

$$\max_{\mathbf{w}} \text{var}(\mathbf{z})$$

$$\|\mathbf{w}\|_2^2 = 1$$

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

$$\mathcal{J}(\mathbf{w}) = \frac{(\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)^T \mathbf{S}^{-1} (\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)}{\tilde{\mathbf{S}}_1^2 + \tilde{\mathbf{S}}_2^2}$$

$$\tilde{\mathbf{m}}_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$$

$$\tilde{\mathbf{S}}_1^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T$$

$$\mathbf{w}^T (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T) \mathbf{w} = \mathbf{w}^T \mathbf{S} \mathbf{w}$$

$$\frac{(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2}{\mathbf{w}^T \mathbf{S}_1^2 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2^2 \mathbf{w}}$$

$$\mathbf{w}^T (\mathbf{S}_1^2 + \mathbf{S}_2^2) \mathbf{w}$$

$$\max_{\mathbf{w}} (\mathbf{w}^T \mathbf{S}_B \mathbf{w})$$

$$\mathbf{S}_B^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

$$\max_{\lambda} \lambda \mathbf{w}^T \mathbf{S}_B \mathbf{w} = \max_{\lambda} \lambda k$$

$$\mathcal{Q}(\theta) = E[\mathcal{L}(\theta)] = \sum_{i=1}^n \log \pi(1 - \delta_i^t) + \log N(\mathbf{x}_i | \mu_0, \Sigma_0, \delta_i^t) + \delta_i^t \log \pi + \delta_i^t \log N(\mathbf{x}_i | \mu_1, \Sigma_1, \delta_i^t)$$

$$\frac{d\mathcal{Q}}{d\alpha} = 0 \rightarrow \sum_{i=1}^n \frac{-1}{1-\alpha} (1-\delta_i^t) + \frac{\delta_i^t}{\alpha} = 0$$

$$\frac{n - \sum \delta_i^t}{1-\alpha} = \frac{\sum \delta_i^t}{\alpha} \rightarrow \alpha = \frac{\sum \delta_i^t}{n}$$

$$\mathcal{Q}(\theta) = E_{\mathbf{z}_i} [\mathcal{L}(\theta)] = \sum_{i=1}^n \sum_{k=1}^K E[\mathbf{z}_{ik}] (\log \alpha_k + \log N(\mathbf{x}_i | \mu_k, \Sigma_k))$$

Mercer's Theorem:

$$\iint g(\mathbf{x}) k(\mathbf{x}, \mathbf{y}) g(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0 \rightarrow \text{Valid kernel}$$

$$F_T(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}) = F_{T-1}(\mathbf{x}) + \alpha_T h_T(\mathbf{x})$$

$$P(\mathbf{z}_{ik} = 1 | \mathbf{x}_i, \theta^t) = P(\mathbf{z}_{ik} = 1 | \mathbf{x}_i, \theta^t)$$

$$\sum_{i=1}^n \omega_i^{(t-1)} e^{-\mathbf{y}_i^T \mathbf{h}_T(\mathbf{x}_i)} = \sum_{i: \mathbf{y}_i \neq \mathbf{h}_T(\mathbf{x}_i)} \omega_i^{(t-1)} e^{-\mathbf{y}_i^T \mathbf{h}_T(\mathbf{x}_i)} + \sum_{i: \mathbf{y}_i = \mathbf{h}_T(\mathbf{x}_i)} \omega_i^{(t-1)} e^{-\mathbf{y}_i^T \mathbf{h}_T(\mathbf{x}_i)}$$

$$\frac{dE}{d\alpha_T} = e^{\alpha_T} \sum_{i \in M} \omega_i^{(T-1)} - e^{-\alpha_T} \sum_{i \in C} \omega_i^{(T-1)} = 0$$

$$\alpha_T = \frac{1}{2} \ln \frac{\sum_{i \in C} \omega_i^{(T-1)}}{\sum_{i \in M} \omega_i^{(T-1)}} \quad \alpha_T = \frac{1}{2} \ln \frac{1 - \epsilon_T}{\epsilon_T}$$

$$P(\mathbf{x} | \theta) = (1-\alpha) N(\mathbf{x} | \mu_0, \Sigma_0) + \alpha N(\mathbf{x} | \mu_1, \Sigma_1)$$

$$P(\mathbf{x}, \mathbf{z} | \theta) = [(1-\alpha) N(\mathbf{x} | \mu_0, \Sigma_0)]^{1-\mathbf{z}} [\alpha N(\mathbf{x} | \mu_1, \Sigma_1)]^{\mathbf{z}}$$

$$E_{\mathbf{z}} [\mathcal{L}(\theta)] = E[\sum_{i=1}^n \log P(\mathbf{x}_i, \mathbf{z}_i | \theta)] = \sum_{i=1}^n E_{\mathbf{z}} [\log P(\mathbf{x}_i, \mathbf{z}_i | \theta)]$$

$$= \sum_{i=1}^n E_{\mathbf{z}} [(1-\mathbf{z}_i) \log((1-\alpha) N(\mathbf{x}_i | \mu_0, \Sigma_0)) + \mathbf{z}_i \log(\alpha N(\mathbf{x}_i | \mu_1, \Sigma_1))] + \sum_{i=1}^n \mathbf{z}_i \log \alpha$$

$$= \sum_{i=1}^n \log((1-\alpha) N(\mathbf{x}_i | \mu_0, \Sigma_0)) + \log N(\mathbf{x}_i | \mu_1, \Sigma_1) (1 - E[\mathbf{z}_i]) + E[\mathbf{z}_i] \log \alpha$$

$$+ E[\mathbf{z}_i] \log N(\mathbf{x}_i | \mu_1, \Sigma_1)$$

$$E[\mathbf{z}_i] = E[\mathbf{z}_i | \mathbf{x}] = P(\mathbf{z}_i = 1 | \mathbf{x}) = \frac{P(\mathbf{x} | \mathbf{z}_i = 1) P(\mathbf{z}_i = 1)}{P(\mathbf{x})} = \frac{N(\mathbf{x} | \mu_1, \Sigma_1) \alpha^t}{(1-\alpha) N(\mathbf{x} | \mu_0, \Sigma_0) + \alpha N(\mathbf{x} | \mu_1, \Sigma_1)}$$

$$\mathcal{L}(\theta) = \log P(\mathbf{D}, \mathbf{H} | \theta) = \sum_{i=1}^n \log P(\mathbf{x}_i, \mathbf{z}_i | \theta)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} (\log \alpha_k + \log N(\mathbf{x}_i | \mu_k, \Sigma_k))$$

$$\frac{\partial \mathcal{Q}(\theta)}{\partial \alpha_k} = 0$$

$$\hat{\alpha}_k = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik}}{n}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \mathbf{x}_i}{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik}}$$

$$P(\mathbf{x} | \theta) = \sum_{k=1}^K \alpha_k N(\mathbf{x} | \mu_k, \Sigma_k)$$

$$\sum_{k=1}^K \alpha_k = 1 \quad \alpha_k \geq 0 \quad P(\mathbf{x}, \mathbf{z} | \theta) = \prod_{k=1}^K (\alpha_k N(\mathbf{x}_i | \mu_k, \Sigma_k))^{\mathbf{z}_{ik}}$$

$$E[\mathbf{z}_{ik}] = P(\mathbf{z}_{ik} = 1 | \mathbf{x}_i, \theta^t) = \frac{P(\mathbf{x}_i | \mu_k, \Sigma_k) \alpha_k^t}{\sum_{k=1}^K P(\mathbf{x}_i | \mu_k, \Sigma_k) \alpha_k^t}$$

$$\frac{\partial \mathcal{Q}(\theta)}{\partial \alpha_j} = \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j} - n = 0 \Rightarrow \alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{n}$$

$$\alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}$$

$$\mathcal{Q}(\theta) = \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} (\log \alpha_k + \log N(\mathbf{x}_i | \mu_k, \Sigma_k)) = n (\sum_{k=1}^K \alpha_k - 1)$$

سوال ۷ خوشه‌بندی (۱۵ نمره)

می‌خواهیم مجموعه داده‌های \mathcal{X} را به k خوشه تقسیم کنیم. همانطور که می‌دانید، الگوریتم k -means با انتخاب k مرکز دسته به صورت تصادفی آغاز می‌شود. فرض کنید که راهبر زیر را برای انتخاب k مرکز دسته اولیه داشته باشیم:

مرکز دسته اول را تصادفی انتخاب کرده و مرکز دسته‌های بعدی را به گونه‌ای انتخاب می‌کنیم که از مرکز دسته‌هایی که قبلاً انتخاب شده‌اند بیشترین فاصله را داشته باشد.

به کمک یک مثال نشان دهید که این راهبر الزاماً به دسته‌بندی پهنه منجر نمی‌شود. به عبارت دیگر نشان دهید که ممکن است هزینه راهبر پیشنهادی از هزینه الگوریتم پهنه خوشه‌بندی k -means که در عبارت زیر نمایش داده شده، بیشتر است:

$$\phi = \sum_{\mathbf{x} \in \mathcal{X}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|^2$$

$$\langle \vec{w}, \mathbf{x}_1 - \mathbf{x}_2 \rangle = 0 \rightarrow \vec{w}^T (\mathbf{x}_2 - \mathbf{x}_1) = \vec{w}^T \mathbf{x}_2 - \vec{w}^T \mathbf{x}_1$$

$$= \vec{w}^T \mathbf{x}_2 + b - \vec{w}^T \mathbf{x}_1 - b = 0$$

$$\max_{\mathbf{w}, b} \text{margin}(\mathbf{w}, b) = \frac{1}{\|\mathbf{w}\|_2} |\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2)| = \frac{1}{\|\mathbf{w}\|_2} |\mathbf{w}^T \mathbf{x}_1 + b - \mathbf{w}^T \mathbf{x}_2 - b| = \frac{1}{\|\mathbf{w}\|_2} |\mathbf{w}^T \mathbf{x}_1 + b - \mathbf{w}^T \mathbf{x}_2 - b|$$

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|_2} \quad \min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2$$

$$s.t. \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad s.t. \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$\mathcal{L}(\mathbf{w}, b, \lambda) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \lambda \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b + 1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} + \sum_{i=1}^n \lambda \mathbf{y}_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = - \sum_{i=1}^n \lambda \mathbf{y}_i \mathbf{x}_i$$

$$- \sum_{i=1}^n \lambda \mathbf{y}_i b = 0 \Rightarrow \sum_{i=1}^n \lambda \mathbf{y}_i = 0 \Rightarrow \sum_{i=1}^n \lambda \mathbf{y}_i = 0$$

$$\mathcal{J}(\lambda) = \frac{1}{2} (\sum_{i=1}^n \lambda \mathbf{y}_i \mathbf{x}_i)^T (\sum_{i=1}^n \lambda \mathbf{y}_i \mathbf{x}_i) - \sum_{i=1}^n \lambda + \sum_{i=1}^n \lambda \mathbf{y}_i (\sum_{j=1}^n \lambda \mathbf{y}_j \mathbf{x}_j^T \mathbf{x}_i + b)$$

$$\mathcal{J}(\lambda) = \frac{1}{2} (\sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j) - \sum_{i=1}^n \lambda + \sum_{i=1}^n \lambda \mathbf{y}_i (\sum_{j=1}^n \lambda_j \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j + b)$$

$$\mathbf{w}^* = \sum_{i=1}^n \lambda_i^* \mathbf{y}_i^* \mathbf{x}_i^* \quad \lambda_k^* = \frac{1}{\sum_{j=1}^n \lambda_j^*} \mathbf{y}_k (\mathbf{w}^T \mathbf{x}_k + b) = 0$$

$$\min_{\mathbf{w}, b, \lambda} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \lambda_i \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \lambda_i \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b)$$

$$s.t. \mathbf{y}_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i \quad \epsilon_i \geq 0$$

$$\mathcal{L}_{\epsilon_j} = C + (-\lambda_j) - \mu_j = 0 \Rightarrow \mu_j = C - \lambda_j \geq 0 \Rightarrow 0 \leq \lambda_j \leq C$$

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^n \lambda_i$$

$$\sum_{i=1}^n \lambda_i \mathbf{y}_i = 0, \quad 0 \leq \lambda_i \leq C$$

$$P(\mathbf{x}_i | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp(-\frac{1}{2} (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k))$$

$$\frac{\partial \mathcal{Q}(\theta)}{\partial \alpha_j} = \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j} - n = 0 \Rightarrow \alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{n}$$

$$\alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}$$

$$\frac{\partial \mathcal{Q}(\theta)}{\partial \alpha_j} = \sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j} - n = 0 \Rightarrow \alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{n}$$

$$\alpha_j = \frac{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}{\sum_{i=1}^n \sum_{k=1}^K \mathbf{z}_{ik} \frac{1}{\alpha_j^t}}$$

فرض کنید $k=2$ است و داده‌ها در فضای درجه‌ی یک باشند. حال فرض کنید m نمونه در نقطه‌ی (۰، -۱) و m نمونه در نقطه‌ی (۱، ۰) قرار دارند و یک نقطه \mathbf{x} در مکان (۰، ۰) قرار دارد. انتخاب بهینه مرکز در خوشه این است که نقطه‌ی (۱، ۰) و (۰، -۱) به عنوان دو مرکز انتخاب شوند. اما الگوریتم پهنه‌بندی (گرسن) عملاً نقطه‌ی (۰، ۰) را به عنوان یکی از مرکز انتخاب می‌کند که هزینه آن با بزرگ شدن m به سمت بی‌نهایت می‌رود.

بنابراین می‌توان تابع توزیع آمیخته را محاسبه کرد:

$$p(x|\theta) = (1 - \alpha)P(x; \lambda_0) + \alpha P(x; \lambda_1)$$

چنانچه تابع احتمال کامل را بنویسیم خواهیم داشت:

$$p(x, z|\theta) = [(1 - \alpha)P(x; \lambda_0)]^{1-z} [\alpha P(x; \lambda_1)]^z$$

مرحله Expectation

$$\begin{aligned} E_z[L(\theta)] &= E_z \left[\sum_{i=1}^n \log p(x_i, z_i|\theta) \right] = \sum_{i=1}^n E_z[\log p(x_i, z_i|\theta)] \\ &= \sum_{i=1}^n E_z[(1 - z_i) \log(1 - \alpha) + (1 - z_i) \log P(x_i; \lambda_0) + z \log \alpha + z \log P(x_i; \lambda_1)] \\ &= \sum_{i=1}^n (1 - E[z_i]) \log(1 - \alpha) + (1 - E[z_i]) \log P(x_i; \lambda_0) + E[z_i] \log \alpha \\ &\quad + E[z_i] \log P(x_i; \lambda_1) \end{aligned}$$

از طرفی داریم:

$$\begin{aligned} E[z] &= E[z|x] = p(z = 1|x) = \frac{p(x|z = 1)p(z = 1)}{p(x)} \\ &= \frac{P(x; \lambda_1) \alpha^{[t]}}{(1 - \alpha^{[t]}) P(x; \lambda_0^{[t]}) + \alpha^{[t]} P(x; \lambda_1^{[t]})} = \delta^{[t]} \end{aligned}$$

بدین ترتیب بدست می‌آید:

$$\begin{aligned} Q(\theta) &= E_z[L(\theta)] \\ Q(\theta) &= \sum_{i=1}^n \left((1 - \delta_i^{[t]}) \log(1 - \alpha) + (1 - \delta_i^{[t]}) \log P(x_i; \lambda_0) + \delta_i^{[t]} \log \alpha \right. \\ &\quad \left. + \delta_i^{[t]} \log P(x_i; \lambda_1) \right) \end{aligned}$$

$$\frac{dQ}{d\alpha} = 0$$

$$\begin{aligned} \sum_{i=1}^n -\frac{1}{1 - \alpha} (1 - \delta_i^{[t]}) + \frac{1}{\alpha} \delta_i^{[t]} \\ \frac{1 - \alpha}{n - \sum \delta_i^{[t]}} = \frac{\sum \delta_i^{[t]}}{\alpha} \\ \alpha = \frac{\sum \delta_i^{[t]}}{n} \end{aligned}$$

$$\frac{dQ}{d\lambda_0} = 0$$

$$\begin{aligned} \frac{dQ}{d\lambda_0} &= \sum_{i=1}^n (1 - \delta_i^{[t]}) \frac{d}{d\lambda_0} \log P(x_i; \lambda_0) = 0 \\ \frac{dQ}{d\lambda_0} &= \sum_{i=1}^n (1 - \delta_i^{[t]}) \left(\frac{x_i}{\lambda_0} - 1 \right) = 0 \\ \lambda_0 &= \frac{\sum x_i - \sum \delta_i^{[t]} x_i}{1 - \frac{\sum \delta_i^{[t]}}{n}} \\ \lambda_0 &= \frac{\text{mean} - \frac{\sum \delta_i^{[t]} x_i}{n}}{1 - \alpha^{[t]}} \end{aligned}$$

$$\frac{dQ}{d\lambda_1} = 0$$

$$\begin{aligned} \frac{dQ}{d\lambda_1} &= \sum_{i=1}^n \delta_i^{[t]} \frac{d}{d\lambda_1} \log P(x_i; \lambda_1) = 0 \\ \frac{dQ}{d\lambda_1} &= \sum_{i=1}^n \delta_i^{[t]} \left(\frac{x_i}{\lambda_1} - 1 \right) = 0 \\ \lambda_1 &= \frac{\sum \delta_i^{[t]} x_i}{\sum \delta_i^{[t]}} \\ \lambda_1 &= \frac{\sum \delta_i^{[t]} x_i}{\alpha^{[t]}} \end{aligned}$$

سوال ۴ یادگیری جمعی (Ensemble Learning) (۱۰+۱۵) (نمره)

الف) (نمره ۱۰) برای یک مسأله رگرسیون، دو مدل $y = h_1(x)$ و $y = h_2(x)$ را با استفاده از کمینه کردن مجموع مجذور خطا آموزش داده‌ایم. حال می‌خواهیم با ترکیب این دو مدل، یک مدل قوی‌تر بسازیم. فرض کنید مجموعه داده‌ی $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ را در اختیار داریم، اگر مدل جدید را به صورت $y = \alpha h_1(x) + \beta h_2(x)$ در نظر بگیریم، مقدار بهینه‌ی ضرایب α و β را بدست آورید. توجه داشته باشید که $\alpha + \beta = 1$ لزوماً نیست. بدست آوردن دستگاه دو معادله و دو مجهول برای α و β کفایت می‌کند و نیازی به حل دستگاه نیست.

ب) (اختیاری ۱۵) اگر خطای دو مدل $h_1(x)$ و $h_2(x)$ به ترتیب e_1 و e_2 باشد و از رابطه‌ی $h(x) = \frac{1}{2} h_1(x) + \frac{1}{2} h_2(x)$ برای ترکیب این دو مدل استفاده کنیم، نشان دهید که خطای (مجموع مجذور خطا) مدل $h(x)$ از رابطه‌ی زیر صدق می‌کند:

$$e \leq \frac{1}{4} (e_1 + e_2) + \frac{1}{2} \sqrt{e_1 e_2}$$

$$\begin{aligned} \min_{\alpha, \beta} \sum_{i=1}^n (\alpha h_1(x_i) + \beta h_2(x_i) - y_i)^2 \\ \frac{\partial L}{\partial \alpha} = \sum_{i=1}^n 2 h_1(x_i) (\alpha h_1(x_i) + \beta h_2(x_i) - y_i) = 0 \\ \Rightarrow \alpha \sum_{i=1}^n h_1^2(x_i) + \beta \sum_{i=1}^n h_1(x_i) h_2(x_i) = \sum_{i=1}^n h_1(x_i) y_i \\ \frac{\partial L}{\partial \beta} = 0 \Rightarrow \alpha \sum_{i=1}^n h_1(x_i) h_2(x_i) + \beta \sum_{i=1}^n h_2^2(x_i) = \sum_{i=1}^n h_2(x_i) y_i \end{aligned}$$

$$\begin{aligned} e &= \sum_{i=1}^n \left(\frac{h_1 + h_2}{2} - y_i \right)^2 = \frac{1}{4} \sum_{i=1}^n (h_1 - y_i + h_2 - y_i)^2 \\ &= \frac{1}{4} \sum_{i=1}^n (h_1 - y_i)^2 + (h_2 - y_i)^2 + 2(h_1 - y_i)(h_2 - y_i) \\ &= \frac{1}{4} (e_1 + e_2) + \frac{1}{2} \sum_{i=1}^n (h_1 - y_i)(h_2 - y_i) \\ &\leq \frac{1}{4} (e_1 + e_2) + \frac{1}{2} \sum_{i=1}^n |h_1 - y_i| |h_2 - y_i| \end{aligned}$$

$$\begin{aligned} &= \frac{1}{4} (e_1 + e_2) + \frac{1}{2} \sqrt{e_1 e_2} \\ &= \frac{1}{4} (e_1 + e_2) + \frac{1}{2} \sqrt{e_1 e_2} \end{aligned}$$

$$\begin{aligned} J &= \frac{1}{n_1 n_2} \sum_{y_i \in Y_1} \sum_{y_j \in Y_2} (y_i - y_j)^2 = \frac{1}{n_1 n_2} \sum_{y_i \in Y_1} \sum_{y_j \in Y_2} y_i^2 + y_j^2 - 2 y_i y_j \\ &= \frac{1}{n_1 n_2} \left(n_2 \sum_{y_i \in Y_1} y_i^2 + n_1 \sum_{y_j \in Y_2} y_j^2 - 2 \sum_{y_i \in Y_1} \sum_{y_j \in Y_2} y_i y_j \right) \\ &= \frac{1}{n_1} \sum_{y_i \in Y_1} y_i^2 + \frac{1}{n_2} \sum_{y_j \in Y_2} y_j^2 - \frac{2}{n_1 n_2} \sum_{y_i \in Y_1} \sum_{y_j \in Y_2} y_i y_j \\ &= \frac{1}{n_1} \sum_{y_i \in Y_1} y_i^2 + \frac{1}{n_2} \sum_{y_j \in Y_2} y_j^2 - 2 \sum_{y_i \in Y_1} \frac{1}{n_1} y_i \sum_{y_j \in Y_2} \frac{1}{n_2} y_j \\ J &= \frac{1}{n_1} \sum_{y_i \in Y_1} y_i^2 + \frac{1}{n_2} \sum_{y_j \in Y_2} y_j^2 - 2 m_1 m_2 \\ &= \frac{1}{n_1} \sum_{y_i \in Y_1} y_i^2 + \frac{1}{n_2} \sum_{y_j \in Y_2} y_j^2 - 2 m_1 m_2 + (m_1^2 + m_2^2) - (m_1^2 + m_2^2) \\ &= \frac{1}{n_1} \left(\sum_{y_i \in Y_1} y_i^2 - n_1 m_1^2 \right) + \frac{1}{n_2} \left(\sum_{y_j \in Y_2} y_j^2 - n_2 m_2^2 \right) + (m_1^2 + m_2^2 - 2 m_1 m_2) \\ &= \frac{1}{n_1} \left(\sum_{y_i \in Y_1} (y_i - m_1)^2 \right) + \frac{1}{n_2} \left(\sum_{y_j \in Y_2} (y_j - m_2)^2 \right) + (m_1 - m_2)^2 \\ J &= (m_1 - m_2)^2 + \frac{1}{n_1} S_1^2 + \frac{1}{n_2} S_2^2 \end{aligned}$$

سوال ۶ کاهش بُعد (۱۵) (نمره)

مجموعه داده‌ی زیر را در نظر بگیرید. با استفاده از PCA می‌خواهیم داده‌ها را به فضای یک بُعدی ببریم. جهت برداری که روش PCA برای نگاشت به فضای یک بُعدی بدست می‌آورد را محاسبه کنید.

$$D = \left\{ \begin{bmatrix} 1 \\ 0 \\ 8 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \\ 4 \end{bmatrix}, \begin{bmatrix} 4 \\ 0 \\ 6 \end{bmatrix} \right\}$$

$$\begin{aligned} X &= \begin{bmatrix} 10 & 8 \\ 8 & 2 \\ 2 & 0 \\ 4 & 0 \end{bmatrix} \quad \mu = \begin{bmatrix} 6 \\ 4 \end{bmatrix} \\ S &= \frac{1}{4} X^T X = \frac{1}{4} \begin{bmatrix} 40 & 24 \\ 24 & 40 \end{bmatrix} = \begin{bmatrix} 10 & 6 \\ 6 & 10 \end{bmatrix} \\ \det(S - \lambda I) &= 0 \Rightarrow \lambda^2 - 20\lambda + 64 = 0 \Rightarrow \lambda_1 = 16, \lambda_2 = 4 \\ u &= \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \\ Su &= \lambda_1 u \Rightarrow 10u_1 + 6u_2 = 16u_1 \Rightarrow u_1 = u_2 \\ \|u\| &= 1 \Rightarrow u_1^2 + u_2^2 = 1 \Rightarrow 2u_1^2 = 1 \Rightarrow u_1 = \pm \sqrt{2}/2 \end{aligned}$$

نشان دهید اگر توزیع $P(\theta|\alpha)$ یک conjugate prior برای پارامتر θ باشد، آن‌گاه توزیع مخلوط $\sum_{d=1}^D \lambda_d P(\theta|\alpha_d)$ هم یک conjugate prior برای θ خواهد بود.

$$\begin{aligned} D &= \{\alpha_1, \dots, \alpha_n\} \\ P(\theta|D) &\propto P(D|\theta) P(\theta) = P(D|\theta) \prod_{d=1}^D \lambda_d P(\theta|\alpha_d) \\ &= \prod_{d=1}^D \lambda_d P(D|\theta) P(\theta|\alpha_d) = \prod_{d=1}^D \lambda_d P(\theta|\alpha_d) \end{aligned}$$

سوال ۲ خوشه‌بندی مبتنی بر بهینه‌سازی (۲۴ نمره)

می‌خواهیم با استفاده از یک روش مبتنی بر بهینه‌سازی، داده‌ها را خوشه‌بندی کنیم. برای این منظور، می‌خواهیم کوچک‌ترین ابرکره‌ای که کل داده‌ها را درون خود جای می‌دهد، بیابیم. برای یافتن این ابرکره، مسأله بهینه‌سازی زیر را باید حل کنیم:

$$\min_{r, c} r^2$$

$$\text{subject to: } \|x_i - c\|^2 \leq r^2, \quad i = 1, \dots, n$$

که در رابطه‌ی بالا، r و c به ترتیب شعاع و مرکز ابرکره است. x_i نمونه‌ی i -ام و n تعداد کل نمونه‌ها می‌باشد.
(الف) (۵ نمره) مسأله‌ی بهینه‌سازی بالا را به گونه‌ای تغییر دهید تا این امکان وجود داشته باشد که بعضی از نمونه‌ها خارج از ابرکره قرار بگیرند. (واهنمایی: مشابه soft-margin SVM عمل کنید).

(ب) (۵ نمره) تابع لاگرانژین مسأله‌ی بهینه‌سازی قسمت (الف) را تشکیل دهید.

(ج) (۱۰ نمره) مسأله دوگان را برای مسأله‌ی قسمت (الف) بدست آورید.

(د) (۴ نمره) توضیح دهید که در این مدل، ضرایب لاگرانژ برای کدام نمونه‌ها صفر و برای کدام نمونه‌ها بزرگتر از صفر می‌شود.

$$p(x) = \alpha \lambda_1 e^{-\lambda_1 x} + (1 - \alpha) \lambda_2 e^{-\lambda_2 x}$$

توزیع مخلوط نمایی^۲ زیر را در نظر بگیرید. با فرض داشتن مجموعه داده‌ی $\{x_1, \dots, x_n\}$ می‌خواهیم با استفاده از روش EM، پارامترهای مدل بالا را که شامل α ، λ_1 و λ_2 است، بدست آوریم (نمونه‌ها i.i.d. هستند).

(الف) (۵ نمره) اگر متغیر تصادفی پنهان z_i را به صورت زیر تعریف می‌کنیم:

$$z_i = \begin{cases} 0 & x_i \text{ comes from the first component} \\ 1 & \text{otherwise} \end{cases}$$

تابع \log -complete likelihood را تشکیل دهید:

$$\log p(x_1, z_1, \dots, x_n, z_n | \lambda_1, \lambda_2, \alpha) = ?$$

(ب) (۱۰ نمره) مرحله E: امید ریاضی تابع \log -complete likelihood نسبت به متغیرهای پنهان را بدست آورید. برای این منظور، مقدار $E_{z_i | x_i, \alpha^t, \lambda_1^t, \lambda_2^t} [z_i]$ را نیز باید محاسبه کنید.

(ج) (۱۰ نمره) مرحله M: روابط به روزسانی پارامتر α و λ_1 را بدست آورید. برای سادگی در روابط، امید ریاضی $E_{z_i | x_i, \alpha^t, \lambda_1^t, \lambda_2^t} [z_i]$ را که در قسمت قبل بدست آمده با نماد γ_i^t نشان دهید.

(الف)

$$\log P(x_1, z_1, \dots, x_n, z_n | \lambda_1, \lambda_2, \alpha) \stackrel{i.i.d.}{=} \sum_{i=1}^n \log P(x_i, z_i | \lambda_1, \lambda_2, \alpha)$$

$$P(x_i, z_i | \lambda_1, \lambda_2, \alpha) = (\alpha \lambda_1 e^{-\lambda_1 x_i})^{1-z_i} ((1-\alpha) \lambda_2 e^{-\lambda_2 x_i})^{z_i}$$

$$\Rightarrow \log P(x_i, z_i | \lambda_1, \lambda_2, \alpha) = (1-z_i) (\log \alpha + \log \lambda_1 - \lambda_1 x_i) + z_i (\log (1-\alpha) + \log \lambda_2 - \lambda_2 x_i)$$

$$\Rightarrow \log \text{ complete likelihood} = \sum_{i=1}^n (1-z_i) (\log \alpha + \log \lambda_1 - \lambda_1 x_i) + z_i (\log (1-\alpha) + \log \lambda_2 - \lambda_2 x_i)$$

(ب)

$$Q(\theta, \theta^t) = E[\log P(x_1, z_1, \dots, x_n, z_n | \lambda_1, \lambda_2, \alpha)]$$

$$= \sum_{i=1}^n (1 - E[z_i]) (\log \alpha + \log \lambda_1 - \lambda_1 x_i) + E[z_i] (\log (1-\alpha) + \log \lambda_2 - \lambda_2 x_i)$$

$$E[z_i] = P(z_i = 1 | x_i, \alpha^t, \lambda_1^t, \lambda_2^t) = \frac{P(x_i | z_i = 1, \alpha^t, \lambda_1^t, \lambda_2^t) P(z_i = 1 | \alpha^t, \lambda_1^t, \lambda_2^t)}{P(x_i | \alpha^t, \lambda_1^t, \lambda_2^t)}$$

$$= \frac{\lambda_2^t e^{-\lambda_2^t x_i} (1-\alpha^t)}{\alpha^t \lambda_1^t e^{-\lambda_1^t x_i} + (1-\alpha^t) \lambda_2^t e^{-\lambda_2^t x_i}} = \gamma_i^t$$

(ج)

$$\frac{\partial Q}{\partial \alpha} = \sum_{i=1}^n (1 - \gamma_i^t) \frac{1}{\alpha} + \gamma_i^t \frac{-1}{1-\alpha} = 0$$

$$\Rightarrow (1-\alpha) \cdot \sum_{i=1}^n (1 - \gamma_i^t) = \alpha \sum_{i=1}^n \gamma_i^t$$

$$\Rightarrow \alpha = 1 - \frac{\sum_{i=1}^n \gamma_i^t}{n}$$

$$\frac{\partial Q}{\partial \lambda_1} = \sum_{i=1}^n (1 - \gamma_i^t) \left(\frac{1}{\lambda_1} - x_i \right) = 0$$

$$\Rightarrow \lambda_1 = \frac{n - \sum_{i=1}^n \gamma_i^t}{\sum_{i=1}^n (1 - \gamma_i^t) x_i}$$

سوال ۲ (الف)

$$\min_{r, c, e} r^2 + \alpha \sum_{i=1}^n e_i$$

$$\text{s.t. } \|x_i - c\|^2 \leq r^2 + e_i \quad i=1, \dots, n$$

$$e_i \geq 0 \quad i=1, \dots, n$$

(ب)

$$L(r, c, e, \lambda, \mu) = r^2 + \alpha \sum_{i=1}^n e_i + \sum_{i=1}^n \lambda_i (r^2 + e_i - \|x_i - c\|^2)$$

$$-\sum_{i=1}^n \mu_i e_i$$

$$\mu_i \geq 0$$

$$\lambda_i \geq 0$$

(ج)

$$g(\lambda, \mu) = \min_{r, c, e} L(r, c, e, \lambda, \mu)$$

$$\frac{\partial L}{\partial r} = 2r + \sum_{i=1}^n 2r \lambda_i = 0 \Rightarrow \sum_{i=1}^n \lambda_i = 1 \quad (1)$$

$$\frac{\partial L}{\partial c} = \sum_{i=1}^n \lambda_i (2c - 2x_i) = 0 \Rightarrow c \sum_{i=1}^n \lambda_i = \sum_{i=1}^n x_i \Rightarrow c = \sum_{i=1}^n x_i \quad (2)$$

$$\frac{\partial L}{\partial e_i} = \alpha - \mu_i \Rightarrow \mu_i = \alpha \Rightarrow \mu_i = \alpha - \lambda_i \quad (3)$$

$$\Rightarrow g(\lambda, \mu) = r^2 + \alpha \sum_{i=1}^n e_i - \sum_{i=1}^n \lambda_i r^2 - \sum_{i=1}^n \lambda_i e_i + \sum_{i=1}^n \lambda_i \|x_i - c\|^2 - \sum_{i=1}^n (\alpha - \lambda_i) e_i$$

$$\max_{\lambda} \sum_{i=1}^n \lambda_i \|x_i - c\|^2 - \sum_{j=1}^n z_j \lambda_j$$

$$\max_{\lambda} \sum_{i=1}^n \lambda_i x_i^T x_i - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j x_i^T x_j$$

$$\text{s.t. } 0 \leq \lambda_i \leq \alpha \quad i=1, \dots, n$$

$$\sum_{i=1}^n \lambda_i = 1$$

(>) نمونه‌هایی که روی سطح ابرکره قرار دارند ضرایب لاگرانژ آن‌ها بزرگتر از صفر است.

