پدر آمده و دارای این ویژگی است به عنوان subtree نود فرزند قرار داده می شود. حالا برای قسمت بازگشتی به هر نود فرزند همانند یک ریشه نگاه می کنیم. و دوباره تابع decision_tree را روی آن فرا می خوانیم. این بازگشت تا جایی ادامه پیدا می کند. تا عمق درخت به حداکثر عمق تعیین شده برسد. در این صورت Node های آخر تبدیل به برگ می شوند. و مقدار آن ها برابر بیشترین مقدار تکرار شده ستون هدف می شود.

بعد از ساخت درخت به سراغ تست آن می رویم. تابع predict مسافر و درخت ساخته شده را میگیرد و نتیجه را برمی گرداند. طرز کار آن بازگشتی است. به این صورت که از ریشه شروع کرده و هر بار مقدار value فرزند های آن را نگاه می کند. هر کدام که مقدار value در آن ویژگی ستون خودش برابر بود. از همان نود فرزند دوباره تابع را فرا می خواند تا جایی که به برگ درخت برسد و مقدار پیشبینی را برگرداند.

تابع score نیز برای همه نمونه های دیتاست آزمون نتیجه را برمی گرداند. در یک لیست ریخته؛ تبدیل به دیتافریم می کند. و دقت تابع را همراه با پیش بینی ها برمی گرداند. در تابع confusion_matrix() نیز

پیش بینی ها را گرفته و به کمک دستور crosstab در کتابخانه pandas به راحتی ماتریس آشفتگی آن را برمی گرداند.

برای قسمت دوم سوال که پیاده سازی جنگل تصادفی است. آرگومان های حداکثر عمق، تعداد ویژگی ها و تعداد درختان را می گیرد. در یک حلقه for به تعداد درختان ورودی ابتدا به شکل رندوم تعدادی ویژگی را انتخاب می کنیم. و بقیه ویژگی ها را از دیتاست دور انداخته. سپس bootstrap aggregation انجام می دهیم. و دیتاست آماده شده را به تابع decision_tree() می دهیم. تابع همه درخت ها را در یک لیست به اسم forest ریخته و برمی گرداند.

حالا برای درخت نیز باید تابع forest_predict() بسازیم. میتوان به راحتی از تابع predict() درخت تصادفی استفاده کرد. و پیش بینی ها را در یک لیست ریخت و مد آن را برگرداند. در آخر نیز از مدل خود run می گیریم.

سوال 3) توضیح کد

ابتدا کتابخانه های مورد نظر را import می کنیم. داده ها به ترتیب هستند و نیاز به بر خوردن دارند. برای این کار ابتدا دیتاست را به دو قسمت X, y تقسیم می کنیم. سپس X, y را زیپ کرده و در یک لیست آن را shuffle می کنیم. و دوباره unzip می کنیم. دیتاست را به نسبت 1 به 4 تقسیم می کنیم.

تابع KNN را تعریف می کنیم. برای هر نقطه در دیتاست آزمون، فاصله از هر نقطه در دیتاست آموزش را به کمک تابع np.linalg.norm() حساب کرده و در یک لیست می ریزیم. لیست را مرتب کرده و از بین K داده اول، لیبل مد آن را به همراه لیبل همه k نقطه برمی گردانیم.

برای مشخص شدن دقت مدل تابع score() را تعریف می کنیم. نتایج و لیبل دیتاست تست را می گیرد و تعداد درست ها بر تعداد کل نمونه هارا برمی گرداند. تابع confusion_matrix() نیز لیبل پیش بینی ها را گرفته و به کمک دستور crosstab در ماتریس آشفتگی آن را برمی گرداند.

حالا برای قسمت ب سوال 3 تابع result_prob() را تعریف می کنیم. این تابع برای هر پیش بینی، تعداد لیبل های پیش بینی را به تعداد k تقسیم کرده و احتمال لیبل خود را برمی گرداند. حالا این تابع را برای K های مختلف فراخوانده و نمودار هیستوگرام بر اساس درست بودن یا غلط بودن پیش بینی ها رسم می کنیم.

برای قسمت دوم که یادگیری بر اساس معیار است. ابتدا از پکیج LMNN داده ها را fit و سپس transform می کنیم. این کار را برای K (پارامتر روش LMNN) های مختلف انجام داده و نمودار آن را رسم

می‌کنیم. سپس برای $k = 15$ داده‌های انتقال یافته را به مدل KNN داده و دقت و ماتریس آشفتگی آن را چاپ می‌کنیم. مشابه این اعمال را برای LFDA انجام می‌دهیم.

منابع:

[Large Margin Nearest Neighbors \(LMNN\) Algorithm \(open-instruction.com\)](https://open-instruction.com/algorithm/lmnn/)

[Deep Large Margin Nearest Neighbor for Gait Recognition \(degruyter.com\)](https://degruyter.com/article/view/111111)