



دانشگاه تهران  
پردیس دانشکده‌های فنی  
دانشکده برق و کامپیوتر



## یادگیری ماشین

### تمرین شماره ۱

نام و نام خانوادگی  
امیرحسین چمی

آبان ۱۴۰۴

## فهرست

۴	..... فهرست شکل ها
۵	..... فهرست جدول ها
۶	..... چکیده
۷	..... پرسش ۱ – MINIMUM RISK CLASSIFICATION WITH REJECTION OPTION
۷	..... ۱-۱. اثبات نامساوی
۹	..... پرسش ۲ – NAIVE BAYES
۹	..... ۱-۲. بخش A: $PY X1, X2(y x1 = 0, x2 = 0)$
۱۰	..... ۲-۲. بخش B: $PY X1yx1 = 0$
۱۱	..... ۳-۲. بخش C: $PY X2yx2 = 1$
۱۲	..... ۴-۲. بخش D: تحلیل بخش های قبلی
۱۲	..... تحلیل بخش $a$
۱۲	..... تحلیل بخش $b$
۱۳	..... تحلیل بخش $c$
۱۴	..... پرسش ۳ – MULTINOMIAL NAIVE BAYES
۱۴	..... ۱-۳. بخش A: $P(y = 1 x)$
۱۵	..... ۲-۳. بخش B: LINEARITY OF DECISION BOUNDARY
۱۶	..... ۳-۳. بخش C: LOGISTIC FUNCTION
۱۷	..... پرسش ۴ – MLE
۱۷	..... ۱-۴. بخش A:
۱۷	..... ۲-۴. بخش B:
۲۰	..... پرسش ۵ – MAP
۲۰	..... ۱-۵. بخش A: CONJUGACY OF THE DIRICHLET PRIOR
۲۱	..... ۲-۵. بخش B: POSTERIOR DISTRIBUTION
۲۲	..... پرسش ۶ – GAUSSIAN DISCRIMINANT ANALYSIS

۲۲	۱-۶. بخش A: محاسبه ی میانگین و کوواریانس
۲۳	۲-۶. بخش B: DISCRIMINANT FUNCTION
۲۳	۳-۶. بخش C: DECISION REGIONS
۲۴	۴-۶. بخش D: خط مرز تصمیم گیری
۲۵	۵-۶. بخش E: تاثیر تابع پیشین در نواحی تصمیم گیری
۲۵	۶-۶. بخش F: تکرار سوال برای سه کلاس
۲۵	۱-۶-۶. بخش a: محاسبه ی میانگین و کوواریانس
۲۶	۲-۶-۶. بخش b: Discriminant function
۲۶	۳-۶-۶. بخش c: Decision Regions
۲۷	۷-۶. بخش I: تکرار سوال با داده های کمتر
۲۷	۱-۷-۶. بخش a: محاسبه ی میانگین و کوواریانس
۲۸	۲-۷-۶. بخش b: Discriminant function
۲۸	۳-۷-۶. بخش c: Decision Regions
۳۳	پرسش ۷ - GAUSSIAN NAIVE BAYES
۳۳	۱-۷. بخش D: عملکرد مدل
۳۴	پرسش ۸ - BAG OF WORDS — NAIVE BAYES CLASSIFIER
۳۴	۱-۸. بخش A: NAIVE BAYES CLASSIFIER — MLE
۳۵	۲-۸. بخش B: NAIVE BAYES CLASSIFIER — MAP WITH DIRICHLET PRIOR
۳۵	۳-۸. بخش C: EFFECT OF B ON CLASSIFICATION ACCURACY
۳۶	۴-۸. بخش D: EFFECT OF TRAINING SET SIZE
۳۸	مراجع
۳۸	REFERENCES

- شکل ۱. نمودار دیتاست سوال ۶ ..... ۲۲
- شکل ۲. ناحیه های تصمیم گیری دو کلاس ..... ۲۴
- شکل ۳. تاثیر احتمال پیشین بر مرز تصمیم‌گیری ..... ۲۵
- شکل ۴. ناحیه های تصمیم‌گیری سه کلاس ..... ۲۷
- شکل ۵. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -10$  ..... ۲۹
- شکل ۶. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -5$  ..... ۲۹
- شکل ۷. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -4$  ..... ۳۰
- شکل ۸. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -3$  ..... ۳۰
- شکل ۹. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -2$  ..... ۳۱
- شکل ۱۰. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -1$  ..... ۳۱
- شکل ۱۱. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = 0$  ..... ۳۲
- شکل ۱۲. تغییرات دقت نسبت به بتا ..... ۳۶
- شکل ۱۳. تاثیر اندازه‌ی داده‌ی آموزش بر دقت مدل ..... ۳۸

## فهرست جدول‌ها

- جدول ۱. تاثیر بتا بر دقت مدل ..... ۳۶
- جدول ۲. تاثیر اندازه‌ی داده‌ی آموزش بر دقت مدل ..... ۳۷

در این تمرین روش های مختلف طبقه بندی را مورد بررسی قرار می دهیم. سوالات ۱ الی ۵ بر مفاهیم تئوری تمرکز دارند. سوال ۱ مفهوم تصمیم گیری با حداقل ریسک و استفاده گزینهی reject را بررسی می کند و می خواهد دو نامساوی مربوط به قانون تصمیم گیری را اثبات کنیم. سوال ۲ تحلیل تئوری از naive bayes است و باید مستقیماً محاسبات انجام دهیم. در واقع در این تمارین naive bayes داریم اثر مستقل فرض کردن ویژگی ها را می بینیم. سوال ۳ دربارهی naive bayes چند متغیره است و اثبات می کند که مرز تصمیم گیری در آن خطی است. سوال ۴ دربارهی MLE در حالت چندمتغیره است. این سوال صرفاً دنبال کردن روش حل مسائل MLE است. سوال ۵ تخمین MAP با توزیع پیشین Dirichlet را مطرح می کند و تاثیر پارامتر های آلفا را روی احتمال بررسی می کند.

سوالات ۶ الی ۸ وارد فاز عملی می شوند. سوال ۶ یک دیتاست ارائه می کند و ابتدا با فرض مشترک بودن کوواریانس و سپس کوواریانس متفاوت برای هر کلاس، مرز تصمیم را بررسی می کند. همچنین در این سوال اثر تغییر احتمال پیشین را بررسی و گزینهی reject با هزینهی لاندا را تحلیل می کنیم. سوال ۷ شامل پیاده سازی کامل gaussian naive bayes است که تخمین پارامتر ها، طبقه بندی داده ها و ارزیابی عملکرد مدل را شامل می شود. سوال ۸ برای حل یک مسئله هم از MLE و هم از MAP استفاده می کند و می بینیم که MLE اصلاً عملکرد خوبی ندارد و علت این عملکرد ضعیف را بررسی می کنیم. به نظر من تمرین یک در مجموع تعادل خوبی از مسائل تئوری و پیاده سازی کد دارد.

## پرسشی ۱ – Minimum Risk Classification with Rejection Option

در این سوال باید یک نامساوی را در مسئله ی طبقه بندی حداقل ریسک ( طبقه بندی بیزی) اثبات کنیم.

### ۱-۱. اثبات نامساوی

مطابق صورت سوال cost به صورت زیر تعریف شده است.

$$h(\omega_i|\omega_j) = \begin{cases} \cdot & \text{if } \omega_i = \omega_j \quad (i.e., \text{CorrectClassification}) \\ \lambda_r & \text{if } \omega_i = \omega. \quad (i.e., \text{Rejection}) \\ \lambda_s & \text{Otherwise} \quad (i.e., \text{SubstitutionError}) \end{cases} \quad (1)$$

مطابق آنچه در کلاس درس دیدیم، Expected loss به صورت زیر تعریف می شود.

$$R(\alpha_i|x) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|x) \quad (2)$$

در رابطه ی (۲)،  $\alpha_i$  اقدام اتخاذ شده است. برای مسئله ی ۱ داریم:  $\alpha_i = \omega_i$   
در این مسئله ی مینیمم ریسک، برای انتخاب کلاس  $\omega_i$  باید داشته باشیم:

$$R(\alpha_i|x) \leq R(\alpha_j|x), \forall j \quad (3)$$

با جایگذاری مقادیر داریم.

$$\lambda_s(1 - P(\omega_i|x)) \leq \lambda_s(1 - P(\omega_j|x)), \forall j \quad (4)$$

با تقسیم طرفین نامساوی بر  $\lambda_s > 0$  خواهیم داشت.

$$1 - P(\omega_i|x) \leq 1 - P(\omega_j|x), \forall j \quad (5)$$

در نهایت

$$P(\omega_i|x) \geq P(\omega_j|x), \forall j \quad (6)$$

نامساوی اول اثبات شد.

در این مسئله ی مینیمم ریسک، ریسک رد کردن باید مساوی یا بیشتر از ریسک بقیه ی حالات باشد در غیر این صورت تصمیم گیری ما بهینه نیست. در واقع داریم:

$$R(\omega. |x) \geq R(\omega_i|x) \quad (7)$$

حال هر بخش را جداگانه محاسبه می کنیم. مقدار جمع ریسک برای طبقه بندی درست به علاوه ی *substitution error* به صورت زیر است.

$$R(\omega_i|x) = \sum_{j \neq i} \lambda(\omega_i|\omega_j) P(\omega_j|x) = \lambda_r P(\omega_i|x) + \lambda_s \sum_{j \neq i} P(\omega_j|x) \quad (8)$$

مقدار ریسک برای رد کردن به صورت زیر است.

$$R(\omega_i|x) = \lambda_r \quad (9)$$

از طرفی می دانیم جمع تمامی احتمالات مساوی ۱ می شود. پس:

$$\sum_{j=1}^C P(\omega_j|x) = 1 \quad (10)$$

$$1 - P(\omega_i|x) = \sum_{j \neq i} P(\omega_j|x) \quad (11)$$

می توانیم رابطه ی (۱۱) را در (۸) اعمال کنیم.

$$R(\omega_i|x) = \lambda_s P(\omega_j|x) = \lambda_s [1 - P(\omega_i|x)] \quad (12)$$

حال از (۹) و (۱۲) در (۷) استفاده می کنیم.

$$\lambda_r \geq \lambda_s [1 - P(\omega_i|x)] \quad (13)$$

$$P(\omega_i|x) \geq 1 - \frac{\lambda_r}{\lambda_s} \quad (14)$$

نامساوی دوم اثبات شد.



## پیش ۲ – Naive Bayes

در این سوال محاسبات یک طبقه بندی Naive Bayes را انجام می دهیم. سپس با توجه به محاسبات، رابطه‌ی بین کلاس ها و ویژگی ها را تحلیل و بررسی می کنیم. ابتدا این مبحث را از جزوه ی درس و [۱] مرور کردم.

### ۱-۲. بخش a: $P_{Y|X_1, X_2}(y|x_1 = 0, x_2 = 0)$

ابتدا مقادیر احتمال هر حالت را با جایگذاری مقادیر در توزیع ها محاسبه می کنم.

ویژگی اول:

$$P_{X_1|Y}(x_1 = 0 | Y = c) = 1 - \theta_c = \begin{cases} 0.5 & \text{if } c = 0 \\ 0.25 & \text{if } c = 1 \\ 0.5 & \text{if } c = 2 \end{cases} \quad (15)$$

ویژگی دوم:

$$P_{X_2|Y}(x_2 = 0 | Y = c) = \frac{1}{\sqrt{2\pi} \sigma_c} e^{-\frac{(0 - \mu_c)^2}{2\sigma_c^2}} = \begin{cases} 0.2419707245 & \text{if } c = 0 \\ 0.3989422804 & \text{if } c = 1 \\ 0.2419707245 & \text{if } c = 2 \end{cases} \quad (16)$$

همچنین با استفاده از قانون احتمال کل می توانیم احتمال مشترک ویژگی ها را محاسبه کرد.

$$P(X_1, X_2) = \sum_{c=0}^2 P(x_1, x_2 | Y = c) P(Y = c) \quad (17)$$

با توجه به naive بودن مسئله، قانون احتمال کل به صورت زیر خواهد بود.

$$P(X_1, X_2) = \sum_{c=0}^2 P(x_1 | Y = c) P(x_2 | Y = c) P(Y = c) \quad (18)$$

با جایگذاری مقادیر در رابطه ی (۱۸) به ازای هر کلاس داریم:

$$\begin{cases} 0.5 * 0.2419707245 * 0.5 & \text{if } c = 0 \\ 0.25 * 0.3989422804 * 0.25 & \text{if } c = 1 \\ 0.5 * 0.2419707245 * 0.25 & \text{if } c = 2 \end{cases} = \begin{cases} 0.06049268113 & \text{if } c = 0 \\ 0.02493389253 & \text{if } c = 1 \\ 0.03024634056 & \text{if } c = 2 \end{cases} \quad (19)$$

در مجموع احتمال مشترک ویژگی ها از جمع مقادیر رابطه‌ی (۱۹) محاسبه می شود.

$$P(X_1 = 0, X_2 = 0) = 0.1156729142 \quad (20)$$

در جزوه‌ی درس و منبع [۱] دیدیم که رابطه‌ی posterior بر اساس قانون بیز به صورت زیر است. مساوی دوم

در رابطه ی (۲۰) بر اساس naive بودن نوشته شده.

$$P(Y|X_1, X_2) = \frac{P(X_1, X_2|Y) \cdot P(Y)}{P(X_1, X_2)} = \frac{P(X_1|Y) \cdot P(X_2|Y) \cdot P(Y)}{P(X_1, X_2)} \quad (21)$$

حال با جایگذاری داریم.

$$\begin{aligned} P(Y|x_1 = \cdot, x_2 = \cdot) &= \frac{P(x_1 = \cdot | Y) \cdot P(x_2 = \cdot | Y) \cdot P(Y)}{P(x_1 = \cdot, x_2 = \cdot)} \\ &= \begin{cases} \frac{0.5 * 0.2419707245 * 0.5}{0.1156729142} & \text{if } c = 0 \\ \frac{0.25 * 0.3989422804 * 0.25}{0.1156729142} & \text{if } c = 1 \\ \frac{0.5 * 0.2419707245 * 0.25}{0.1156729142} & \text{if } c = 2 \end{cases} \\ &= \begin{cases} 0.5229632325 & \text{if } c = 0 \\ 0.2155551513 & \text{if } c = 1 \\ 0.2614816163 & \text{if } c = 2 \end{cases} \end{aligned} \quad (22)$$

که اگر بخواهیم در فضای  $\mathbb{R}^3$  نمایش بدهیم، داریم.

$$P(Y|x_1 = \cdot, x_2 = \cdot) = [0.5229632325, 0.2155551513, 0.2614816163] \quad (23)$$

جمع این مقادیر برابر ۱ می شود. (برای اینکه این جمع برابر یک باشد، مقادیر را تا ده الی یازده رقم اعشار محاسبه کردم).

## ۲-۲. بخش b: $P_{Y|X_1}(y|x_1 = 0)$

با استفاده از قانون بیز داریم:

$$P_{Y|X_1}(Y|x_1 = \cdot) = \frac{P(x_1 = \cdot | Y)P(Y)}{P(x_1 = \cdot)} \quad (24)$$

عبارت  $P(x_1 = \cdot | Y)$  را در رابطه‌ی (۱۵) محاسبه کردم. عبارت  $P(Y)$  را هم در صورت مسئله داریم. اما مخرج کسر با استفاده از قانون احتمال کل محاسبه می شود.

$$\begin{aligned} P(x_1 = \cdot) &= \sum_{c=\cdot}^2 P(x_1 = \cdot | Y = c)P(Y = c) \\ &= 0.5 * 0.5 + 0.25 * 0.25 + 0.5 * 0.25 = 0.4375 \end{aligned} \quad (25)$$

حال با جایگذاری در (۲۴) خواهیم داشت:

$$P_{Y|X_1}(Y|x_1 = \cdot) = \frac{P(x_1 = \cdot | Y)P(Y)}{P(x_1 = \cdot)} = \begin{cases} \frac{0.5 * 0.5}{0.4375} & \text{if } c = 0 \\ \frac{0.25 * 0.25}{0.4375} & \text{if } c = 1 \\ \frac{0.5 * 0.25}{0.4375} & \text{if } c = 2 \end{cases} \quad (26)$$

$$= \begin{cases} 0.5714285714 & \text{if } c = 0 \\ 0.1428571429 & \text{if } c = 1 \\ 0.2857142857 & \text{if } c = 2 \end{cases}$$

اگر بخواهیم این مقادیر را در فضای  $\mathbb{R}^3$  بنویسیم:

$$P_{Y|X_1}(Y|x_1 = \cdot) = [0.5714285714, 0.1428571429, 0.2857142857] \quad (27)$$

جمع این مقادیر برابر ۱ می شود. (برای اینکه این جمع برابر یک باشد، مقادیر را تا ده رقم اعشار محاسبه کردم.)

## ۳-۲. بخش c: $P_{Y|X_2}(y|x_2 = 1)$

با استفاده از قانون بیز داریم:

$$P_{Y|X_2}(Y|x_2 = 1) = \frac{P(x_2 = 1 | Y)P(Y)}{P(x_2 = 1)} \quad (28)$$

عبارت  $P(Y)$  در صورت مسئله داده شده.  $P(x_2 = 1 | Y)$  را به صورت زیر محاسبه می کنیم.

$$P(x_2 = 1 | Y) = \frac{1}{\sqrt{2\pi} \sigma_c} e^{-\frac{(1-\mu_c)^2}{2\sigma_c^2}} = \begin{cases} 0.05399096651 & \text{if } c = 0 \\ 0.2419707245 & \text{if } c = 1 \\ 0.3989422804 & \text{if } c = 2 \end{cases} \quad (29)$$

مخرج کسر با استفاده از قانون احتمال کل محاسبه می شود

$$P(x_2 = 1) = \sum_{c=0}^2 P(x_2 = 1 | Y = c)P(Y = c) \quad (30)$$

$$= 0.5 * 0.05399096651 + 0.25 * 0.2419707245 + 0.25$$

$$* 0.3989422804 = 0.1872237345$$

حال مقادیر (۳۰) و (۲۹) را در (۲۸) قرار می دهیم.

$$P_{Y|X_r}(Y|x_r = 1) = \frac{P(x_r = 1 | Y)P(Y)}{P(x_r = 1)} = \begin{cases} \frac{0.05399096651 * 0.05}{0.1872237345} & \text{if } c = 0 \\ \frac{0.2419707245 * 0.25}{0.1872237345} & \text{if } c = 1 \\ \frac{0.39896422804 * 0.25}{0.1872237345} & \text{if } c = 2 \end{cases} \quad (31)$$

$$= \begin{cases} 0.1441883601 & \text{if } c = 0 \\ 0.3231036988 & \text{if } c = 1 \\ 0.5327079409 & \text{if } c = 2 \end{cases}$$

اگر بخواهیم این مقادیر را در فضای  $\mathbb{R}^3$  بنویسیم:

$$P_{Y|X_r}(Y|x_r = 1) = [0.1441883601, 0.3231036988, 0.5327079409] \quad (32)$$

جمع این سه مقدار با تقریب بسیار خوبی برابر ۱ می شود.

## ۴-۲. بخش d: تحلیل بخش های قبلی

از هر بخش می توان چیزهایی یافت و تحلیل هایی ارائه کرد. ابتدا باید اشاره کرد که در بخش های قبلی از قانون بیز استفاده شد. مخرج کسر در قانون بیز با توجه به ثبات بودن برای تمامی کلاس ها (بسته به هر بخش) می تواند از تحلیل کنار گذاشته شود زیرا صرفا به عنوان یک ضریب نرمالیزاسیون عمل می کند.

یکی از راه های تحلیل مقادیر هر بخش این است که بینیم هر prior چه تاثیری بر likelihood گذاشته تا خروجی posterior را داده است.

### تحلیل بخش a

از محاسبات پایانی بخش a یعنی روابط (۲۲) و (۲۳) می بینیم که کلاس ۰ بیشترین و کلاس ۱ کمترین احتمال پسین را دارد. علت آن این است که توزیع احتمال  $X_1 = 0$  بیشترین احتمال را به کلاس های ۰ و ۲ و کمترین را به کلاس ۰ می دهد و در مقابل برای  $X_2 = 0$  با توجه به مقدار  $\mu$  احتمال کلاس ۱ کمی بیشتر می شود. از طرفی prior بیشترین احتمال را به کلاس ۰ می دهد.

ترکیب این ها باعث می شود که کلاس ۰ همچنان محتمل ترین کلاس باشد. کلاس ۱ به دلیل احتمال کمی که  $X_1 = 0$  و prior به آن می دهند، تاثیر احتمال بالایی که ویژگی  $X_2 = 0$  به کلاس ۱ می دهد را کم می کند و باعث می شود کم احتمال ترین کلاس باشد.

### تحلیل بخش b

در این بخش تنها ویژگی  $X_1 = 0$  مورد استفاده قرار می گیرد. در محاسبات احتمال پسین در (۲۶) و (۲۷) می بینیم که کلاس ۰ بیشترین احتمال و کلاس ۱ کمترین احتمال را می گیرند. می توان گفت که کلاس ۰ قوی تر و کلاس ۱ ضعیف تر از بخش a ظاهر می شوند و کلاس ۲ کمترین تغییر (حدود ۰.۲ افزایش) را دارد. حضور

$X_1 = 0$  به نفع کلاس ۰ و ۲ است و به ضرر کلاس ۱ می باشد. در مقابل می توان گفت، حضور  $X_2 = 0$  کمی به نفع کلاس ۱ بود.

### تحلیل بخش c

در این بخش تنها ویژگی  $X_2 = 1$  مورد استفاده قرار می گیرد. در محاسبات احتمال پسین در (۳۱) و (۳۲) می بینیم که کلاس ۲ بیشترین و کلاس ۰ کمترین احتمال را دارد. مقدار  $X_2 = 1$  دقیقاً برابر با میانگین کلاس ۲ می باشد و باعث می شود احتمال کلاس ۲ به شدت افزایش یابد و به نفع آن باشد. در مقابل هر چی از مقدار میانگین  $\mu = 1$  دور می شیم (به علت ساختار نمودار توزیع گوسی) احتمال کلاس های دیگر کمتر می شود. کما اینکه می بینیم کمترین احتمال به کلاس ۰ که میانگینش ۱- است می رسد.

پس به طور خلاصه  $X_2 = 1$ ، به نفع کلاس ۲ بوده و آن را غالب می کند. کلاس ۱ را در رتبه ی دوم و کلاس ۰ را به پایین ترین مقدار می رساند.

## پرسش ۳ – Multinomial Naive Bayes

این سوال هم مانند سوال قبل درباره‌ی naive bayes می باشد ولی محاسبات سوال قبل را ندارد و صرفاً اثبات روابط است. مباحث مربوط به این سوال را از جزوه‌ی درس و تدریس دکتر ابولقاسمی و منبع [۱] مرور کردم.

### ۱-۳. بخش a: $P(y = 1|x)$

با استفاده از قانون بیز داریم:

$$P(y = 1|x) = \frac{P(x|y = 1)P(y = 1)}{P(x|y = 0)P(y = 0) + P(x|y = 1)P(y = 1)} \quad (33)$$

اگر document مورد بررسی  $L$  کلمه داشته باشد. آنگاه  $L = \sum_{i=1}^d c_i$  است. یعنی دقیقاً برابر جمع تعداد دفعات تکرار همه‌ی کلمات است (بدیهی است). توزیع این document مطابق توزیع چندجمله‌ای (multinomial) به صورت زیر است.

برای کلاس ۰:

$$P(x|y = 0) = \frac{(\sum_{i=1}^d c_i)!}{\prod_{i=1}^d c_i!} \prod_{i=1}^d P_{i|0}^{c_i} \quad (34)$$

برای کلاس ۱:

$$P(x|y = 1) = \frac{(\sum_{i=1}^d c_i)!}{\prod_{i=1}^d c_i!} \prod_{i=1}^d P_{i|1}^{c_i} \quad (35)$$

توجه کنید که notation های  $P_{i|0}$  و  $P_{i|1}$  را خود سوال معرفی کرده است.

حال (۳۴) و (۳۵) را در (۳۳) جایگذاری می‌کنم.

$$P(y = 1|x) = \frac{\frac{(\sum_{i=1}^d c_i)!}{\prod_{i=1}^d c_i!} \prod_{i=1}^d P_{i|1}^{c_i} \cdot P_1}{\frac{(\sum_{i=1}^d c_i)!}{\prod_{i=1}^d c_i!} \prod_{i=1}^d P_{i|1}^{c_i} \cdot P_1 + \frac{(\sum_{i=1}^d c_i)!}{\prod_{i=1}^d c_i!} \prod_{i=1}^d P_{i|0}^{c_i} \cdot P_0} \quad (36)$$

حال با ساده سازی عبارات مشترک صورت و مخرج داریم:

$$P(y = 1|x) = \frac{\prod_{i=1}^d P_{i|1}^{c_i} \cdot P_1}{\prod_{i=1}^d P_{i|1}^{c_i} \cdot P_1 + \prod_{i=1}^d P_{i|0}^{c_i} \cdot P_0} \quad (37)$$

تا اینجا ساده سازی این عبارت تمام شده؛ اما من در درس تئوری آشکارسازی تخمین (دانشجوی مخابرات هستم) دیدم که برای راحتی بیشتر از معادل های  $P_1 = \pi$  و  $P_0 = 1 - \pi$  در اکثر مواقع استفاده می شود. برای همین رابطه‌ی (۳۷) را به این صورت هم می توان نوشت.

$$P(y = 1|x) = \frac{\pi \cdot \prod_{i=1}^d P_{i|1}^{c_i}}{\pi \cdot \prod_{i=1}^d P_{i|1}^{c_i} + (1 - \pi) \prod_{i=1}^d P_{i|0}^{c_i}} \quad (38)$$

واضحا، روابط (37) و (38) عین یکدیگرند.

### ۲-۳. بخش b: Linearity of Decision Boundary

در مبحث Two-Category Classification از جزوه ی درس برای likelihood ratio نامساوی زیر را دیدیم.

$$\frac{P(x|\omega_1)}{P(x|\omega_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(\omega_2)}{P(\omega_1)} \quad (39)$$

در درس تئوری آشکارسازی و تخمین بارها دیدیم که وقتی مقدار loss ها مشخص نشده اند، باید آن ها را اصطلاحاً uniform loss در نظر بگیریم (از اصطلاحات تئوری آشکارسازی تخمین). به این معنی که کسر مربوط به loss ها در نامساوی (39) برابر ۱ می شود.

حالا مفروضات مسئله را در رابطه ی (39) اعمال می کنیم.

$$\frac{P(x|y_1)}{P(x|y_2)} > \frac{1 - \pi}{\pi} \quad (40)$$

با استفاده از (34) و (35) در (40) داریم:

$$\prod_{i=1}^d \frac{P_{i|1}^{c_i}}{P_{i|0}^{c_i}} > \frac{1 - \pi}{\pi} \quad (41)$$

از دو طرف این رابطه لگاریتم می گیریم.

$$\sum_{i=1}^d c_i \log \left( \frac{P_{i|1}}{P_{i|0}} \right) > \log \left( \frac{1 - \pi}{\pi} \right) \quad (42)$$

$$\sum_{i=1}^d c_i \log \left( \frac{P_{i|1}}{P_{i|0}} \right) + \log \left( \frac{\pi}{1 - \pi} \right) > 0 \quad (43)$$

حال اگر  $x = \sum_{i=1}^d c_i$  و  $\theta_i = \log \left( \frac{P_{i|1}}{P_{i|0}} \right)$  و  $\theta_0 = \log \left( \frac{\pi}{1 - \pi} \right)$  و همچنین  $\theta = (\theta_1, \theta_2, \dots, \theta_d)$  می توانیم (43) را به صورت زیر بازنویسی کنیم.

$$\sum_{i=1}^d c_i \log \left( \frac{P_{i|1}}{P_{i|0}} \right) + \log \left( \frac{\pi}{1 - \pi} \right) = \theta^T x + \theta_0 > 0 \quad (44)$$

که با در نظر گرفتن مرز تصمیم گیری به صورت زیر خواهد بود.

$$\theta^T x + \theta_0 = 0 \quad (45)$$

که واضحا این یک رابطه ی خطی نسبت به  $x$  بوده که بایاس  $\theta_0$  دارد. (توجه کنید تغییر متغیر های  $x, \theta_i, \theta., \theta$  با توجه به بخش بعدی انتخاب شده‌اند).

### ۳-۳. بخش c: Logistic Function

در بخش a رابطه ی (۳۴) را دیدیم. می خواهیم به نحوی رابطه ی (۴۵) را در (۳۸) ایجاد کنم. با تقسیم صورت و مخرج آن بر  $\prod_{i=1}^d P_{i|1}^{c_i}$  خواهیم داشت:

$$P(y = 1|x) = \frac{\frac{\pi}{1-\pi} \prod_{i=1}^d \left(\frac{P_{i|1}}{P_{i|0}}\right)^{c_i}}{1 + \frac{\pi}{1-\pi} \prod_{i=1}^d \left(\frac{P_{i|1}}{P_{i|0}}\right)^{c_i}} \quad (46)$$

با توجه به تعریف متغیر ها در (۴۴) می توان نوشت:

$$\prod_{i=1}^d \left(\frac{P_{i|1}}{P_{i|0}}\right)^{c_i} = e^{\theta^T x}, \quad \frac{\pi}{1-\pi} = e^{\theta_0} \quad (47)$$

پس می توانیم (۴۶) را اینگونه بنویسیم.

$$P(y = 1|x) = \frac{e^{\theta^T x + \theta_0}}{1 + e^{\theta^T x + \theta_0}} = \frac{1}{1 + e^{-(\theta^T x + \theta_0)}} \quad (48)$$

پس اثبات شد این احتمال شرطی یک تابع لجستیک است.



۴-۱. بخش a:

در جزوه ی درس دیدیم که likelihood function به صورت زیر است.

$$\mathcal{L}(\theta|D) = \prod_{i=1}^N p(X_i|\theta) \quad (49)$$

با توجه به توزیعی که در صورت سوال داده شده، اگر حتی یکی از نمونه ها خارج از لوزی باشد مقدار likelihood function صفر می شود. پس باید تک تک نمونه ها داخل لوزی باشند.

$$\|X_i\|_1 \leq \theta, \quad \forall i \quad (50)$$

واضح است که اگر بزرگ ترین نمونه در شرط بالا صدق کند، بقیه ی نمونه ها هم صدق خواهند کرد. پس می توان این شرط به صورت زیر نوشت.

$$\max_i \|X_i\|_1 \leq \theta \quad (51)$$

با توجه به مستقل بودن نمونه ها و با توجه به (49) می توان نوشت.

$$\mathcal{L}(\theta|D) = \prod_{i=1}^N p(X_i|\theta) = \prod_{i=1}^N \frac{1}{2\theta^2} = \frac{1}{(2\theta^2)^N} \quad (52)$$

تابع بدست آمده یک تابع اکیدا نزولی است که بیشترین مقدار آن به ازای کمترین مقدار  $\theta$  بدست می آید. از طرفی در (51) کمترین مقدار  $\theta$  را یافتیم. پس:

$$\hat{\theta} = \max_{i=1, \dots, N} \|X_i\|_1 \quad (53)$$

این تخمین MLE می باشد.

۴-۲. بخش b:

توزیع  $\text{gamma}(\alpha, \beta)$  به صورت زیر است.

$$f(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-\frac{x}{\beta}} \quad (54)$$

تابع likelihood آن مطابق (49) به صورت زیر محاسبه می شود.

$$\begin{aligned}\mathcal{L}(\beta|D) &= \prod_{i=1}^N p(X_i|\beta) = \prod_{i=1}^N \frac{1}{\Gamma(\alpha)\beta^\alpha} X_i^{\alpha-1} \exp\left(-\frac{X_i}{\beta}\right) \\ &= \left(\frac{1}{\Gamma(\alpha)}\right)^N \beta^{-\alpha N} \left(\prod_{i=1}^N X_i^{\alpha-1}\right) \exp\left(-\frac{1}{\beta} \sum_{i=1}^N X_i\right)\end{aligned}\quad (55)$$

در این رابطه  $D$  بردار شامل نمونه هاست. می خواهیم ماکزیمم این تابع را پیدا کنیم. می دانیم لگاریتم تابعی اکیدا صعودی است پس می توانیم به جای خود  $\mathcal{L}(\beta|D)$  ماکزیمم لگاریتم آن را بیابیم.

$$\log(\mathcal{L}(\beta|D)) = N \log\left(\frac{1}{\Gamma(\alpha)}\right) - \alpha N \log(\beta) + (\alpha - 1) \sum_{i=1}^N \log(X_i) - \frac{1}{\beta} \sum_{i=1}^N X_i \quad (56)$$

حال از (56) مشتق می گیریم و مساوی صفر قرار می دهیم تا نقطه ای که ماکزیمم در آن رخ می دهد را بیابیم. این مشتق را با  $A$  نمایش می دهیم.

$$A = \frac{d \log(\mathcal{L}(\beta|D))}{d\beta} = -\frac{\alpha N}{\beta} + \frac{1}{\beta^2} \sum_{i=1}^N X_i = 0 \quad (57)$$

$$-\alpha N \beta + \sum_{i=1}^N X_i = 0 \quad (58)$$

در نتیجه خواهیم داشت:

$$\hat{\beta} = \frac{1}{\alpha N} \sum_{i=1}^N X_i \quad (59)$$

طبق (59) مقدار تخمین ML برای  $\beta$  بدست آمد. بخشی که در ادامه می نویسم نه در صورت سوال خواسته شده و نه در کلاس درس مطرح شده اما من می نویسم. در درس تئوری آشکارسازی و تخمین، خواندیم که حتما باید مشتق دوم را حساب کنیم تا از مقعر بودن تابع و ماکزیمم بودن نقطه ی بدست آمده در تخمین مطمئن شویم.

مشتق دوم با استفاده از (57) به صورت زیر بدست می آید.

$$\frac{dA}{d\beta} = \frac{d^2 \log(\mathcal{L}(\beta|D))}{d\beta^2} = \frac{\alpha N}{\beta^2} - \frac{2}{\beta^3} \sum_{i=1}^N X_i \quad (60)$$

مقدار تخمین (59) را در (60) قرار می دهیم.

$$\frac{\alpha N}{\beta^2} - \frac{2}{\beta^3} (\alpha N \beta) = -\frac{\alpha N}{\beta^2} < 0 \quad (61)$$

توجه کنید  $\beta^2$  و  $N$  ذاتا مقادیر مثبت هستند و  $\alpha$  هم فقط مقادیر مثبت می گیرید. (می توانید در مورد توزیع گاما در منبع [۲] بخوانید.) پس مشتق دوم منفی شد در نتیجه تابع مقعر بوده و نقطه ی بدست آمده در تخمین، ماکزیمم است.

در این سوال مفهوم مزدوج بودن توزیع پیشین را مرور و از تخمین MAP استفاده می کنیم.

### ۵-۱. بخش a: Conjugacy of the Dirichlet Prior

توزیع دیریکله با پارامتر  $u_i$  به صورت زیر تعریف شده.

$$\mathbb{P}(p_1, p_2, \dots, p_\epsilon) = \frac{\Gamma(\sum_{i=1}^{\epsilon} u_i)}{\prod_{i=1}^{\epsilon} \Gamma(u_i)} \prod_{i=1}^{\epsilon} p_i^{u_i-1} \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right) \quad (62)$$

توزیع دیریکله توزیع پیشین است. توزیع (multinomial) likelihood به صورت زیر است.

$$P(n_1, \dots, n_\epsilon | p) = \frac{N!}{\prod_{i=1}^{\epsilon} n_i!} \prod_{i=1}^{\epsilon} p_i^{n_i} \quad (63)$$

در (۶۳) مقدار  $N = 1000$  تعداد دفعات پرتاب تاس است.

هدف این است که نشان دهیم توزیع پسین هم توزیع دیریکله دارد، پس روابط را در قانون بیز می گذاریم و ادامه می دهیم.

$$P(p | \{n_i\}) = \frac{p(\{n_i\} | p) \mathbb{P}(p)}{p(\{n_i\})} \quad (64)$$

توجه کنید که مخرج در (۶۴) صرفاً نقش یک ضریب نرمالیزاسیون دارد. پس می توانیم از آن صرف نظر کنیم.

$$\begin{aligned} p(\{n_i\} | p) \mathbb{P}(p) &= \left( \frac{N!}{\prod_{i=1}^6 n_i!} \prod_{i=1}^6 p_i^{n_i} \right) \left( \frac{\Gamma(\sum_{i=1}^6 u_i)}{\prod_{i=1}^6 \Gamma(u_i)} \prod_{i=1}^6 p_i^{u_i-1} \delta\left(\sum_{i=1}^6 p_i - 1\right) \right) \quad (65) \\ &= \frac{N! \Gamma(\sum_{i=1}^{\epsilon} u_i)}{(\prod_{i=1}^{\epsilon} n_i!) (\prod_{i=1}^{\epsilon} \Gamma(u_i))} \left( \prod_{i=1}^{\epsilon} p_i^{n_i+u_i-1} \right) \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right) \end{aligned}$$

در رابطه ی (۶۵) کسر پایانی تنها یک ضریب است که نسبت به  $p$  ها ثابت است. پس می توانیم از آن ها صرف نظر کنیم.

$$P(p | \{n_i\}) \propto \left( \prod_{i=1}^{\epsilon} p_i^{n_i+u_i-1} \right) \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right) \quad (66)$$

رابطه ی بدست آمده مشابه توزیع دیریکله در (۶۲) است. پارامتر توزیع دیریکله بدست آمده برای توزیع پسین، به صورت زیر است.

$$\alpha_i = u_i + n_i \quad (67)$$

پس در کل توزیع پسین یک توزیع دیریکله به صورت زیر دارد.

$$\begin{aligned}
 P(p|\{n_i\}) &= \frac{\Gamma(\sum_{i=1}^{\epsilon} \alpha_i)}{\prod_{i=1}^{\epsilon} \Gamma(\alpha_i)} \prod_{i=1}^{\epsilon} p_i^{\alpha_i-1} \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right) \\
 &= \frac{\Gamma(\sum_{i=1}^{\epsilon} (u_i + n_i))}{\prod_{i=1}^{\epsilon} \Gamma(u_i + n_i)} \prod_{i=1}^{\epsilon} p_i^{u_i+n_i-1} \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right)
 \end{aligned} \tag{68}$$

## ۵-۲. بخش b: Posterior Distribution

با جایگذاری مقادیر مشخص شده در صورت سوال در (۶۸)، می توانیم به صورت زیر احتمال خواسته شده را بنویسیم.

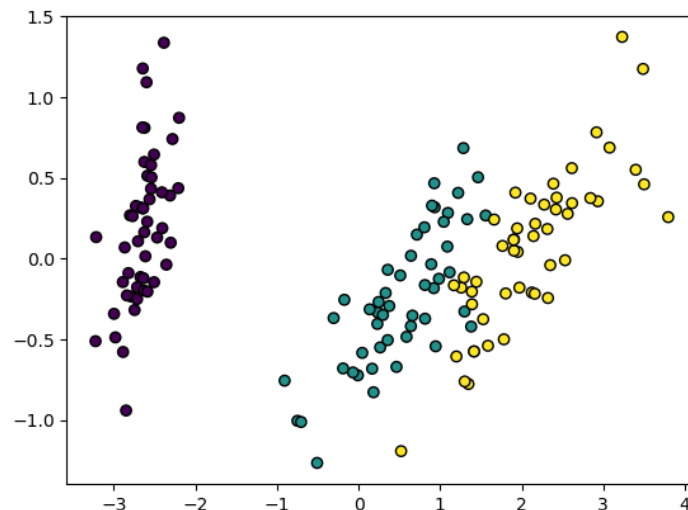
$$\begin{aligned}
 P(p_1, \dots, p_{\epsilon} | n_1, \dots, n_{\epsilon}) &= \frac{\Gamma(\sum_{i=1}^{\epsilon} (i + n_i))}{\prod_{i=1}^{\epsilon} \Gamma(i + n_i)} \prod_{i=1}^{\epsilon} p_i^{i+n_i-1} \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right) \\
 &= \frac{\Gamma(1.21)}{\prod_{i=1}^{\epsilon} \Gamma(u_i + n_i)} \prod_{i=1}^{\epsilon} p_i^{i+n_i-1} \delta\left(\sum_{i=1}^{\epsilon} p_i - 1\right)
 \end{aligned} \tag{69}$$

توجه کنید که خود سوال گفته  $u_i = i$ .

## پرسش ۶ – Gaussian Discriminant Analysis

در این تمرین کامپیوتری پارامترهای توزیع شرطی (روش پارامتریک تخمین توزیع) را محاسبه می کنیم، داده ها را با استفاده از یک discriminant function به کلاس های مختلف کلاس بندی می کنیم و مرز تصمیم گیری را تحلیل می کنیم.

داده های تولید شده برای دیتاست در نمودار زیر رسم شده اند.



شکل ۱. نمودار دیتاست سوال ۶

### ۶-۱. بخش a: محاسبه ی میانگین و کوواریانس

در این بخش باید با استفاده از داده ها، پارامترهای توزیع گوسی bivariate را محاسبه کنیم. خود سوال گفته که ماتریس کوواریانس برای هر دو کلاس مشترک است پس باید برای ساخت ماتریس کوواریانس مشترک، از pooling covariance استفاده شود. یعنی به این شکل:

$$\Sigma = \frac{1}{N_1 + N_2} \left( \sum_{i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T + \sum_{i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T \right) \quad (70)$$

بردار میانگین و ماتریس کوواریانس هر کلاس به صورت زیر است.

Mean of class ۱: [ ۰,۵۳۳۲۰۶۵۷ -۰,۲۴۵۵۴۹۸۳]

Mean of class ۲: [۲,۱۰۹۲۰۸۸۹ ۰,۰۵۴۶۶۴۷۹]

Covariance of class ۱:

[[۰,۳۴۸۳۶۱۵۵ ۰,۱۹۴۷۶۷۳۲]

$$[0,19476732 \ 0,18146878]$$

Covariance of class ۲:

$$[[0,48832634 \ 0,27033776]$$

$$[0,27033776 \ 0,22862663]]$$

و نهایتاً کوواریانس مشترک به صورت زیر است.

Shared covariance matrix ( $\Sigma$ ):

$$[[0,41834394 \ 0,23255254]$$

$$[0,23255254 \ 0,20504771]]$$

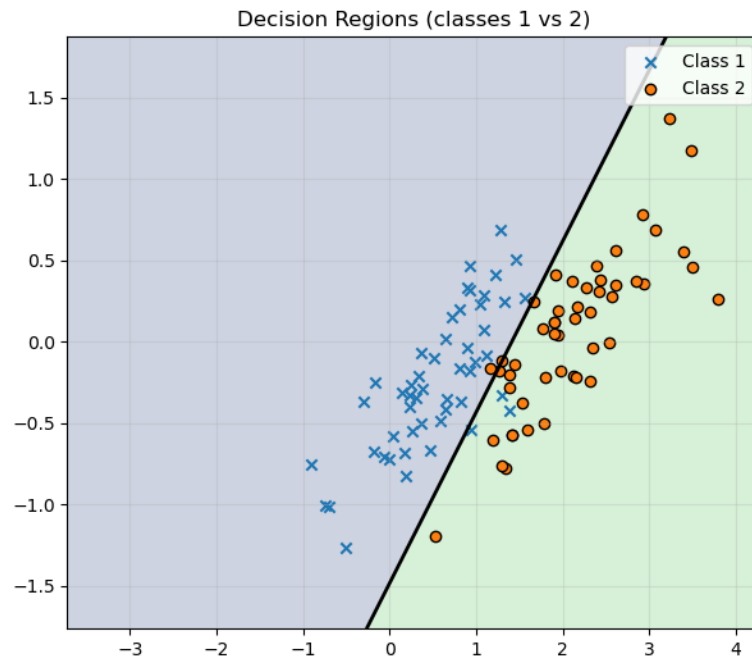
## ۲-۶. بخش b: Discriminant function

با توجه به توزیع چند متغیره و یکسان بودن کوواریانس برای هر دو کلاس، تابع discriminant به صورت زیر تعریف می شود.

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma^{-1}(x - \mu_i) + \ln P(\omega_i) \quad (۷۱)$$

## ۳-۶. بخش c: Decision Regions

ناحیه های تصمیم گیری برای هر کلاس در تصویر مشخص شده اند.



شکل ۲. ناحیه های تصمیم گیری دو کلاس

#### ۴-۶. بخش d: خط مرز تصمیم گیری

مرز ناحیه‌ی تصمیم گیری دقیقاً جایی است که مقدار discriminant function به ازای هر دو کلاس برابر می شود.

$$g_1(x) = g_2(x) \quad (72)$$

$$\begin{aligned} & -\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1) + \ln P(\omega_1) \\ & = -\frac{1}{2}(x - \mu_2)^T \Sigma^{-1}(x - \mu_2) + \ln P(\omega_2) \end{aligned} \quad (73)$$

به علت یکسان بودن ماتریس کوواریانس برای هر دو کلاس، جملات مرتبه دوم  $x$  در دو طرف تساوی، برابر شده و حذف می شوند. در نهایت به رابطه ی خط زیر می رسیم.

$$w^T x + b = 0 \quad (74)$$

در این رابطه  $w^T$  برداری وزن ها و  $b$  بایاس (عرض از مبدا) خط است.

$$w = \Sigma^{-1}(\mu_1 - \mu_2) \quad (75)$$

$$b = -\frac{1}{2}(\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2) + \ln \left( \frac{P(\omega_1)}{P(\omega_2)} \right) \quad (76)$$



با استفاده از کد پایتون نوشته شده هر کدام از مقادیر (۷۵) و (۷۶) به صورت زیر بدست آمد.

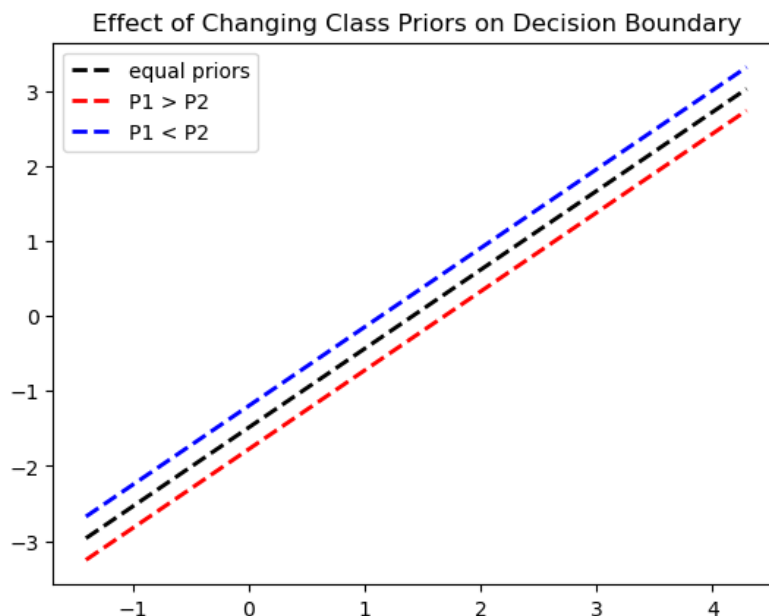
w (decision boundary vector):  $[-۷,۹۹۱۸۵۲۹۵ \quad ۷,۵۹۹۷۴۹۰۲]$

b (decision boundary bias):  $۱۱,۲۸۴۲۳۷۱۲۷۳۶,۱۵۴$

## ۵-۶. بخش e: تاثیر تابع پیشین در نواحی تصمیم گیری

تصویر زیر مرز تصمیم گیری را به ازای زوج  $(prior1, prior2)$  با مقادیر مختلف  $(۰.۵, ۰.۵)$ ،  $(۰.۱, ۰.۹)$  و  $(۰.۹, ۰.۱)$  نشان می دهد.

می بینیم که وقتی احتمال پیشین کلاس ۱ بیشتر است، مرز تصمیم به سمت کلاس ۲ می رود و ناحیه ی تصمیم گیری کلاس ۱ بزرگتر می شود (خطای نوع اول پیش می آید). وقتی احتمال پیشین کلاس ۲ بیشتر است، مرز تصمیم به سمت کلاس ۱ می رود و ناحیه ی تصمیم گیری کلاس ۲ بزرگتر می شود (خطای نوع دوم پیش می آید).



شکل ۳. تاثیر احتمال پیشین بر مرز تصمیم گیری

## ۶-۶. بخش f: تکرار سوال برای سه کلاس

حال باید بخش های a تا c را برای سه کلاس با فرض ماتریس کوواریانس نابرابر تکرار کرد. داده های سه کلاس در شکل ۱ رسم شده اند.

### ۶-۶-۱. بخش a: محاسبه ی میانگین و کوواریانس

بردار میانگین و ماتریس کوواریانس هر کلاس به صورت زیر است.

Mean of class ۰:  $[-۲,۶۴۲۴۱۵۴۶ \quad ۰,۱۹۰۸۸۵۰۵]$

Mean of class ۱: [ ۰,۵۳۳۲۰۶۵۷ -۰,۲۴۵۵۴۹۸۳]

Mean of class ۲: [۲,۱۰۹۲۰۸۸۹ ۰,۰۵۴۶۶۴۷۹]

Covariance of class ۰:

[[۰,۰۴۸۰۴۱۵ ۰,۰۵۴۹۲۱۸۵]

[۰,۰۵۴۹۲۱۸۵ ۰,۲۱۳۳۴۳۳۶]]

Covariance of class ۱:

[[۰,۳۴۸۳۶۱۵۵ ۰,۱۹۴۷۶۷۳۲]

[۰,۱۹۴۷۶۷۳۲ ۰,۱۸۱۴۶۸۷۸]]

Covariance of class ۲:

[[۰,۴۸۸۳۲۶۳۴ ۰,۲۷۰۳۳۷۷۶]

[۰,۲۷۰۳۳۷۷۶ ۰,۲۲۸۶۲۶۶۳]]

#### ۲-۶-۶. بخش b: Discriminant function

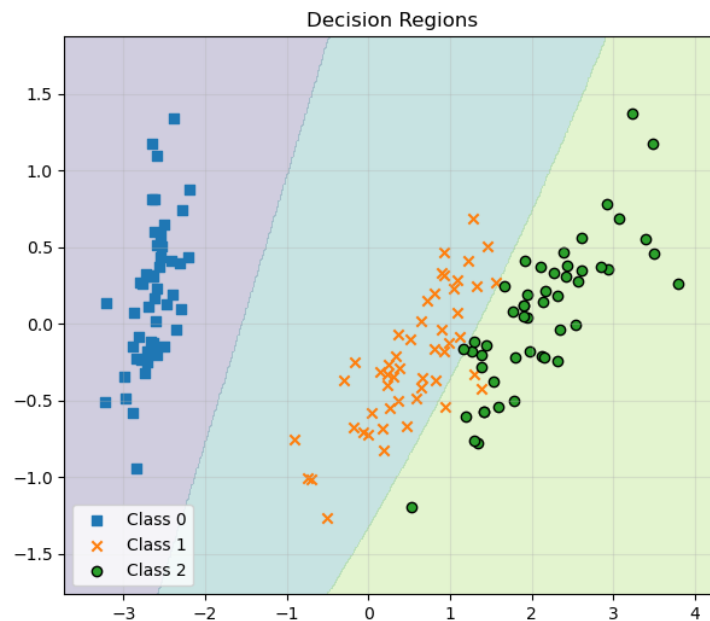
اینجا باید از discriminant function برای حالتی که ماتریس کوواریانس کلاس ها متفاوت است استفاده کنیم. رابطه ی آن به صورت زیر است.

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i) \quad (۷۷)$$

توجه کنید که جمله ی  $\frac{d}{2} \ln 2\pi$  در تمامی کلاس ها یکسان می شود پس می توانیم در نوشتن کد، از آن صرف نظر کنیم.

#### ۳-۶-۶. بخش c: Decision Regions

ناحیه های تصمیم گیری سه کلاس به صورت زیر می شود.



شکل ۴. ناحیه های تصمیم گیری سه کلاس

## ۶-۷. بخش i: تکرار سوال با داده های کمتر

در اینجا باید بخش های  $a$  تا  $c$  را با فرض تفاوت ماتریس کوواریانس توزیع ها و همچنین با تقسیم داده ها و استفاده از ۸۰ درصد نمونه ها و در نظر گرفتن هزینه rejection تکرار کنیم. اول از هر کاری باید داده ها را تقسیم کنیم و ۸۰ درصد مورد نظر را مشخص کنیم در غیر این صورت احتمالا leakage خواهیم داشت.

### ۶-۷-۱. بخش a: محاسبه ی میانگین و کوواریانس

بردار میانگین و ماتریس کوواریانس هر کلاس به صورت زیر است.

Mean of class ۰:  $[-۲,۶۶۱۷۸۵۷۸ \ ۰,۲۰۴۵۹۱۸۹]$

Mean of class ۱:  $[۰,۵۲۵۷۰۷۶ \ -۰,۲۴۰۹۰۶۴۸]$

Mean of class ۲:  $[۲,۱۵۹۱۱۵۷۴ \ ۰,۰۹۰۵۹۵۰۸]$

Covariance of class ۰:

$[[۰,۰۴۲۳۲۱۶۷ \ ۰,۰۴۴۲۹۰۴۳]$

$[۰,۰۴۴۲۹۰۴۳ \ ۰,۲۱۵۶۶۵۰۴]]$

Covariance of class ۱:

$[[۰,۳۳۳۱۰۴۴۸ \ ۰,۱۹۱۱۸۴۳۷]$

$[۰,۱۹۱۱۸۴۳۷ \ ۰,۱۹۸۳۰۴۱۵]]$

Covariance of class ۲:

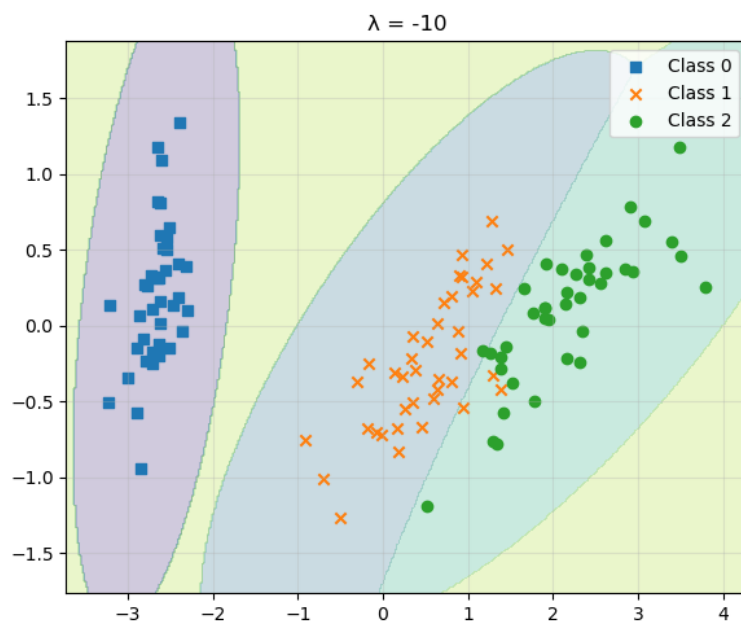
$[[۰,۵۱۲۱۹۶۲۱ \ ۰,۲۶۰۳۴۸۳۵]$

**۶-۷-۲. بخش b: Discriminant function**

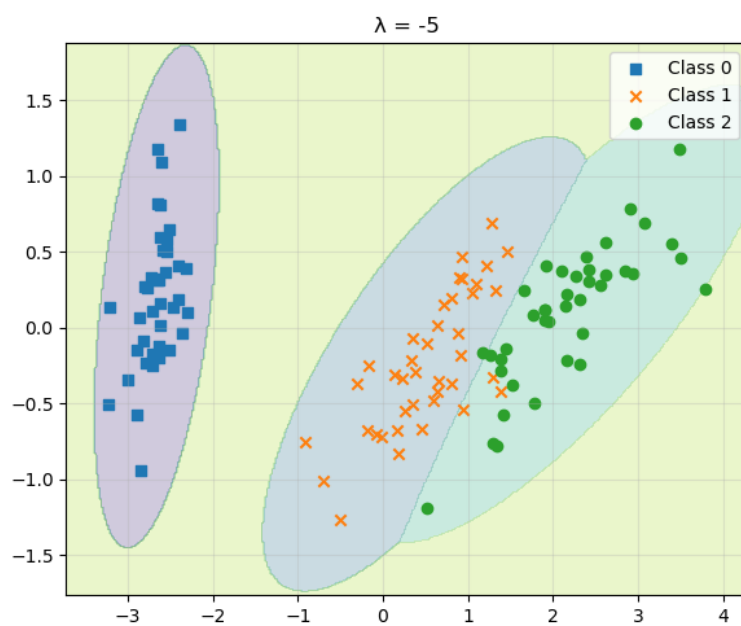
اینجا باید از discriminant function برای حالتی که ماتریس کوواریانس کلاس ها متفاوت است استفاده کنیم ( رابطه‌ی (۷۷) ). همان طور که صورت مسئله خواسته، هزینه ی rejection در نظر گرفته می شود.

**۶-۷-۳. بخش c: Decision Regions**

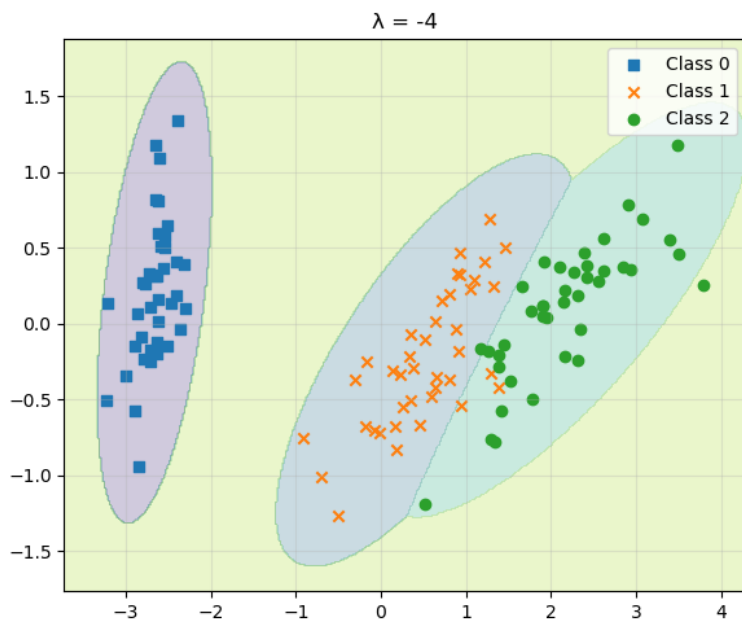
ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda$  های مختلف به صورت زیر می شود. با توجه به اشکال ۵ الی ۱۱، هر چه مقدار لاندا بزرگ تر می شود ناحیه های reject هم بزرگ تر شده و داده هایی که در این ناحیه قرار می گیرند بیشترند. به ازای مقادیر ۱۰- و ۵- برای لاندا، ناحیه های reject بسیار محافظه کارانه هستند. در مقابل به ازای ۲-، ۱- و ۰ نواحی reject بسیار بزرگ شده و بخش زیادی از فضای تصمیم از بین می رود. اما به نظر به ازای مقادیر ۴- و ۳- تعادل خوبی بین ناحیه ی reject و نواحی تصمیم گیری وجود دارد. اگر قرار بر انتخاب لاندا باشد، بهترین مقدار ها ۴- و ۳- هستند.



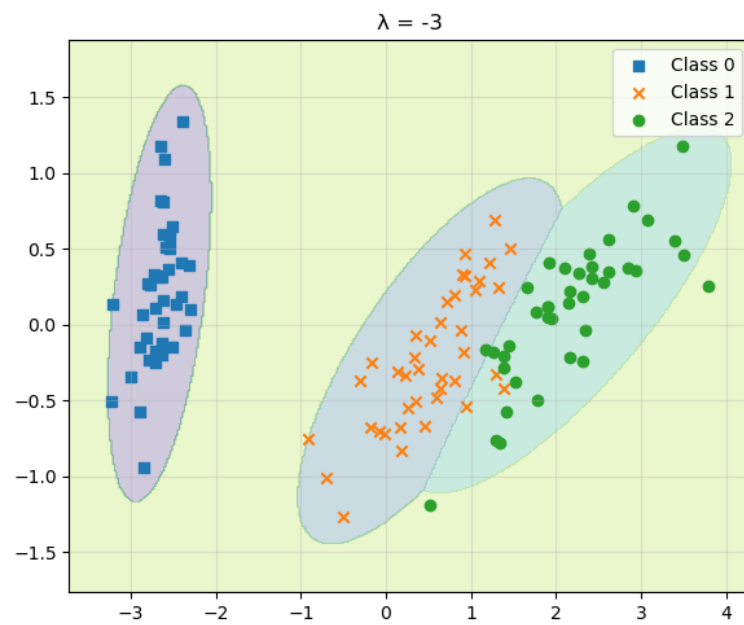
شکل ۵. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -10$



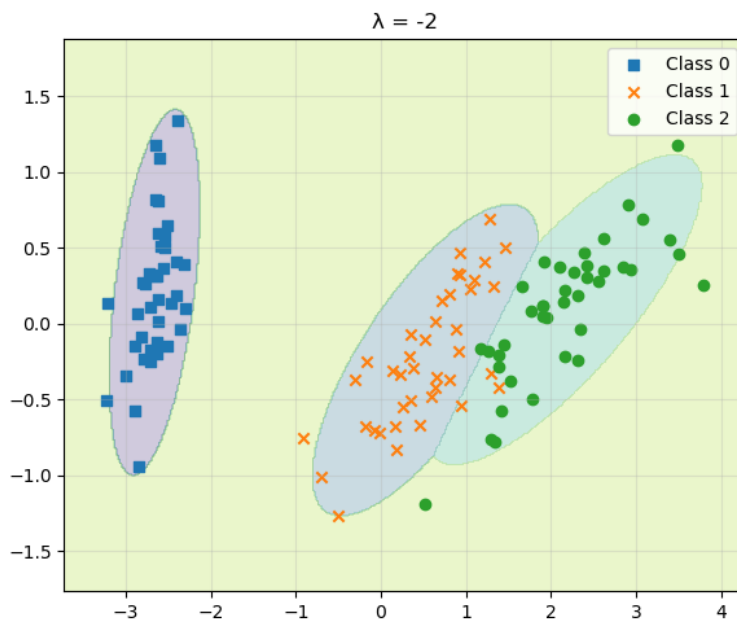
شکل ۶. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -5$



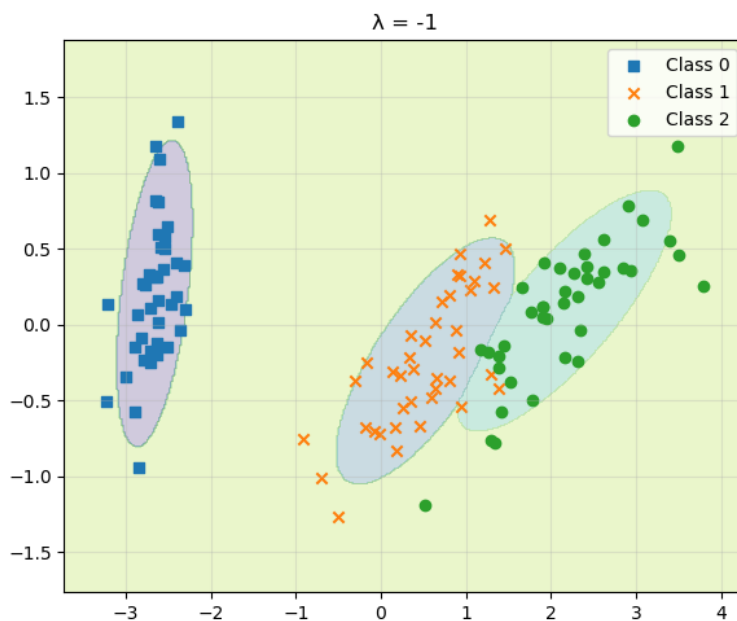
شکل ۷. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -4$



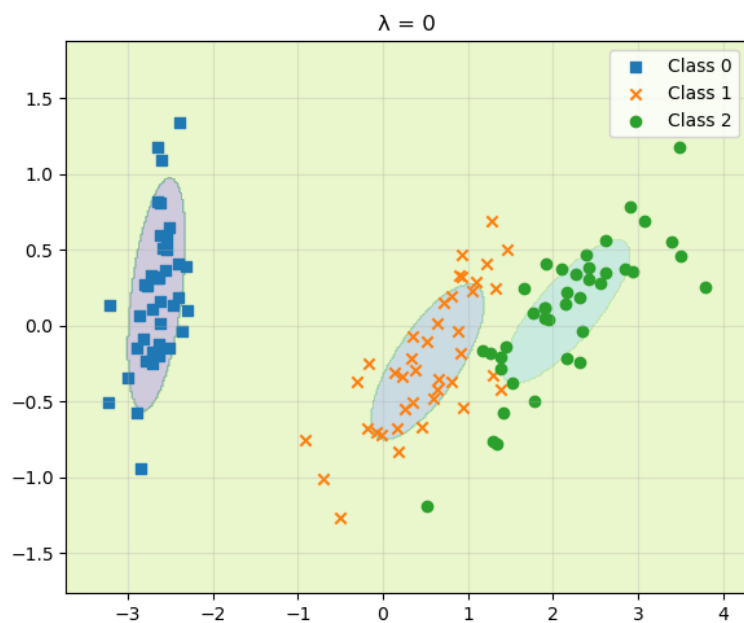
شکل ۸. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -3$



شکل ۹. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -2$



شکل ۱۰. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = -1$



شکل ۱۱. ناحیه های تصمیم گیری سه کلاس به ازای  $\lambda = 0$



## پرسش ۷ – Gaussian Naive Bayes

در این تمرین یک مدل Gaussian Naive Bayes را کامل نوشته، آموزش می دهیم و عملکرد آن را بررسی می کنیم. بخش های a تا c این تمرین کاملاً کامپیوتری بوده و جنبه پرسش تئوری ندارد که نیاز به توضیح داشته باشد.

### ۷-۱. بخش d: عملکرد مدل

تمامی مقادیر زیر در فایل evaluation.txt موجود است. در اینجا صرفاً برای کامل بودن گزارش دوباره نوشته می شود.

Accuracy: ۰,۸۳۵

Precision for class ۱: ۰,۹۵۱

Recall for class ۱: ۰,۹۷۵

Precision for class ۵: ۰,۸۷۵

Recall for class ۵: ۰,۷۷۸

## پرسش ۸ — Bag of Words — Naive Bayes Classifier

این تمرین درباره‌ی پیاده سازی طبقه بند naive bayes با MLE و MAP می شود. مفاهیم مهمی از سوالات تئوری این تمرین می توان دریافت.

### ۸-۱. بخش a: Naive Bayes Classifier — MLE

در MLE با احتمال زیر کار می کنیم.

$$p(w|c) = \frac{\text{count}_{w,c}}{N_c} \quad (78)$$

به مفهوم تعداد ظاهر شدن هر کلمه در هر کلاس به کل کلمات آن کلاس است.

- دقت مجموعه تست شما چقدر است؟ مقدار دقت روی داده های تست:

Accuracy on baseball test: ۰,۳۲

Accuracy on hockey test: ۰,۰۴

Total accuracy: ۰,۱۸

که این دقت بسیار بد است.

- در نهایت طبقه بند شما چه نتیجه‌ای پیش‌بینی کرد؟  
در خروجی زیر می بینیم که بسیاری از پست های بیسبال را به اشتباه هاکی و بسیاری از پست های هاکی را به اشتباه بیسبال طبقه بندی کرده است. که این مورد نشان دهنده ی بایاس این طبقه بند است. علت اصلی این مورد را در سوال بعدی توضیح خواهیم داد.

Predictions for baseball\_test: ['hockey', 'hockey', 'baseball', 'baseball', 'hockey',  
'baseball', 'baseball', 'hockey', 'hockey', 'hockey'] ...

Predictions for hockey\_test: ['hockey', 'baseball', 'baseball', 'baseball', 'baseball',  
'baseball', 'baseball', 'baseball', 'baseball', 'baseball'] ...

- چرا استفاده از MLE در این شرایط ایده‌ی خوبی نیست؟  
(۱) اگر در داده های آموزش یکی از کلاس ها، کلمه ای هیچ گاه ظاهر نشود آنگاه MLE احتمال آن کلمه را به شرط آن کلاس صفر در نظر می گیرد.

$$p(w|c) = 0 \quad (79)$$

در نتیجه برای هر داده ی تستی که شامل آن کلمه ی دیده نشده باشد، خواهیم داشت:

$$p(w|c) = 0 \Rightarrow \log(p(w|x)) = -\infty \quad (80)$$

این باعث می شود مدل آن کلاس را رد کند و روی داده های تست دقت بدی داشته باشد.

۲) از آنجا که smoothing در این بخش انجام نشده (سوال نخواستہ است) احتمالات تحت تاثیر چند کلمه پرتکرار قرار می گیرد و مدل به سمت یک کلاس متمایل می شود. این دقیقا در نتایج هم دیده شد.

## ۸-۲. بخش b: Naive Bayes Classifier — MAP with Dirichlet Prior

نتیجه ی محاسبه ی  $\arg \max$  (posterior) به صورت زیر است.

$$p(w|c) = \frac{count_{w,c} + \beta}{N_c + \beta V} \quad (۸۱)$$

که در آن  $V$  تعداد کلمات متمایز (۵۸۲۲) است و  $\beta$  تعداد دیده شدن کلمه در داده های آموزش است.

## ۸-۳. بخش c: Effect of $\beta$ on Classification Accuracy

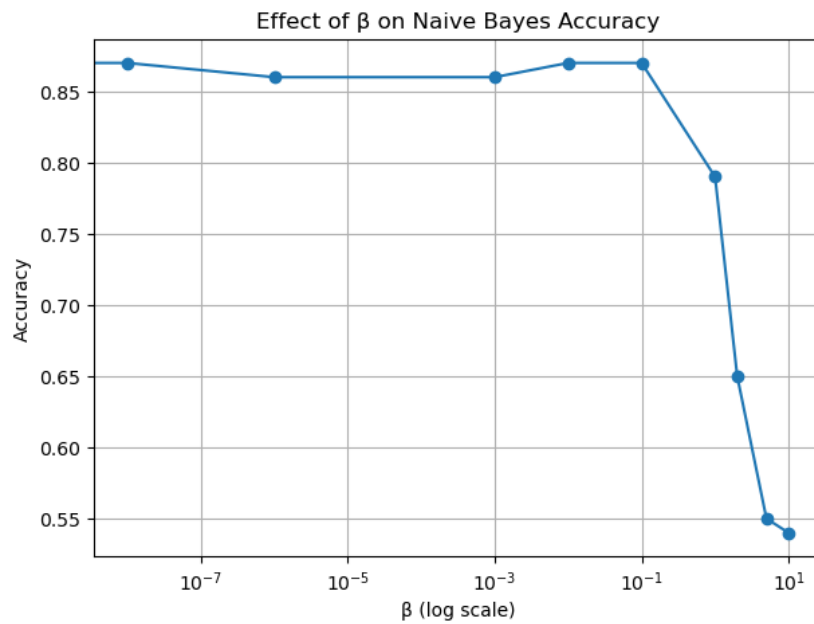
در ابتدا باید گفت که اضافه شدن  $\beta$  به کسر احتمال باعث smoothing می شود. یعنی حالا اگر کلمه ای هیچگاه در داده های آموزش دیده نشود، احتمال آن صفر نیست و ضریبی از  $\beta$  است. در جدول ۱ و شکل ۶ تاثیر مقادیر مختلف  $\beta$  روی دقت مدل را می بینیم.  $\beta$  های کوچک (از  $1e-08$  الی ۰) باعث smoothing می شوند. یعنی صرفا احتمال های صفر را اصلاح می کنند و باعث پایداری بیشتر و دقت بیشتر مدل روی داده های تست نسبت به MLE می شوند.

در مقابل  $\beta$  های بزرگ دقت مدل را به شدت کاهش می دهند. در واقع با زیاد شدن  $\beta$  شمارش هر کلمه بیش از حد smooth می شود و باعث می شود تفاوت ها واقعی بین کلاس های بیسبال و هاکی از بین برود. هدف ما این بود که با اضافه شدن  $\beta$  مشکل ظاهر نشدن برخی کلمات در آموزش مدل را رفع کنیم. ولی اگر  $\beta$  مقدار زیادی داشته باشد به این مفهوم است که ما خود خواسته احتمال دیده شدن آن کلمه را بیش از حد افزایش دادیم در حالی که این مورد نه هدف ما بوده و نه مطلوب ما. بنابراین  $\beta$  های زیاد قدرت تمایز مدل را کاهش می دهد و دقت افت می کند.

پس  $\beta$  های کوچک مطلوب هستند.

جدول ۱. تاثیر بتا بر دقت مدل

$\beta$	Accuracy
.	۰.۸۶۰۰
$1e-08$	۰.۸۷۰۰
$1e-06$	۰.۸۶۰۰
۰.۰۰۱	۰.۸۶۰۰
۰.۰۱	۰.۸۷۰۰
۰.۱	۰.۸۷۰۰
۱	۰.۷۹۰۰
۲	۰.۶۵۰۰
۵	۰.۵۵۰۰
۱۰	۰.۵۴۰۰



شکل ۱۲. تغییرات دقت نسبت به بتا

## ۸-۴. بخش d: Effect of Training Set Size

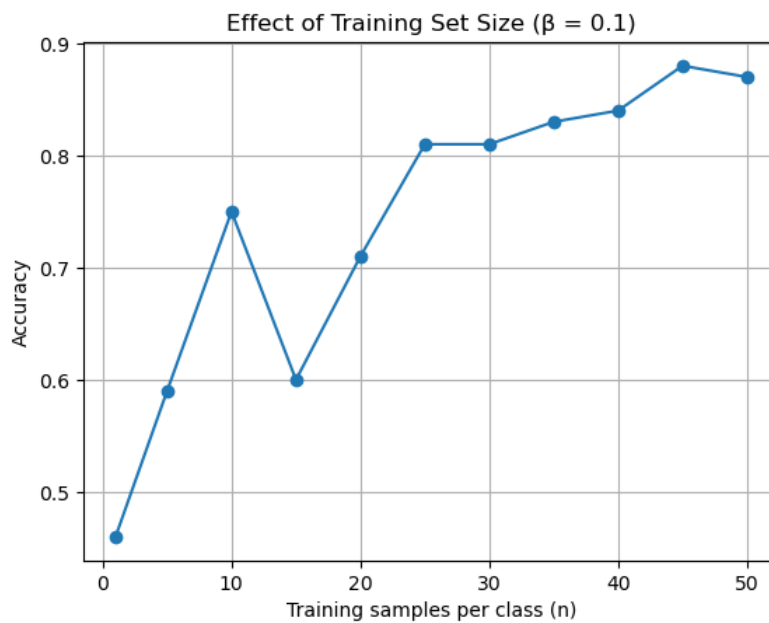
در این بخش تاثیر اندازه‌ی داده‌ی آموزش بر دقت مدل می‌سنجیم. در  $n = 1$  و  $n = 5$  دقت مدل پایین است زیرا با تعداد کمی پست، تقریباً اطلاعات قابل اعتمادی درباره‌ی توزیع کلمات بدست نمی‌آوریم. در نتیجه مدل بدون داده‌ی کافی قادر به یادگیری تمایز بین کلاس‌ها نیست. با افزایش  $n$  دقت مدل به مرور افزایش می‌یابد

کامل تری از هر کلاس بدست تمایزش بهتر می‌شود.  $n = 50$  و  $n = 15$  به دلیل نویز یا از پست‌ها باشد.

n	Accuracy
1	0.4600
5	0.5900
10	0.7500
15	0.6000
20	0.7100
25	0.8100
30	0.8100
35	0.8300
40	0.8400
45	0.8800
50	0.8700

زیرا مدل تصویر توزیع کلمات در می‌آورد و توانایی اکتی که در 15 داریم می‌تواند خاص بودن بعضی

جدول ۲. تاثیر اندازه‌ی داده‌ی آموزش بر دقت مدل



شکل ۱۳. تاثیر اندازه‌ی داده‌ی آموزش بر دقت مدل

پس در کل هر چه اندازه‌ی داده‌های آموزش بیشتر باشد دقت مدل افزایش می‌یابد اما این افزایش دقت در  $n$  های بزرگ آهسته تر رخ می‌دهد. بهترین دقت هم در اینجا به ازای  $n = 45$  رخ داده است.

## مراجع

## References

- [۱] faradars, "naive-bayes-classifier," ۲۳ ۰۱ ۱۴۰۴. [Online]. Available: <https://blog.faradars.org/naive-bayes-classifier/#%d8%af%d8%b3%d8%aa%d9%87-%d8%a8%d9%86%d8%af-%d8%a8%db%8c%d8%b2-%d9%88-%d9%85%d8%af%d9%84-%d8%a7%d8%ad%d8%aa%d9%85%d8%a7%d9%84%d8%a7%d8%aa%db%8c-۱>.
- [۲] "Gamma distribution," Wikipedia, [Online]. Available: [https://en.wikipedia.org/wiki/Gamma\\_distribution](https://en.wikipedia.org/wiki/Gamma_distribution).

جزوه‌ی درس [۳]