



## تمرین چهارم

شماره دانشجویی: ۹۸۱۰۱۰۷۴

محمدجواد هزاره

## سوال ۱

(آ) بردار  $\hat{n}$  را بردار یکه در راستای متناظر با  $w$  ای در نظر می گیریم که  $B$  را بدست می آورد. به عبارتی خواهیم داشت:

$$\forall i \in \{1, \dots, m\} : y_i \langle B\hat{n}, x_i \rangle \geq 1 \quad (*)$$

فرض کنیم الگوریتم  $k$  بار اجرا می شود. هدف پیدا کردن کرانی برای  $k$  است. حال اگر در مرحله  $t$  ام، داده ی  $j$  ام اشتباه دسته بندی شده باشد، آنگاه  $w^{t+1} = w^t + y_j x_j$  بنابرین داریم:

$$\begin{aligned} \langle w^{t+1}, B\hat{n} \rangle &= B \langle w^t + y_j x_j, \hat{n} \rangle \\ &= B (\langle w^t, \hat{n} \rangle + \langle y_j x_j, \hat{n} \rangle) \end{aligned}$$

$$\begin{aligned} \Rightarrow \langle w^{t+1}, \hat{n} \rangle &= \langle w^t, \hat{n} \rangle + \frac{1}{B} \langle y_j x_j, B\hat{n} \rangle \\ &\geq \langle w^t, \hat{n} \rangle + \frac{1}{B} \end{aligned}$$

که خط آخر با توجه به  $(*)$  نتیجه شده است. حال با توجه به استقرا و این که  $w^1 = 0$ ، برای هر  $t$  در بازه ی ۱ تا  $k$  خواهیم داشت: (تمامی نرم ها  $L_2$  هستند).

$$\langle w^{t+1}, \hat{n} \rangle \geq \frac{t}{B} \xrightarrow{\|\hat{n}\|=1} \|w^{t+1}\| \geq \frac{t}{B} \quad (*)$$

هم چنین برای  $\|w^{t+1}\|$  داریم:

$$\begin{aligned} \|w^{t+1}\|^2 &= \|w^t + y_j x_j\|^2 \\ &= \|w^t\|^2 + y_j^2 \|x_j\|^2 + 2 \langle w^t, y_j x_j \rangle \quad (\diamond) \end{aligned}$$

و از آنجایی که فرض کرده بودیم داده‌ی  $j$  ام اشتباه دسته‌بندی شده است،  $\langle w^t, y_j x_j \rangle \leq 0$  خواهد بود و همچنین مطابق تعریف الگوریتم برچسب‌ها یک یا منفی یک بودند که نتیجه می‌دهد  $y_j^2 = 1$ . مطابق تعریف  $R$  برای هر  $i$ ، خواهیم داشت  $\|x_i\| \leq R$ ؛ با توجه به این موارد داریم:

$$(\diamond) \implies \|w^{t+1}\|^2 \leq \|w^t\|^2 + R^2$$

با استفاده از استقرا خواهیم داشت:

$$\|w^{t+1}\|^2 \leq t R^2 \quad (**)$$

با کنار هم گذاشتن  $(*)$  و  $(**)$  و قرار دادن  $k$  به جای  $t$  خواهیم داشت:

$$\frac{k^2}{B^2} \leq \|w^{t+1}\|^2 \leq k R^2 \implies \boxed{k \leq B^2 R^2}$$

بنابراین حداکثر دفعات تکرار مراحل الگوریتم برابر  $B^2 R^2$  خواهد بود.

ب) کفایت مقیاس فضای ویژگی‌ها را در  $\frac{1}{\eta}$  ضرب کنیم. داده‌های مسئله  $(\tilde{x}_i, y_i) = (\eta x_i, y_i)$  خواهند بود و در مرحله‌ی آپدیت کردن وزن‌ها داریم  $w^{t+1} = w^t + \eta y_i x_i = w^t + y_i \tilde{x}_i$  بنابراین کفایت پارامترهای  $B$  و  $R$  را در فضای ویژگی‌های جدید که  $\tilde{x}_i$  است محاسبه کنیم. برای  $R$  در فضای جدید خواهیم داشت:

$$\tilde{R} = \max_i \|\tilde{x}_i\| = \eta \max_i \|x_i\| = \eta R$$

برای محاسبه‌ی  $\tilde{B}$  داریم:

$$y_i \langle B \hat{n}, x_i \rangle \geq 1$$

$$y_i \langle B \hat{n}, \eta x_i \rangle \geq \eta$$

$$\frac{1}{\eta} (y_i \langle B \hat{n}, \eta x_i \rangle) \geq 1$$

$$y_i \langle \frac{B}{\eta} \hat{n}, \eta x_i \rangle \geq 1$$

بنابراین  $\tilde{B} = \frac{B}{\eta}$  و اگر این‌گونه نباشد، یعنی بتوان  $\tilde{B}$  دیگری یافت که در رابطه‌ی مورد نظر صدق کند و از مقدار داده شده کمتر باشد، فرض کم‌ترین بودن  $B$  نقض می‌شود. بنابراین مقدار ارائه شده کم‌ترین مقداری است که می‌توان برای  $\tilde{B}$  پیدا کرد. بنابراین در این حالت مراحل الگوریتم حداکثر  $\tilde{B}^2 \tilde{R}^2 = \frac{B^2}{\eta^2} \eta^2 R^2 = B^2 R^2$  که مشابه همان حالت (آ) است تکرار خواهد شد.

## سوال ۲

ToDo (آ)

(ب) همانطور که گفته شده، نشان می‌دهیم ضرب داخلی دو بردار صفر است. با توجه به روابطی که برای ضرایب داریم خواهیم داشت:

$$\begin{aligned} A &= \sum_{i=1}^m y_i h_t(x_i) D_{t+1}(i) \\ &= \sum_{i=1}^m \frac{y_i h_t(x_i) D_t(i)}{Z_t} e^{-\alpha_t y_i h_t(x_i)} \end{aligned}$$

اگر  $\mathcal{C}$  مجموعه‌ی اندیس داده‌هایی باشد که درست دسته‌بندی شده و  $\mathcal{M}$  مجموعه‌ی اندیس داده‌هایی که غلط دسته‌بندی شده‌اند، آنگاه  $A$  را می‌توان به صورت زیر نوشت:

$$A = \frac{1}{Z_t} \left( \sum_{i \in \mathcal{C}} D_t(i) e^{-\alpha_t} + \sum_{j \in \mathcal{M}} -D_t(j) e^{\alpha_t} \right)$$

که با توجه به این که  $\alpha_t = \ln(\sqrt{\frac{1-\epsilon_t}{\epsilon_t}})$  ادامه می‌دهیم:

$$\begin{aligned} A &= \frac{1}{Z_t} \left( \sum_{i \in \mathcal{C}} D_t(i) \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} - \sum_{j \in \mathcal{M}} D_t(j) \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \right) \\ &= \frac{1}{Z_t \sqrt{\epsilon_t(1-\epsilon_t)}} \left( \sum_{i \in \mathcal{C}} \epsilon_t D_t(i) + \sum_{j \in \mathcal{M}} \epsilon_t D_t(j) - \sum_{j \in \mathcal{M}} D_t(j) \right) \\ &= \frac{1}{Z_t \sqrt{\epsilon_t(1-\epsilon_t)}} \left( \epsilon_t \sum_{i=1}^m D_t(i) - \sum_{j \in \mathcal{M}} D_t(j) \right) \end{aligned}$$

با توجه به این که ضرایب نرمالایز شده‌اند داریم  $\sum_{i=1}^m D_t(i) = 1$ ، همچنین برای خطا داریم:

$$\epsilon_t = \sum_{i=1}^m D_t(i) \mathbf{I}(h(x_i) \neq y_i) = \sum_{i \in \mathcal{M}} D_t(i)$$

بنابراین:

$$A = \frac{1}{Z_t \sqrt{\epsilon_t(1-\epsilon_t)}} \left( \epsilon_t - \sum_{j \in \mathcal{M}} D_t(j) \right) = 0$$

بنابراین بردار ضرایب  $D_{t+1}$  بر بردار شامل  $y_i h_t(x_i)$  عمود بوده یا به عبارتی uncorrelated هستند.

### سوال ۳

(آ) نرخ خطا را به صورت زیر بازنویسی می‌کنیم:

$$R(h) = \mathbb{P}\{Y \neq h(X)\} = 1 - \mathbb{P}\{Y = h(X)\}$$

بنابراین برای کمینه کردن خطا کافیت احتمال برابر شدن  $Y$  با  $h(X)$  را بیشینه کنیم. این عبارت را نیز به صورت زیر می‌توان باز کرد که  $D = \{x \in \mathbb{R}^n \mid h(x) = 0\}$  و  $D' = \mathbb{R}^n - D$  و  $n$  بعد بردار ویژگی‌ها است.

$$\mathbb{P}\{Y = h(x)\} = \mathbb{P}\{Y = 1, h(X) = 1\} + \mathbb{P}\{Y = 0, h(X) = 0\}$$

$$= \mathbb{P}\{Y = 1, X \in D'\} + \mathbb{P}\{Y = 0, X \in D\}$$

$$= \int_{D'} \mathbb{P}\{Y = 1 \mid X = x\} p_X(x) dx + \int_D \mathbb{P}\{Y = 0 \mid X = x\} p_X(x) dx \quad (*)$$

که اگر تعریف کنیم  $m(x) = \mathbb{P}\{Y = 1 \mid X = x\}$ ، آن‌گاه خواهیم داشت  $\mathbb{P}\{Y = 0 \mid X = x\} = 1 - m(x)$ . با جایگذاری این تعاریف در رابطه‌ی بالا ادامه می‌دهیم:

$$\begin{aligned} \mathbb{P}\{Y = h(x)\} &\stackrel{(*)}{=} \int_{D'} m(x) p_X(x) dx + \int_D (1 - m(x)) p_X(x) dx \\ &= \int_{D'} m(x) p_X(x) dx + \int_D (1 - m(x)) p_X(x) dx + \left( \int_D m(x) p_X(x) dx - \int_D m(x) p_X(x) dx \right) \\ &= \int_{\mathbb{R}^n} m(x) p_X(x) dx + \int_D (1 - 2m(x)) p_X(x) dx \\ &= C + \int_D (1 - 2m(x)) p_X(x) dx \end{aligned}$$

بنابراین فقط عبارت دوم در رابطه‌ی بالا به  $h$  بستگی دارد که اگر بخواهیم احتمال مورد نظر را بیشینه کنیم، باید  $X$  هایی عضو  $D$  باشند که عبارت زیر انتگرال برای آن‌ها نامنفی شود. بنابراین برای  $h^*$  داریم:

$$\left. \begin{array}{l} \forall x \in D : h^*(x) = 0 \\ \forall x \in D : 1 - 2m(x) \geq 0 \end{array} \right\} \implies h^*(x) = 0 \iff m(x) \leq \frac{1}{2}$$

بنابراین تابع  $h^*$  به صورت زیر خواهد بود:

$$h^*(x) = \begin{cases} 1 & m(x) > \frac{1}{2} \\ 0 & m(x) \leq \frac{1}{2} \end{cases}$$

ب) با توجه به قسمت (آ) برای خطای توابع داده شده داریم:

$$\begin{cases} R(h^*) = 1 - C - \int_D (1 - 2m(x))p_X(x) dx & D = \{x \mid m(x) \leq \frac{1}{2}\} \\ R(\hat{h}) = 1 - C - \int_B (1 - 2m(x))p_X(x) dx & B = \{x \mid \hat{m}(x) \leq \frac{1}{2}\} \end{cases}$$

بنابراین:

$$E = R(\hat{h}) - R(h^*) = \int_D (1 - 2m(x))p_X(x) dx - \int_B (1 - 2m(x))p_X(x) dx$$

با اضافه و کم کردن  $2\hat{m}(x)$  به انتگرال‌ده‌ها ادامه می‌دهیم:

$$\begin{aligned} E &= \int_D (1 - 2m(x) + 2\hat{m}(x) - 2\hat{m}(x))p_X(x) dx - \int_B (1 - 2m(x) + 2\hat{m}(x) - 2\hat{m}(x))p_X(x) dx \\ &= \left( \int_D (1 - 2\hat{m}(x))p_X(x) dx - \int_B (1 - 2\hat{m}(x))p_X(x) dx \right) + 2 \int_D (\hat{m}(x) - m(x))p_X(x) dx \\ &\quad - 2 \int_B (\hat{m}(x) - m(x))p_X(x) dx \end{aligned}$$

حاصل داخل پرانتز در عبارت بالا منفی خواهد بود؛ چرا که  $B$  بهترین مجموعه‌ای است که می‌توان با توجه به  $\hat{m}$  انتخاب کرد. اگر در استدلال‌های قسمت (آ) تابع  $m$  را با  $\hat{m}$  جایگذاری کنیم، به این نتیجه می‌رسیم که  $B$  مجموعه‌ای است که حاصل انتگرال داخل پرانتز روی  $B$  را بیشینه می‌کند. بنابراین همین انتگرال روی هر مجموعه‌ی دیگری مانند  $D$  مقداری کمتر یا مساوی حاصل انتگرال روی  $B$  خواهد داشت. بنابراین:

$$E \leq 2 \int_D (\hat{m}(x) - m(x))p_X(x) dx - 2 \int_B (\hat{m}(x) - m(x))p_X(x) dx \quad (\star)$$

برای ادامه فرض کنیم  $F = \{x \mid x \in D \wedge x \notin B\}$  و  $G = \{x \mid x \in B \wedge x \notin D\}$ . انتگرال‌های بالا به ازای  $x$ ‌های مشترک در  $B$  و  $D$  یکدیگر را خنثی می‌کنند و فقط  $x$ ‌هایی که در  $F$  و  $G$  هستند باقی خواهند ماند. از

طرفی داریم:

$$\begin{cases} \forall x \in F : m(x) \leq \frac{1}{2}, \hat{m}(x) > \frac{1}{2} & \implies |\hat{m}(x) - m(x)| = \hat{m}(x) - m(x) \\ \forall x \in G : m(x) > \frac{1}{2}, \hat{m}(x) \leq \frac{1}{2} & \implies |\hat{m}(x) - m(x)| = -(\hat{m}(x) - m(x)) \end{cases}$$

بنابراین در ادامه‌ی (\*) خواهیم داشت:

$$\begin{aligned} E &\leq 2 \int_F (\hat{m}(x) - m(x)) p_X(x) dx - 2 \int_G (\hat{m}(x) - m(x)) p_X(x) dx \\ &\leq 2 \int_F |\hat{m}(x) - m(x)| p_X(x) dx + 2 \int_G |\hat{m}(x) - m(x)| p_X(x) dx \\ &\leq 2 \int_{F \cup G} |\hat{m}(x) - m(x)| p_X(x) dx \\ &\leq 2 \int_{R^n} |\hat{m}(x) - m(x)| p_X(x) dx \end{aligned}$$

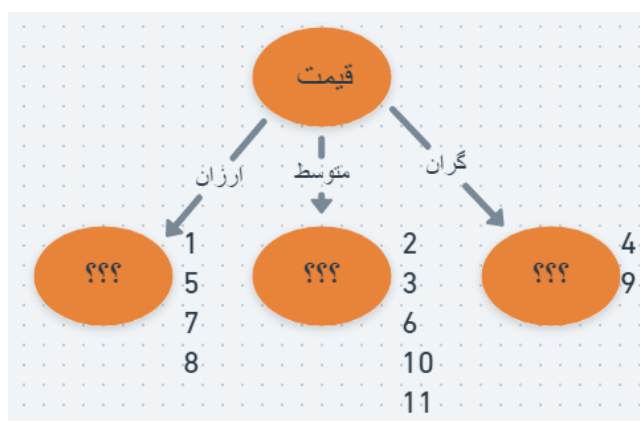
که در مرحله‌ی آخر از مثبت بودن انتگرال‌ده استفاده شده و مقادیری مثبت به حاصل اضافه شده است که تغییری در جهت نابرابری ایجاد نمی‌کند.

## سوال ۴

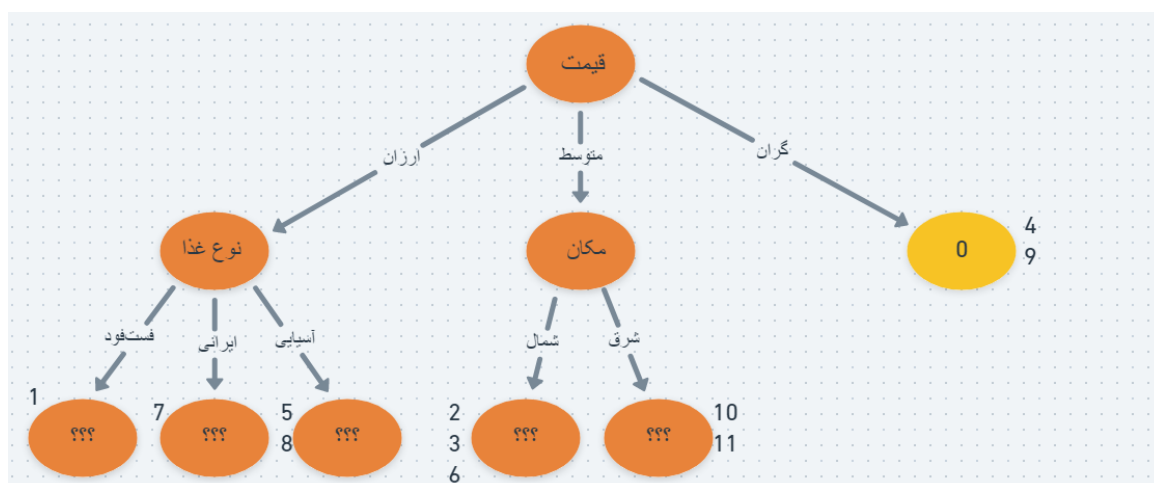
آ) برای ساخت درخت از هیوریستیک Information Gain استفاده می‌کنیم؛ به این صورت که در هر راس، ویژگی‌ای را انتخاب می‌کنیم که بین داده‌هایی که به آن راس رسیده‌اند این پارامتر برای آن بیشینه باشد. در مرحله‌ی اول و برای راس داریم:<sup>۱</sup>

$$IG(\text{محدودیت}) \approx 0.016, \quad IG(\text{محل}) \approx 0.016, \quad IG(\text{قیمت}) \approx 0.189, \quad IG(\text{نوع غذا}) \approx 0.052$$

بنابراین در ریشه ویژگی قیمت انتخاب خواهد شد و درخت به صورت زیر تقسیم‌بندی می‌شود:



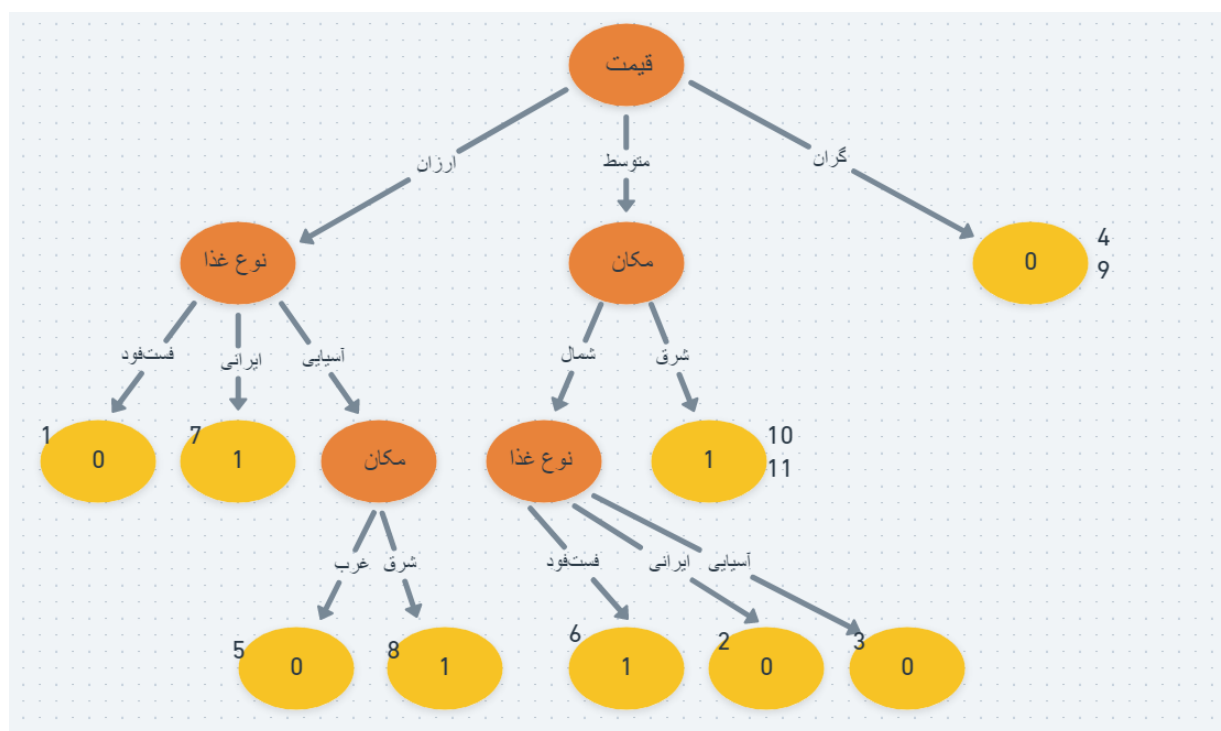
در مرحله‌ی بعد داده‌هایی که به سمت گران دسته‌بندی می‌شوند همه یک برجسب خورده‌اند بنابراین این راس به برگ تبدیل شده و مقدار برجسب 0 می‌گیرد. برای راس ارزان، داده‌ها IG برابری دارند بنابراین یکی از آن‌ها را به صورت تصادفی انتخاب کرده‌ایم که ویژگی نوع غذا بوده است. برای راس متوسط نیز IG ویژگی مکان بیش‌تر بوده و این ویژگی انتخاب شده است. درخت بعد از این مرحله به شکل زیر در می‌آید:



در این مرحله راس‌هایی که یک داده برای آن‌ها مانده برجست همان داده را می‌گیرند. داده‌های ۱۰ و ۱۱ نیز

<sup>۱</sup> برای محاسبه‌ی Information Gain از ماشین حساب آنلاین استفاده شده است.

برچسب یکسانی دارند بنابراین راس آن‌ها نیز همان برچسب آن‌ها یعنی 1 را خواهد گرفت. برای مابقی راس‌ها از میان ویژگی‌های باقی مانده، آن ویژگی که IG بیش‌تری داشته باشد را انتخاب می‌کنیم. برای داده‌های ۵ و ۸، ویژگی‌های باقی مانده «مکان» و «محدودیت» هستند که هر دو نیز IG یکسانی دارند بنابراین به صورت تصادفی مکان انتخاب شده است. برای داده‌های ۲، ۳ و ۶، نوع غذا بیش‌ترین IG را دارد. پس از این تقسیم‌بندی راس‌های باقی مانده یک داده خواهند داشت و در نتیجه به برگ تبدیل می‌شوند. درخت تصمیم در نهایت به صورت زیر در خواهد آمد:



ب) با استفاده از درختی که در قسمت (آ) بدست آوردیم داریم:

$$\left\{ \begin{array}{l} \text{رضایت مندی}(12) = 0 \\ \text{رضایت مندی}(13) = 1 \\ \text{رضایت مندی}(14) = 1 \\ \text{رضایت مندی}(15) = 1 \\ \text{رضایت مندی}(16) = 1 \end{array} \right.$$



ج) برای بدست آوردن  $F_1$  Score نیاز به  $recall$  و  $precision$  داریم. برای این پارامترها نیز خواهیم داشت:

$$\begin{cases} recall = \frac{TP}{TP + FN} = \frac{2}{2 + 0} = 1 \\ precision = \frac{TP}{TP + FP} = \frac{2}{2 + 2} = 0.5 \end{cases}$$

بنابراین برای  $F_1$  Score داریم:

$$F_1 \text{ Score} = \frac{2 \times precision \times recall}{precision + recall} = \frac{2 \times 0.5 \times 1}{0.5 + 1} = \frac{2}{3}$$

## سوال ۵

(آ) با توجه به استقلال داده‌ها از یکدیگر، درست‌نمایی را می‌توان به صورت زیر نوشت:

$$\mathbb{P}\{\mathcal{D} \mid \pi_1, \dots, \pi_K\} = \prod_{i=1}^N \mathbb{P}\{\mathcal{D}_i \mid \pi_1, \dots, \pi_K\} \quad (*)$$

که با توجه به این که  $\mathcal{D} = \{\phi_n, t_n\}_{n=1}^N$ ، برای احتمال دیدن یک داده خواهیم داشت:

$$\mathbb{P}\{t_i \mid \pi_1, \dots, \pi_K\} = \pi_j, \quad (t_i)_j = 1 \quad (*)$$

بنابراین احتمال دیدن هر داده برابر با خواهد بود با  $\pi_k$  که  $k$  کلاسی است که به آن تعلق دارد. حال اگر  $N_i$  تعداد داده‌هایی باشد که به کلاس  $i$ ام تعلق دارند، با توجه به (\*) و (\*) برای احتمال دیدن  $\mathcal{D}$  خواهیم داشت:

$$\mathbb{P}\{\mathcal{D} \mid \pi_1, \dots, \pi_K\} = \prod_{i=1}^K \pi_i^{N_i}$$

برای پیدا کردن تخمین گر  $\pi_i$  از لگاریتم تابع درست‌نمایی استفاده می‌کنیم. بنابراین داریم:

$$\ln(\mathbb{P}\{\mathcal{D} \mid \pi_1, \dots, \pi_K\}) = \sum_{j=1}^K N_j \ln(\pi_j)$$

که  $N_j$  تعداد داده‌هایی است که در کلاس  $j$ ام دسته‌بندی شده‌اند. برای پیدا کردن  $\pi_i$ هایی که عبارت بالا را بیشینه می‌کنند، شرط دیگری نیز داریم و آن هم یک شدن جمع همه‌ی آن‌ها یا به عبارتی  $\sum_{j=1}^K \pi_j = 1$  است. برای پیدا کردن بیشینه از روش ضرایب لاگرانژ استفاده می‌کنیم. تابع لاگرانژ به صورت زیر خواهد بود:

$$L(\pi_1, \dots, \pi_K, \lambda) = \sum_{j=1}^K N_j \ln(\pi_j) - \lambda \left( \sum_{j=1}^K \pi_j - 1 \right)$$

بنابراین برای مشتقات داریم:

$$\begin{cases} \frac{\partial L}{\partial \pi_j} = \frac{N_j}{\pi_j} - \lambda \\ \frac{\partial L}{\partial \lambda} = - \sum_{j=1}^K \pi_j + 1 \end{cases}$$

معادلات بالا را برابر صفر قرار داده و با بازنویسی رابطه‌ی اول خواهیم داشت:

$$N_j = \lambda \pi_j^* \xrightarrow{\sum_{j=0}^K} \lambda = N$$

$$\implies \boxed{\pi_j^* = \frac{N_j}{N}}$$

(ب) فرض کنیم  $\mathcal{C}_i$  مجموعه‌ی شامل اندیس داده‌هایی باشد که کلاس مربوط به آن‌ها  $i$  است. همچنین  $\mu$  بردار شامل  $\mu_k$ ها بوده و  $\pi$  بردار شامل  $\pi_k$ ها باشد. با توجه به استقلال داده‌ها برای درست‌نمایی خواهیم داشت:

$$\mathbb{P}\{\mathcal{D} \mid \mu, \Sigma, \pi, \Phi\} = \prod_{i=1}^N \mathbb{P}\{t_i \mid \mu, \Sigma, \pi, \phi_i\}$$

که با توجه به تعریف  $\mathcal{C}_i$ ها خواهیم داشت:

$$\mathbb{P}\{\mathcal{D} \mid \mu, \Sigma, \pi, \Phi\} = \prod_{j=1}^K \left( \prod_{i \in \mathcal{C}_j} \mathbb{P}\{\mathcal{C}_j \mid \mu, \Sigma, \pi, \phi_i\} \right)$$

که احتمال داخل پرانتز بالا را نیز می‌توان به صورت زیر نوشت که  $d$  بعد ویژگی‌هاست:

$$\begin{aligned} \mathbb{P}\{\mathcal{C}_j \mid \mu, \Sigma, \pi, \phi_i\} &= \mathbb{P}(\mathcal{C}_j) \mathbb{P}\{\phi_i \mid \mathcal{C}_j, \mu, \Sigma\} \\ &= \pi_j (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\phi_i - \mu_j)^T \Sigma^{-1} (\phi_i - \mu_j) \right) \end{aligned}$$

از بیشینه کردن لگاریتم درست‌نمایی استفاده می‌کنیم:

$$\begin{aligned} f(\mu, \Sigma) &= \ln(\mathbb{P}\{\mathcal{D} \mid \mu, \Sigma, \pi, \Phi\}) \\ &= \sum_{j=1}^K N_j \ln(\pi_j) - \frac{Nd}{2} \ln(2\pi) - \frac{N}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{j=1}^K \left( \sum_{i \in \mathcal{C}_j} (\phi_i - \mu_j)^T \Sigma^{-1} (\phi_i - \mu_j) \right) \end{aligned}$$

برای مشتق تابع بالا نسبت به  $\mu$  و  $\Sigma$  خواهیم داشت:

$$\begin{cases} \frac{\partial f}{\partial \mu_j} = -\frac{1}{2} \sum_{i \in \mathcal{C}_j} 2 \Sigma^{-1} (\phi_i - \mu_j) \\ \frac{\partial f}{\partial \Sigma} = -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} \sum_{j=1}^K \left( \sum_{i \in \mathcal{C}_j} (\phi_i - \mu_j) (\phi_i - \mu_j)^T \Sigma^{-1} \Sigma^{-1} \right) \end{cases}$$

با صفر قرار دادن معادلات بالا برای  $\mu_j^*$  و  $\Sigma^*$  خواهیم داشت:

$$N_j \mu_j^* = \sum_{i \in \mathcal{C}_j} \phi_i \implies \boxed{\mu_j^* = \frac{\sum_{i \in \mathcal{C}_j} \phi_i}{N_j}}$$

$$\boxed{\Sigma^* = \frac{\sum_{j=1}^K \left( \sum_{i \in \mathcal{C}_j} (\phi_i - \mu_j)(\phi_i - \mu_j)^T \right)}{N}}$$