



به نام خدا
درس یادگیری
ماشین
تمرین دوم

حدیثه مصباح
۸۱۰۱۰۲۲۵۳



پاسخ ۱.

۱-۱. مفهوم bias-variance trade off را با توجه به h_n در روش پارزن و k_n در روش KNN توضیح

دهید.

۱-۲. نشان دهید که مدل KNN برای ($K \neq 1$) تابع توزیع نامناسبی را تعریف میکند که انتگرال آن در تمام

فضا واگرا میباشد.

(۱.۱)

بایاس و واریانس دو مفهوم مهم در یادگیری ماشین هستند که بر دقت پیش‌بینی‌های مدل تأثیر می‌گذارند. بایاس به خطای ناشی از ساده‌سازی بیش از حد مدل اشاره دارد که باعث می‌شود مدل نتواند روابط پیچیده میان داده‌ها را درک کند. واریانس به خطای ناشی از حساسیت بیش از حد مدل به داده‌های آموزشی اشاره دارد که موجب می‌شود مدل به نویز و نوسانات جزئی داده‌ها واکنش نشان دهد.

در مدل‌سازی، ما با دو مشکل عمده روبرو هستیم: underfitting و overfitting. زمانی رخ می‌دهد که مدل نتواند داده‌های آموزشی را به خوبی بیاموزد و معمولاً با بایاس بالا و واریانس پایین همراه است. Overfitting زمانی اتفاق می‌افتد که مدل بیش از حد به داده‌های آموزشی وابسته شود و نتواند به خوبی به داده‌های جدید تعمیم داده شود، که این معمولاً با بایاس پایین و واریانس بالا همراه است.

در روش KNN، تعداد نزدیکترین همسایگان (k_n) تعیین می‌کند که چگونه مدل به داده‌های آموزشی واکنش نشان می‌دهد. هرچه k_n کوچک‌تر باشد، مدل حساس‌تر به داده‌های آموزشی است و ممکن است واریانس بیشتری داشته باشد. با افزایش k_n ، مدل کمتر حساس می‌شود و بایاس بیشتری خواهد داشت.

در روش پارزن، پارامتر h_n که عرض کرنل را تعیین می‌کند، بر تعادل بین بایاس و واریانس تأثیر می‌گذارد. یک h_n بزرگ‌تر منجر به یک تخمین هموارتر و بایاس بیشتر می‌شود، در حالی که یک h_n کوچک‌تر باعث می‌شود مدل انعطاف‌پذیرتر باشد اما واریانس بیشتری داشته باشد.

بنابراین، در انتخاب مدل و تنظیم فرایارامترها، یافتن تعادل مناسب بین بایاس و واریانس برای ساخت مدلی که بتواند به خوبی به داده‌های دیده نشده تعمیم دهد، حیاتی است. این مصالحه یکی از ملاحظات کلیدی در یادگیری ماشین است.

(۱.۲)

مدل KNN (k-nearest neighbors) یک روش غیرپارامتریک در یادگیری ماشین است که برای دسته‌بندی و رگرسیون استفاده می‌شود. اگر بخواهیم از KNN برای تخمین تابع چگالی احتمال استفاده کنیم، باید دقت داشته باشیم که تابع توزیعی که KNN تولید می‌کند (برای $K \neq 1$) ممکن است به‌طور کلی به عنوان یک تابع توزیع احتمال مناسب نباشد. این به دلیل آن است که تابع توزیع تخمین زده شده توسط KNN ممکن است انتگرالی واگرا داشته باشد وقتی $K \neq 1$ است.

$$\int p(x)dx \approx \sum_{i=1}^N p(x_i) \cdot V_i = \sum_{i=1}^N \frac{K}{NV_i} \cdot V_i = k \neq 1$$

این مشکل از این واقعیت نشأت می‌گیرد که KNN برای هر نقطه‌ای در فضای ویژگی یک کره می‌کشد و تعداد ثابتی از نزدیک‌ترین همسایه‌ها را در نظر می‌گیرد، نه یک فاصله ثابت. وقتی $K \neq 1$ ، حجم این کره‌ها به تعداد نقاط داده‌ای که در آن‌ها وجود دارد بستگی دارد، و این حجم‌ها می‌توانند بسیار نامتناسب باشند.

برای نمونه، در مناطقی که داده‌ها خیلی پراکنده هستند، KNN ممکن است کره‌های بزرگی را تولید کند که برای حفظ تعداد ثابت K از همسایه‌ها مورد نیاز است. در نتیجه، این مناطق دورافتاده وزن بیشتری در تخمین چگالی احتمال خواهند داشت، که می‌تواند منجر به تخمینی واگرا شود، زیرا وقتی انتگرال تابع توزیع چگالی احتمال را برای تمام فضا محاسبه کنیم، ممکن است مجموع وزن‌های مناطق پراکنده به بی‌نهایت میل کند.

در مقابل، در مناطق پرجمعیت داده‌ای، کره‌ها می‌توانند بسیار کوچک باشند، که این موضوع نیز می‌تواند باعث شود که تخمین چگالی احتمال در این مناطق بیش از حد بالا رود. این ناهماهنگی در اندازه کره‌ها می‌تواند به این معنی باشد که تابع چگالی احتمال تخمین زده شده توسط KNN برای $K \neq 1$ نمی‌تواند به درستی بر روی تمام فضای ویژگی انتگرال‌گیری شود و در نتیجه، یک تابع توزیع احتمال مناسب نخواهد بود.



پاسخ ۲.

توزیع یکنواخت $p(x)$ و پنجره پارزن $\varphi(x)$ به صورت زیر تعریف شده است.

$$p(x) \sim U(0, a)$$
$$\varphi(x) = \begin{cases} e^{-x} ; x > 0 \\ 0 ; x \leq 0 \end{cases}$$

۲-۱. نشان دهید که میانگین چنین تخمینی از پنجره پارزن به صورت زیر میشود.

$$\bar{p}_n(x) = \begin{cases} 0 ; x < 0 \\ \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}} \right) ; 0 \leq x \leq a \\ \frac{1}{a} \left(e^{\frac{a}{h_n}} - 1 \right) e^{-\frac{x}{h_n}} ; a \leq x \end{cases}$$

۲-۲. $\bar{p}_n(x)$ را بر حسب x برای $a = 1$ و $h_n = \{1, \frac{1}{4}, \frac{1}{16}\}$ رسم کنید.

۲-۳. h_n چه قدر باشد تا در بازه $0 < x < a$ مقدار بایاس کمتر از ۱ درصد باشد.

۲-۴. در شرایط h_n بخش ۳-۵ و مقدار $a = 1$ ، $\bar{p}_n(x)$ را در بازه $0 < x < 0.05$ رسم کنید.



(۲.۳)

پ) با بس - صورت r_i تعریف می‌شود:

$$E(p(x) - \hat{p}(x)) = p(x) - \bar{p}(x)$$

$$\text{bias}(x) = \frac{|p(x) - \bar{p}(x)|}{p(x)} = \frac{\frac{1}{a} - \bar{p}(x)}{\frac{1}{a}} = 1 - ap(x) = e^{-\frac{x}{h_n}}$$

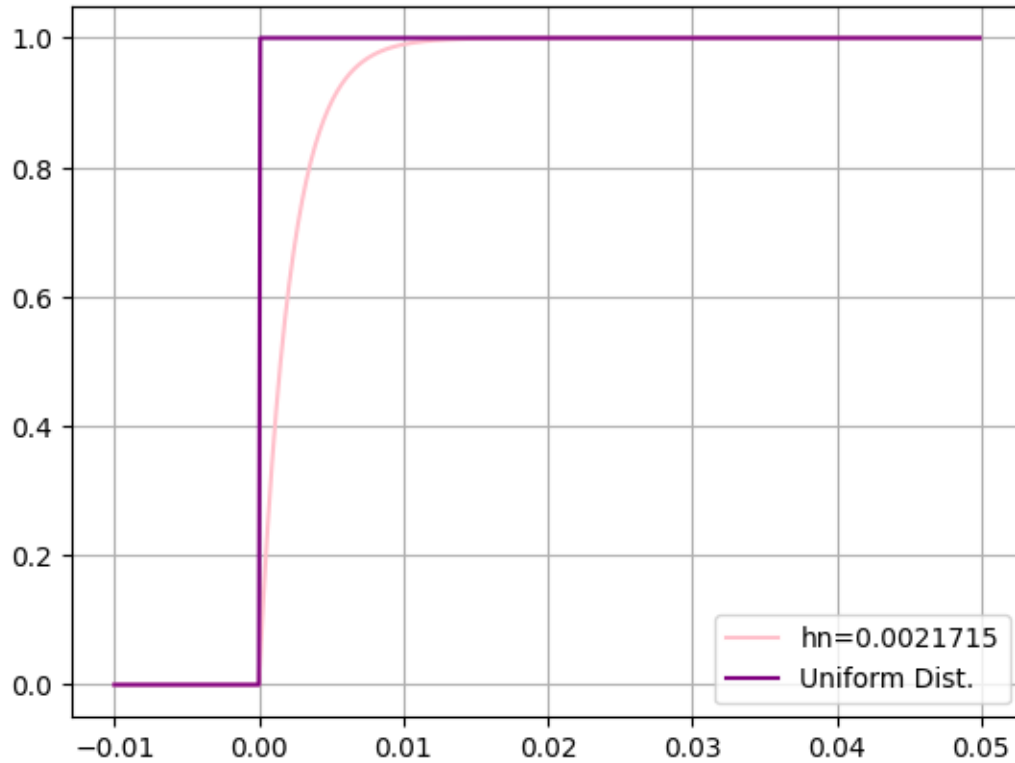
$p(x) \leq \frac{1}{a}$ $0 < x < a$ $\text{چون } p(x) \text{ یک تابع چگالی است}$

$$\rightarrow h_n = \text{bias}\left(\frac{a}{100}\right) \leq 0.01 \leftrightarrow \exp\left(\frac{a}{-100 h_n}\right) \leq 0.01$$

Equation

$$e^{-\frac{a}{100 h_n}} \leq 0.01 \xrightarrow{-\text{Ln} 0.01} -\frac{a}{100 h_n} \leq \text{Ln}(0.01) \rightarrow h_n \leq \frac{a}{100 \text{Ln} 0.01}$$

(۲.۴)



پاسخ 3.

توزیع نرمال $p(x) \sim N(\mu, \sigma^2)$ و پنجره پارزن $\varphi(x) \sim N(0,1)$ را در نظر بگیرید. نشان دهید که

تخمین پنجره پارزن $p(x) = \frac{1}{nh_n} \sum_{i=1}^n \varphi\left(\frac{x-x_i}{h_n}\right)$ برای h_n های کوچک دارای ویژگی های زیر است :

- $p_n(x) \sim N(\mu, h_n^2 + \sigma)$
- $p_n(x) - \tilde{p}_n(x) \cong \frac{1}{2} \left(\frac{h_n}{\sigma}\right)^2 \left[1 - \left(\frac{x-\mu}{\sigma}\right)^2\right] p(x)$
- $var[p_n(x)] \cong \frac{1}{2nh_n\sqrt{\pi}} p(x)$



$$P_n(x) \sim N(\mu, h_n^2 + \sigma)$$

$$\hat{P}_n(m) = E[P_n(m)] = \frac{1}{n} \sum_{i=1}^n E\left[\frac{1}{h_n} \varphi\left(\frac{x - x_i}{h_n}\right)\right] \rightarrow$$

$$\begin{aligned} \text{iid} \rightarrow \varphi &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x}{\sigma})^2} \\ P_n(m) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} \rightarrow \hat{P}_n(m) = \frac{1}{h_n} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-x'}{h_n}\right)^2\right) \times \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2}\left(\frac{x'-\mu}{\sigma}\right)^2\right) dx' \\ e^{\text{جوابی}} \rightarrow & \frac{-1}{2} \left[\frac{x'^2(\sigma^2 + h_n^2) - 2x'(x\sigma^2 + \mu h_n^2) + x^2\sigma^2 + h_n^2\mu^2}{h_n^2\sigma^2} \right] \end{aligned}$$

$$\alpha = \frac{x\sigma^2 + \mu h_n^2}{\sqrt{\sigma^2 + h_n^2}} \quad \beta = x^2\sigma^2 + h_n^2\mu^2 \quad y = x' \sqrt{\sigma^2 + h_n^2}$$

$$\textcircled{1} \quad \frac{1}{2} \left(\frac{y^2 - 2x'(\frac{x\sigma^2 + \mu h_n^2}{\sqrt{\sigma^2 + h_n^2}}) + \beta}{h_n^2\sigma^2} \right) = \frac{1}{2} \left(\frac{(y-\alpha)^2 - \alpha^2 + \beta}{h_n^2\sigma^2} \right)$$

$$\textcircled{2} \quad dx' \rightarrow \frac{1}{\sqrt{\sigma^2 + h_n^2}} dy$$

$$\textcircled{1} \textcircled{2} \quad \frac{1}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}h_n} \frac{1}{\sqrt{\sigma^2 + h_n^2}} e^{-\frac{1}{2} \frac{(y-\alpha)^2}{h_n^2\sigma^2}} \int e^{-\frac{1}{2} \frac{(y-\alpha)^2}{h_n^2\sigma^2}} dy$$

$$\star \text{ Case} = \frac{1}{\sigma^2 + h_n^2} (x^2 + 2\mu x + \mu^2) \times \frac{(x-\mu)^2}{\sigma^2 h_n^2} \frac{1}{\sqrt{2\pi}h_n\sigma}$$

$$\text{حل نهایی} \rightarrow \hat{P}_n(m) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{\sigma^2 + h_n^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2 + h_n^2}}$$

$$\text{نتیجه} \quad N(\mu, \sigma^2 + h_n^2) \quad \text{یعنی} \quad \hat{P}_n(m)$$

$$P_n(m) - \bar{P}_n(m) \leq \frac{1}{2} \left(\frac{h_n}{\sigma} \right)^2 \left[1 - \left(\frac{m - \mu}{\sigma} \right)^2 \right] P(m)$$

① برای این فرمول (میانگین) (میانگین)

$$P(m) - \bar{P}_n(m) \leq \frac{1}{\sqrt{n} \sigma} e^{-\frac{1}{2} \left(\frac{m - \mu}{\sigma} \right)^2} - \frac{1}{\sqrt{h_n^2 + \sigma^2}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{m - \mu}{h_n^2 + \sigma^2} \right)^2\right)$$

$$\frac{\bar{P}_n(m)}{P(m)} \leq \frac{\sigma}{\sqrt{\sigma^2 + h_n^2}} \exp\left[-\frac{(m - \mu)^2}{2} \left(\frac{1}{h_n^2 + \sigma^2} - \frac{1}{\sigma^2} \right)\right]$$

① حاشیای در انتهای

$$P(m) - \bar{P}_n(m) \leq P(m) \left(1 - \frac{\sigma}{\sqrt{\sigma^2 + h_n^2}} \exp\left[-\frac{(m - \mu)^2}{2} \left(\frac{1}{h_n^2 + \sigma^2} - \frac{1}{\sigma^2} \right)\right] \right)$$

$$\sqrt{1 + \left(\frac{h_n}{\sigma} \right)^2}$$

$$\frac{1}{\sqrt{1 + \left(\frac{h_n}{\sigma} \right)^2}} \approx 1 - \frac{1}{2} \left(\frac{h_n}{\sigma} \right)^2$$

$$e^x \approx 1 + \frac{h_n^2 (m - \mu)^2}{2 \sigma^2 (h_n^2 + \sigma^2)}$$

و اینها برای تقریب تابع میانی داریم.

با استفاده از دنباله تیلور داریم:

$$P(m) - \bar{P}_n(m) \leq \left(1 - \left(1 - \frac{1}{2} \left(\frac{h_n}{\sigma} \right)^2 \right) \left(1 + \frac{h_n^2 (m - \mu)^2}{2 \sigma^2 (h_n^2 + \sigma^2)} \right) \right) P(m)$$

$$\leq \frac{1}{2} \left(\frac{h_n}{\sigma} \right)^2 \left[1 - \frac{(m - \mu)^2}{(h_n^2 + \sigma^2)} \right] P(m)$$

$$P(m) - \bar{P}_n(m) \leq \frac{1}{2} \left(\frac{h_n}{\sigma} \right)^2 \left[1 - \left(\frac{m - \mu}{\sigma} \right)^2 \right] P(m)$$

$$\text{Var}(P(m)) \leq \frac{1}{2 n h_n \pi} P(m)$$

$$\text{Var}[P_n(m)] = \text{Var}\left[\frac{1}{n h_n} \sum_{i=1}^n \varphi\left(\frac{x - m_i}{h_n}\right) \right] = \frac{1}{n h_n^2} \text{Var}\left[\varphi\left(\frac{x - m'}{h_n}\right)\right]$$

$$= \frac{1}{n h_n^2} \left(E\left[\varphi^2\left(\frac{x - m'}{h_n}\right)\right] - \left(E\left[\varphi\left(\frac{x - m'}{h_n}\right)\right]\right)^2 \right)$$

$$\begin{aligned}
 \text{معمایه سار} \rightarrow &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x-x'}{h_n}\right)^2\right) \\
 \exp\left[-\frac{1}{2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2\right) \frac{1}{\sqrt{2\pi} \sigma} dx'\right] \\
 &= \frac{1}{\sqrt{2\pi} \frac{h_n^2}{2} + \sigma^2} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\frac{h_n^2}{2} + \sigma^2}\right) \\
 \text{در نتیجه:} & \frac{1}{n h_n^2} E\left(\varphi\left(\frac{x-\mu'}{h_n}\right)\right)^2 = \frac{1}{n h_n^2} \frac{h_n^2}{\sqrt{2\pi} \sqrt{\frac{h_n^2}{2} + \sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2 \frac{h_n^2}{2} + \sigma^2}\right) \propto \\
 & \text{ماتریس همبستگی ارتباط با هم دارد:} \\
 \text{Var}[P_n(\mu)] &= \frac{1}{n h_n^2} \left(E\left[\varphi^2\left(\frac{x-\mu'}{h_n}\right)\right] - \left(E\left[\varphi\left(\frac{x-\mu'}{h_n}\right)\right]\right)^2 \right) \propto \frac{P(\mu)}{2 n h_n \sqrt{\pi}}
 \end{aligned}$$

پاسخ ۴.

در هر ۲ سوال زیر متغیرهای X و θ به ترتیب بیانگر نمونه‌های مشاهده شده و پارامترهای مسئله می‌باشند.

۴-۱. فرض کنید که توزیع پارامتر $P(\theta)$ می‌باشد. مراحل expectations و maximization را برای بیشینه کردن $P(\theta|X)$ بنویسید. (محاسبات را تنها به صورت پارامتری بنویسید).

راهنمایی: از نامساوی جنسن می‌دانیم:

$$\begin{aligned}
 p(\theta | x) &\propto p(x | \theta) p(\theta) \propto \left(\sum_z Q(z) \frac{p(x, z | \theta)}{Q(z)} \right) p(\theta) \\
 \ln\left(\sum_z Q(z) \frac{p(x, z | \theta)}{Q(z)} \right) &\geq \sum_z Q(z) \{ \ln(p(x, z | \theta)) - \ln(Q(z)) \}
 \end{aligned}$$

۴-۲. متغیر تصادفی X با ۴ حالت طبق جدول زیر مفروض است. فرض کنید θ یک عدد حقیقی در بازه $[0, 1]$ و احتمال هر حالت مطابق زیر می‌باشد.



State	Probability
A	$\frac{1}{3}$
B	$\frac{1}{3} (1 - \theta)$
C	$\frac{2}{3} (\theta)$
D	$\frac{1}{3} (1 - \theta)$

با فرض انجام n آزمایش روی X ، حالت های A, B, C, D به تعداد n_a, n_b, n_c, n_d بار به دست آمده است. متاسفانه مقدار متغیر های n_c و n_a ناشناخته است. فرض کنید که تابع توزیع θ در ابتدا به صورت زیر نوشته شده است. مراحل E و M را نوشته و θ را بدست آورید.

$$p(\theta) = \frac{\Gamma(v_1 + v_2)}{\Gamma(v_1) \Gamma(v_2)} \theta^{v_1-1} (1 - \theta)^{v_2-1}$$

راهنمایی: توجه کنید که داده های با لیبیل های B و D پس از مشاهده مشخص است.

فرض کنید تابع درست‌نمایی به صورت زیر است:

$$L(\theta; x^1, \dots, x^m) = \prod_{i=1}^m p(x^i; \theta)$$

که در آن θ پارامترهای مدل و x^1, \dots, x^m داده‌های مستقل در مجموعه آموزشی هستند. تابع $p(x^i; \theta)$ احتمال مشاهده داده x را با توجه به پارامترهای فعلی θ مشخص می‌کند.

اگر این تابع درست‌نمایی برای یک مدل مخلوط (مانند مدل با متغیرهای پنهان) باشد، می‌توانیم از الگوریتم EM (Expectation-Maximization) برای برازش پارامترها استفاده کنیم. در این روش، ما ابتدا انتظارات (یا تخمین‌هایی) برای داده‌های پنهان ایجاد می‌کنیم (مرحله E) و سپس با استفاده از این انتظارات، پارامترهای مدل را به‌روزرسانی می‌کنیم (مرحله M).



این فرایند به صورت تکراری ادامه می‌یابد تا زمانی که به یک تخمین ثابت برای پارامترها برسیم که دیگر تغییر قابل توجهی نکند. این تخمین‌های نهایی به ما اجازه می‌دهند تا مدل را بر اساس داده‌های موجود برازش دهیم.

تابع Loglikelihood به شرح زیر است:

$$l(\theta) = \sum_{i=1}^m \log p(x; \theta) = \sum_{i=1}^m \log \sum_z p(x, z; \theta)$$

یافتن برآوردهای صریح درست‌نمایی حداکثری (ML) برای پارامترهای θ ممکن است دشوار باشد. در اینجا، $z^{(i)}$ متغیرهای تصادفی هستند که به آن‌ها دسترسی نداریم و اغلب این‌طور است که اگر $z^{(i)}$ مشاهده شود، تخمین ML آسان خواهد بود. در چنین مواقعی، الگوریتم EM یک روش کارآمد برای تخمین ML ارائه می‌دهد. برای هر Q_i, i مقداری توزیع بر روی z است ($Q_i(z) \geq 0, \sum_z Q_i(z) = 1$) روابط زیر را خواهیم داشت:

$$\begin{aligned} \sum_i \log p(x^{(i)}; \theta) &= \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \\ &= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \\ &\geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \end{aligned}$$

آخرین مرحله این رویکرد این است که از نابرابری Jensen استفاده نماییم. به طور خاص، $f(x) = \log(x)$ یک تابع مقعر است، زیرا $f''(x) = -\frac{1}{x^2} < 0$ در دامنه $x \in R^+$ منفی می‌باشد. همچنین، ترم $\sum_{z^{(i)}} Q_i(z^{(i)}) \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$ در جمع سیگما فقط یک امید از کمیت $\left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$ با توجه به $z^{(i)}$ بر اساس توزیع داده شده توسط Q_i ترسیم شده است. حابل با استفاده از نامساوی Jensen داده شده در صورت سوال خواهیم داشت:

$$f \left(E_{z^{(i)} \sim Q_i} \left[\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right) \geq E_{z^{(i)} \sim Q_i} \left[f \left(\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) \right]$$

بقیه ش توی عکسه وقت تایپ نبود ☺



Subject

Date

ما $Q_i(z^{(i)} | a^{(i)})$ Posterior $z^{(i)}$ بر طبق Bayes می‌توانیم بنویسیم

① E-step ← انتخاب $Q_i(z^{(i)})$ از میان $P(z^{(i)} | a^{(i)})$ که این کار را می‌توانیم به صورت زیر بنویسیم

② M-step ← $Q_i(z^{(i)})$ را به جای $P(z^{(i)} | a^{(i)})$ در θ قرار می‌دهیم و θ را بهینه می‌کنیم

نکته این دو مرحله در EM همگرا می‌شوند

① ابتدا $Q_i(z^{(i)}), P(z^{(i)} | a^{(i)}; \theta)$

$$\textcircled{2} \theta = \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(a^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

توجه: این دو مرحله

تکرار می‌شوند تا $\theta^{(t)}$ و $\theta^{(t+1)}$ به هم نزدیک شوند. EM همگرا می‌شود

که هم Log Likelihood را بهینه می‌کند

$$L(\theta) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(a^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

و این $\theta^{(t+1)}$ را می‌توانیم به جای $\theta^{(t)}$ در $L(\theta)$ قرار دهیم

$$L(\theta^{(t+1)}) \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(a^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i(z^{(i)})}$$

$$\rightarrow \geq \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(a^{(i)}, z^{(i)}; \theta^{(t)})}{Q_i(z^{(i)})} = L(\theta^{(t)})$$

توجه: این دو مرحله را می‌توانیم به صورت زیر بنویسیم

① $Q_i(z^{(i)})$ و $P(z^{(i)} | a^{(i)}; \theta^{(t)})$ را به هم نزدیک می‌کنیم

② θ را بهینه می‌کنیم تا $\theta^{(t+1)}$ را به دست آوریم. در این مرحله $\theta^{(t)}$ را به جای θ قرار می‌دهیم

و در نهایت به $L(\theta^{(t)})$ می‌رسیم. این دو مرحله را می‌توانیم به صورت زیر بنویسیم



در نهایت این باعث می‌شود که احتمال به صورت یکنواخت پخش شود و EM هم‌تکامل می‌شود.
 انعام می‌شود. و یک آزمون می‌تواند این است که بررسی کنیم آیا افزایش $L(\theta)$ بین تکرارهای متوالی
 کاهش یافته یا نه، و اگر نه، تفاوتی اعلام می‌شود اگر EM ثابت باشد.

این را به این شکل می‌نویسند:

$$L(\theta) = \sum_{i=1}^n \sum_{j=1}^K Q_i(z^{(i)}) \log(P(\alpha^{(i)}, z^{(i)}; \theta)) - \log(Q_i(z^{(i)}))$$

ما از این برای اهدافی قبلی خود می‌دانیم که $\theta \geq T(\theta, \theta)$ و همچنین می‌دانیم EM باید بهینه باشد.
 در این مسئله که در آن مرحله θ آن را یک بسته به هم می‌زنیم و بهینه می‌کنیم و بهینه می‌کنیم آن را با هم
 به هم می‌زنیم و بهینه می‌کنیم.

(۴.۲)

مثال چهارم: (۲ صفت)

$$L(\theta) = \log P(n, z | \theta) P(\theta)$$

$$\log P(n, z | \theta) = \log P(\theta) + \log P(n, z | \theta)$$

$$P(n, z | \theta) = P(n | z, \theta) P(z | \theta) \quad \log P(n, z | \theta) = \log P(n | z, \theta) + \log P(z | \theta)$$

① $a, b \rightarrow \log P(a, b | z, \theta) + \log P(z | \theta) \rightarrow nb \log(\frac{1}{3} (1-\theta))$

② $a, d \rightarrow \log P(a, d | z, \theta) + \log P(z | \theta) \rightarrow nd \log(\frac{1}{3} (1-\theta))$

③ $z, a, o, r, c \rightarrow \log P(a, o, r, c | z, \theta) + \log P(z | \theta)$

مثال $\rightarrow (na + nc) (1 + \log \frac{1}{3})$

$$P(z | \theta) P(a, o, r, c | z, \theta) = \frac{1}{1 + \theta}$$

مثال $\rightarrow \frac{1}{1 + \theta} (na + nc) (1 + \log \frac{1}{3})$

④ $a, o, r, c \rightarrow \log P(a, o, r, c | z, \theta) + \log P(z | \theta)$

$$\frac{\theta^2}{1 + \theta^2} (na + nc) \log(\frac{2}{3} | \theta)$$

Kian



Subject

Date

$$\ln(L(\theta)) = \ln(P(\theta)) + (n_b + n_d) \ln\left(\frac{1}{3}(1-\theta)\right) + \frac{1}{1+\hat{\theta}} (n_a + n_c) + \ln\left(\frac{1}{3}\right) \\ + \frac{\hat{\theta}}{1+\hat{\theta}} (n - (n_b + n_d)) \ln\left(\frac{\hat{\theta}}{3}\right) \quad (1)$$

$$\ln(P(\theta)) = \ln\left(\frac{\Gamma(v_1 + v_2)}{\Gamma(v_1) \Gamma(v_2)} \theta^{v_1-1} (1-\theta)^{v_2-1}\right)$$

$$\ln(L(\theta)) = (v_1 - 1) \ln \theta + (v_2 - 1) \ln(1-\theta) \quad (2)$$

$$\frac{\partial}{\partial \theta} \ln(L(\theta)) = \frac{2(v_1 - 1)}{2\theta} + \frac{(1 - v_2)}{1-\theta} + \frac{2\hat{\theta}}{1+2\hat{\theta}} \frac{(n - (n_b + n_d))}{2\theta}$$

$$\frac{v_2 - 1}{1-\theta} = \frac{2v_1 - 2}{2\theta} \frac{(n - (n_b + n_d))}{1+2\hat{\theta}} \left(\frac{\hat{\theta} \times 2}{1+2\hat{\theta}} \right)$$

$$\frac{v_2 - 1}{1-\theta} = \frac{\alpha}{2\theta} \quad \text{where } \alpha = 2\theta v_2 - 2\theta \frac{\alpha - \alpha\theta}{\theta(2v_2 - 2 + \alpha)} = \alpha$$

$$\theta = \frac{\alpha(2v_1 - 2)(n - (n_b + n_d)) \left(\frac{\hat{\theta} \times 2}{1+2\hat{\theta}} \right)}{2v_2 + 2\theta(2v_1 + 2)(n - (n_b + n_d)) \frac{2\hat{\theta}}{1+2\hat{\theta}}}$$

پاسخ ۵.

فرض کنید K بازیکن در روز t وجود دارد. یکی از آنها به تعداد m_t بار بازی می کند و تعداد w_t بازی را می برد. شما تنها تعداد کل این افراد، تعداد کل راند های بازی شده و تعداد بازی های برده شده توسط بازیکن را می دانید اما نمیدانید کدام یک از K بازیکن در کدام روز بازی کرده است. شما می خواهید از یادگیری ماشین برای حل این مسئله استفاده کنید. برای هر یک از K بازیکن شما یک مدل احتمالی می سازید که در آن فرد با احتمال p_k بازی را می برد. بنابراین در روز t اگر بازیکن i ام به تعداد m_t بار بازی کرده باشد، احتمال آن که w_t بازی را ببرد توسط یک توزیع دوجمله ای بیان می شود. (I)

در این مسئله شما باید از یک مدل ترکیب شده با K متغیر تصادفی دوجمله ای با پارامترهای p_1, p_2, \dots, p_K استفاده کنید و برای N روز داده شده به صورت $(m_1, w_1), \dots, (m_n, w_n)$ بکار ببرید. به این صورت که در ابتدا در روز t ما ابتدا یک بازیکن از کل آنها با احتمال π انتخاب می کنیم (C_t). در مرحله بعد در روز t این بازیکن m_t بار بازی می کند و با دانستن اینکه بازیکن C_t است، تعداد برد w_t با یک متغیر دوجمله ای توصیف می شود.

(II)

۵-۱. روابط توصیف شده (I) و (II) را بنویسید.

۵-۲. مرحله E برای بروزرسانی Q با توجه به پارامتر های مرحله قبل بدست آورید. (نشان دهید در دور i ام $Q_t^{(i)}[k]$ چیست)

۵-۳. برای هر مدل ترکیبی مرحله M برای π در دور i ام از رابطه زیر محاسبه می شود.

$$\pi^i[k] = \frac{\sum_{t=1}^n Q_t^i[k]}{n}$$

مرحله M را برای بروزرسانی پارامتر های مدل $p_1^t, p_2^t, \dots, p_K^t$ در دور i ام برحسب داده و مقادیر $Q_t^{(i)}$ بدست آورید. در ابتدا مرحله maximization را برای پارامتر ها نشان داده و مسئله بهینه سازی را حل کنید.

(۵.۱)

برای مدل کردن بخش I ، می توانیم از یک رویکرد احتمالاتی با استفاده از توزیع های دوجمله ای بهره ببریم. در این مسئله، ما K بازیکن داریم. باید تعداد بازی های روزانه و بردها را به گونه ای نمایش دهیم که با فرضیات مسئله سازگار باشد. فرض کنید در روز t یک بازیکن m_t بار بازی کرده و w_t بار برنده شده است. احتمال برنده شدن هر بازیکن با شرط اینکه هر بازیکن k که $(k=1,2,3,\dots,K)$ با احتمال p_t پیروز می شود، مشخص می گردد. با توجه به این مولفه ها، می توانیم احتمال مشاهده w_t برد از m_t بازی برای بازیکن i ام در روز t را با استفاده از توزیع



دوجمله‌ای تعریف کنیم. توزیع دوجمله‌ای، تعداد موفقیت‌ها (در این مسئله، بردها) را در تعداد مشخصی از آزمایش‌ها (بازی‌های انجام شده) بر اساس احتمال موفقیت در هر آزمایش توصیف می‌کند. رابطه‌ی آن به شکل زیر است.

$$P(w_t|m_t, p_i) = \binom{m_t}{w_t} p_i^{w_t} (1 - p_i)^{m_t - w_t}$$

برای تخمین احتمالات k_p هر بازیکن، بر اساس داده‌های مشاهده‌شده از بردها و باخت‌ها، می‌توانیم از یک فرآیند تکراری مانند روش EM (Expectation-Maximization) استفاده کنیم. در این رویکرد، ابتدا تخمین می‌زنیم که کدام بازیکن به احتمال زیاد در هر روز بازی کرده است. این کار با استفاده از داده‌های موجود و یک مدل اولیه انجام می‌شود (مرحله Expectation).

سپس، بر اساس این تخصیص اولیه، تخمین‌های احتمالات k_p را به‌روزرسانی می‌کنیم (مرحله Maximization). این فرآیند به صورت تکراری انجام می‌شود تا زمانی که به یک معیار خاص برای همگرایی برسیم یا تغییرات در تخمین‌ها کمتر از یک آستانه مشخص شوند. به این ترتیب، می‌توان احتمالات برنده شدن هر بازیکن را به طور دقیق‌تری بر اساس داده‌های موجود تخمین زد.

برای مدلسازی بخش II، ما یک فرآیند دو مرحله‌ای را برای هر روز در نظر می‌گیریم. این فرآیند شامل دو مرحله است:

انتخاب بازیکن: در روز t یک بازیکن c_t از بین K بازیکن با احتمال $\pi(c_t)$ انتخاب می‌شود. این احتمال نشان می‌دهد که هر بازیکن چه شانس برای بازی کردن در روز t دارد. به عبارت دیگر، ما برای هر بازیکن یک احتمال اختصاص می‌دهیم که نشان‌دهنده میزان احتمال بازی کردن آن‌ها در روز مشخص است. این احتمالات می‌توانند بر اساس عوامل مختلفی مانند عملکرد قبلی، سطح مهارت، و دیگر معیارهای مرتبط تعیین شوند.

نتایج بازی: وقتی یک بازیکن c_t انتخاب می‌شود، آنها m_t بازی را در روز t انجام می‌دهند. تعداد بازی‌های برده شده، w_t بسته به احتمال برنده شدن بازیکن انتخابی، به عنوان یک متغیر تصادفی دوجمله‌ای مدل می‌شود.



با این توضیحات، می‌توانیم روابط مربوطه را به این صورت توصیف کنیم:

- $c_t \in \{1, 2, \dots, K\}$ نشان‌دهنده بازیکنی است که در روز t بازی می‌کند.

- احتمال انتخاب بازیکن t در روز t با $\pi(c_t)$ نشان داده می‌شود.

- بر اساس اینکه بازیکن c_t در روز t بازی می‌کند، تعداد بردهای w_t از توزیع دوجمله‌ای با پارامترهای m_t (تعداد بازی‌های انجام شده) و $\pi(c_t)$ (احتمال برنده شدن بازیکن c_t) پیروی می‌کند.

- احتمال مشاهده w_t پیروزی در روز t از رابطه‌ای خاص به دست می‌آید که بر اساس توزیع دوجمله‌ای و پارامترهای مذکور تعیین می‌شود.

$$P(w_t | m_t, c_t) = \binom{m_t}{w_t} p_{c_t}^{w_t} (1 - p_{c_t})^{m_t - w_t}$$

احتمال کلی داده‌های مشاهده شده برای تمام N روز به شرح زیر است

$$L = \prod_{t=1}^N \pi(c_t) \binom{m_t}{w_t} p_{c_t}^{w_t} (1 - p_{c_t})^{m_t - w_t}$$

برای مدل کلی که برای N روز در نظر گرفته شده و شامل داده‌های مشاهده $(m_1, w_1), \dots, (m_N, w_N)$ شده است، احتمال داده‌ها با توجه به پارامترهای p_1, p_2, \dots, p_K مدل و π می‌تواند به صورت زیر بیان شود:

$$P((m_1, w_1), \dots, (m_N, w_N) | p_1, \dots, p_K, \pi) = \prod_{t=1}^N \sum_{c_t=1}^K \pi(c_t) \binom{m_t}{w_t} p_{c_t}^{w_t} (1 - p_{c_t})^{m_t - w_t}$$

در این مدل، ما احتمالات مربوط به همه بازی‌ها را در تمام روزها محاسبه می‌کنیم. این کار شامل مجموعه‌ای از احتمالات برای هر بازیکن است که می‌توانستند در هر روز بازی کنند. در اینجا، دو هدف اصلی مدل وجود دارد:

برآورد احتمال برنده شدن P_t برای هر بازیکن و احتمال انتخاب $\pi(c_t)$ برای هر بازیکن: این به معنی تخمین احتمال این است که یک بازیکن خاص در یک روز مشخص انتخاب شده باشد. **دو بخش بعدی در عکس هست وقت نشد**

تایپ شه



5.2

سوال و جمع

وقتی ما الگوریتم EM را برای احتمال می‌زنیم در مرحله (E) ما باید به هر کدام از پارامترهای مدل احتمال را به دست آوریم. \log likelihood را حساب می‌کنیم که شامل $\pi(E_k)$ و P_k برای هر پارامتر که می‌خواهیم. در مرحله M ما پارامترها را بر اساس مقادیر قبلی و مجموع لاری‌ها و احتمال این که هر پارامتر C_i در هر روز t از k باشد را حساب می‌کنیم. $Q_k^{(i)}[k]$ به عنوان احتمال این که پارامتر k در مرحله i ام الگوریتم در روز t از k باشد نشان می‌دهیم.

$$E \text{ مرحله} \rightarrow Q_k^{(i)}[k] = \frac{\pi^{(i)}(k) \cdot \left(\frac{m_b}{w_b}\right) (P_k^{(i)})^{w_b} (1-P_k^{(i)})^{m_b-w_b}}{\sum_{j=1}^k \pi^{(i)}(j) \left(\frac{m_b}{w_b}\right) (P_j^{(i)})^{w_b} (1-P_j^{(i)})^{m_b-w_b}}$$

Posterior

$\pi^{(i)}(k)$ نشان دهنده احتمال انتخاب پارامتر k است.
 $P_k^{(i)}$ به هر پارامتر احتمال روز t از k در مرحله i ام است.
 m_b مجموع هر یک از پارامترها است. هر روز t پارامترها را جمع می‌کنیم تا مطمئن شویم این پارامترها احتمال روزی می‌دهد که پارامترها هستند.

5.3

ما می‌خواهیم پارامترها را به دست آوریم. \log likelihood را حساب می‌کنیم. \log likelihood را حساب می‌کنیم.

$$\log \text{ likelihood} \rightarrow \ln \left(\prod_{i=1}^n \prod_{k=1}^k P_k^{w_b} \left(\frac{m_b}{w_b}\right) (1-P_k)^{m_b-w_b} C_i \right)$$

$$\sum_{i=1}^n \sum_{k=1}^k \left(\ln \left(\frac{m_b}{w_b}\right) + w_b \ln(P_k) + (m_b - w_b) \ln(1-P_k) + \ln(C_i) \right)$$

ماتریس C_i و C_n

$$\frac{\left(\frac{m_b}{w_b}\right) P_i^{w_b} (1-P_i)^{m_b-w_b} C_i}{\sum_{i=1}^k \left(\frac{m_b}{w_b}\right) P_i^{w_b} (1-P_i)^{m_b-w_b} C_i} \times \left(\sum_{i=1}^n \sum_{k=1}^k \left(\ln \left(\frac{m_b}{w_b}\right) + w_b \ln(P_i) + (m_b - w_b) \ln(1-P_i) + \ln(C_i) \right) \right)$$

$\frac{\partial L}{\partial P_i}$

$$m\text{-step} \quad \frac{\partial L}{\partial P_i} \quad \frac{\partial L}{\partial P_i}$$

A B



$$A \rightarrow \sum_{b=1}^N Q_{b,i} \left(\rho_b + w_b \frac{1}{P_i} + (m_b - w_b) \left(\frac{1}{1-P_i} \right) \right) \text{ so}$$
$$\sum_{b=1}^N \frac{Q_{b,i} (w_b)}{P_i} \text{ , } \sum_{b=1}^N \frac{(m_b - w_b) Q_{b,i} (w_b)}{1-P_i}$$

در نتیجه،

$$P_i = \frac{\sum_{b=1}^N Q_{b,i} w_b}{\sum_{b=1}^N (m_b - w_b) Q_{b,i} + \sum_{b=1}^N Q_{b,i} w_b}$$
$$B \rightarrow \sum_{b=1}^N \left(Q_{b,i} \left(\frac{1}{C_i} \right) - \mu \right) \text{ so } C_i \mu = \sum_{b=1}^N Q_{b,i}$$

در نتیجه

$$C_i = \frac{\sum_{b=1}^N Q_{b,i}}{N}$$

پاسخ ۶.



در این سوال هدف پیاده سازی الگوریتم KNN و استفاده از آن به عنوان طبقه بند میباشد.

۱-۶. همانطور که میدانید الگوریتم KNN ساده و شهودی است، هنگام پیش‌بینی، فاصله بین هر یک از نقاط داده موجود را محاسبه می‌کند و آن را همانند نزدیک‌ترین کلاس به آن طبقه‌بندی می‌کند. یک کلاس KNN ساخته و با استفاده از کتابخانه numpy این الگوریتم را پیاده سازی کنید.

۲-۶. مجموعه داده 'iris' را لود کرده و اطلاعات کلی دیتاست شامل تعداد کلاس و تعداد سمپل ها و فرمت داده ها و ... بیان کنید.

۳-۶. Scatter plot مجموعه داده iris را رسم کنید.

۴-۶. مجموعه داده iris را به دو دسته آموزش و ارزیابی تقسیم کنید.

۵-۶. به کمک کلاس KNN که در بخش ۱ پیاده سازی کردید ، مدلی بر روی داده های آموزش به ازای k برابر با ۵ آموزش داده سپس دقت مدل را بر روی داده های آموزش و ارزیابی گزارش کنید.

۶-۶. بخش ۵ را به ازای k های متفاوت (۱ تا ۱۰) تکرار کنید و نمودار دقت بر روی داده های ارزیابی به ازای k های متفاوت را رسم کرده و بهترین k را بر اساس آن گزارش کنید.

با استفاده از کد زیر ما میتوانیم الگوریتم KNN را پیاده سازی کنیم



```
import numpy as np

class KNN:
    def __init__(self, k=3):
        self.k = k

    def fit(self, X, y):
        self.X_train = X
        self.y_train = y

    def predict(self, X):
        y_pred = [self._predict(x) for x in X]
        return np.array(y_pred)

    def _predict(self, x):
        distances = [np.linalg.norm(x - x_train) for x_train in self.X_train]
        k_indices = np.argsort(distances)[:self.k]
        k_nearest_labels = [self.y_train[i] for i in k_indices]
        most_common = np.bincount(k_nearest_labels).argmax()
        return most_common
```

در قسمت بعدی کد ما خصوصیات دیتا ست را بررسی میکنیم که به شرح زیر میباشد

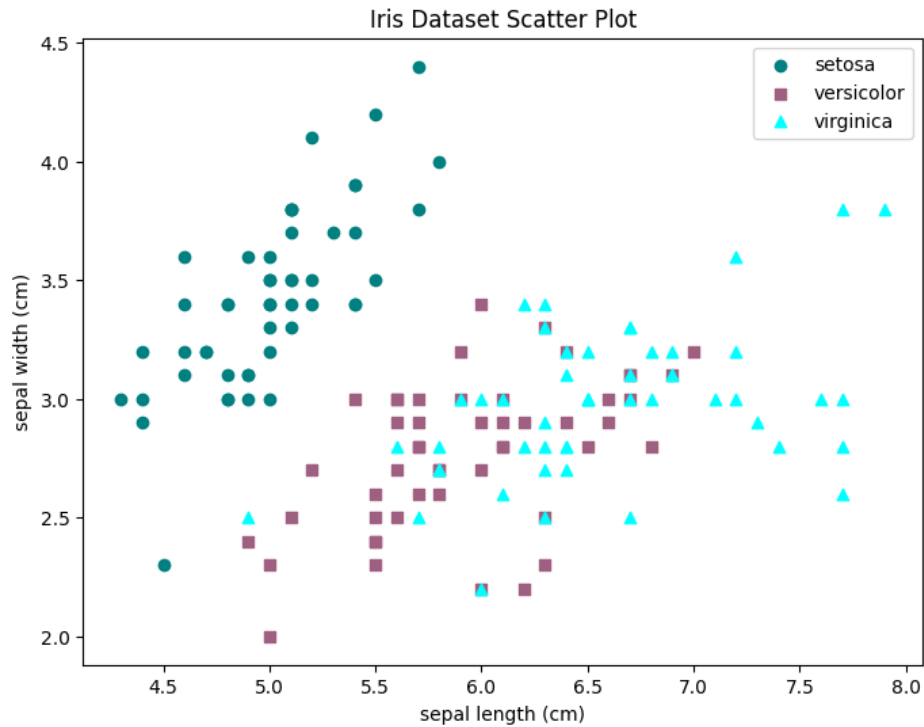
```
Number of classes: 3
Class names: ['setosa' 'versicolor' 'virginica']
Number of samples: 150
Number of features: 4
Feature names: ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

من یک قسمت نرمالایزشن هم به کدم اضافه کردم

```
from sklearn.preprocessing import MinMaxScaler
# Initialize the MinMaxScaler
scaler = MinMaxScaler()

# Fit and transform the data
X = scaler.fit_transform(data)
```

در مرحله بعدی ما scatter plot دیتا ست Iris را میکشیم



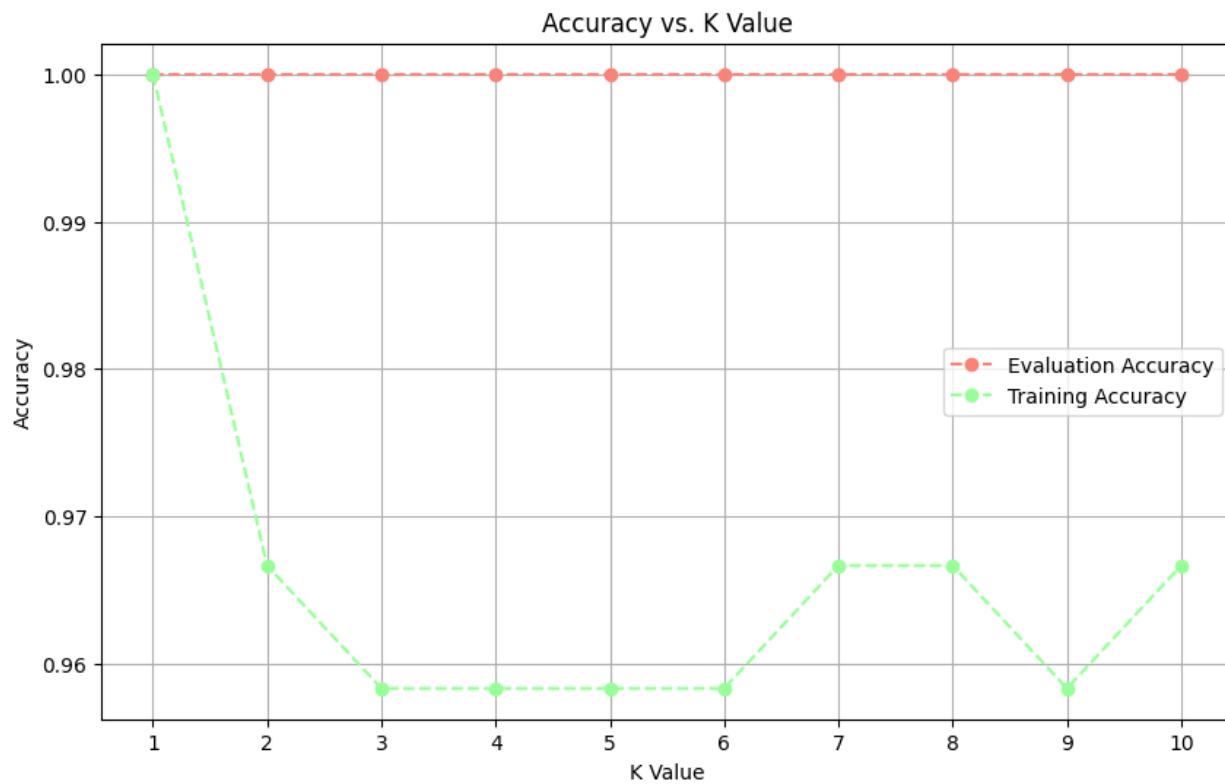
بعد از تقسیم داده به دو بخش آموزش و تست ما با استفاده از الگوریتم KNN که قبلاً خودمان تعریف کرده ایم مدل خود را به ازای $k = 5$ آموزش می‌دهی مکه دقت آن به شرح زیر است

Training accuracy: 95.83%
Evaluation accuracy: 100.00%

در مرحله بعدی به ازای $K = 1$ تا $k = 10$ این کار را تکرار می‌کنیم که دقت آن به شرح زیر است

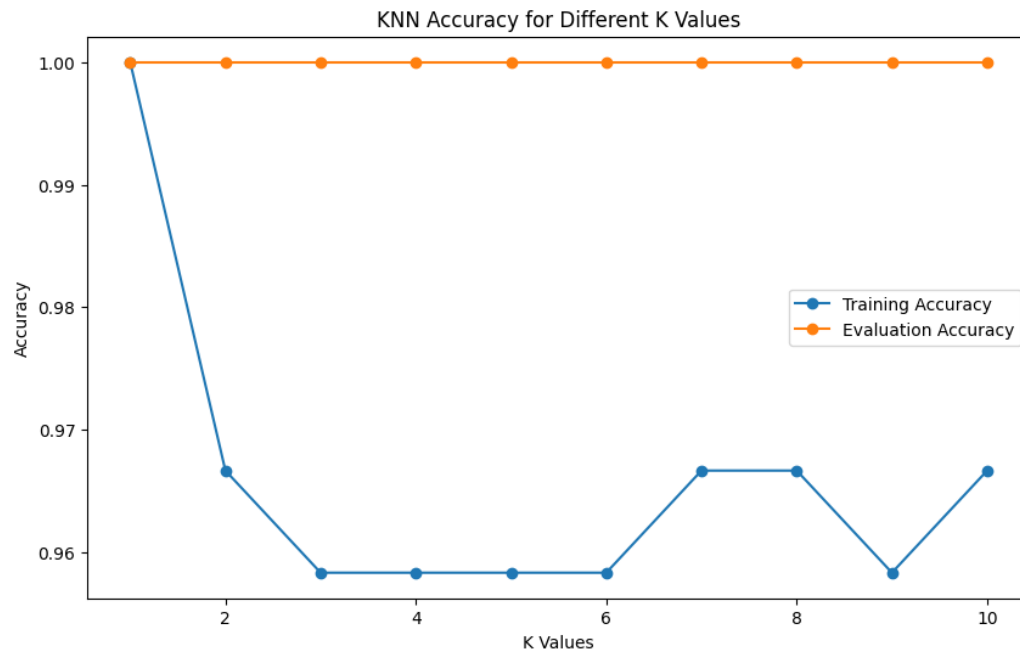
```
K = 1: Training Accuracy = 1.0000, Evaluation Accuracy = 1.0000
K = 2: Training Accuracy = 0.9667, Evaluation Accuracy = 1.0000
K = 3: Training Accuracy = 0.9583, Evaluation Accuracy = 1.0000
K = 4: Training Accuracy = 0.9583, Evaluation Accuracy = 1.0000
K = 5: Training Accuracy = 0.9583, Evaluation Accuracy = 1.0000
K = 6: Training Accuracy = 0.9583, Evaluation Accuracy = 1.0000
K = 7: Training Accuracy = 0.9667, Evaluation Accuracy = 1.0000
K = 8: Training Accuracy = 0.9667, Evaluation Accuracy = 1.0000
K = 9: Training Accuracy = 0.9583, Evaluation Accuracy = 1.0000
K = 10: Training Accuracy = 0.9667, Evaluation Accuracy = 1.0000
```

و در نهایت نمودار آن را میکشیم که به نظر می‌آید به ازای $K = 1$ بهترین نتیجه نشان داده میشود



The best k based on evaluation accuracy is: 1 with an accuracy of 1.0000

و در نهایت من برای این که متوجه درست کار کردن الگوریتم خود بشوم با استفاده از کتابخانه نیز این کد را پیاده کردم که نمودار مشابه ای به دست آمد که کد آن در فایل پایتون هست.



پاسخ ۷.

۷-۱. به کمک دستور زیر مجموعه داده X را تولید کنید.

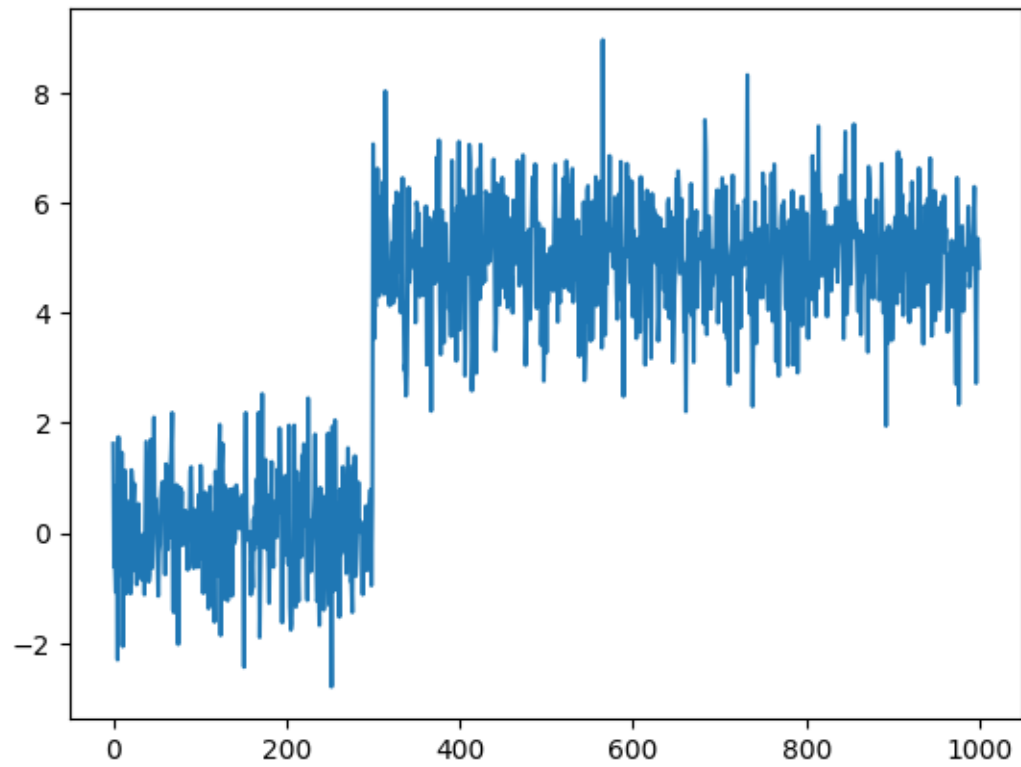
```
import numpy as np
N = 1000
np.random.seed(1)
X = np.concatenate((np.random.normal(0, 1, int(0.3 * N)),
np.random.normal(5, 1, int(0.7 * N))))[:, np.newaxis]
```

۷-۲. توزیع دیتا X را با استفاده از روش پنجره پارزن با کرنل گوسی بدست آورید.

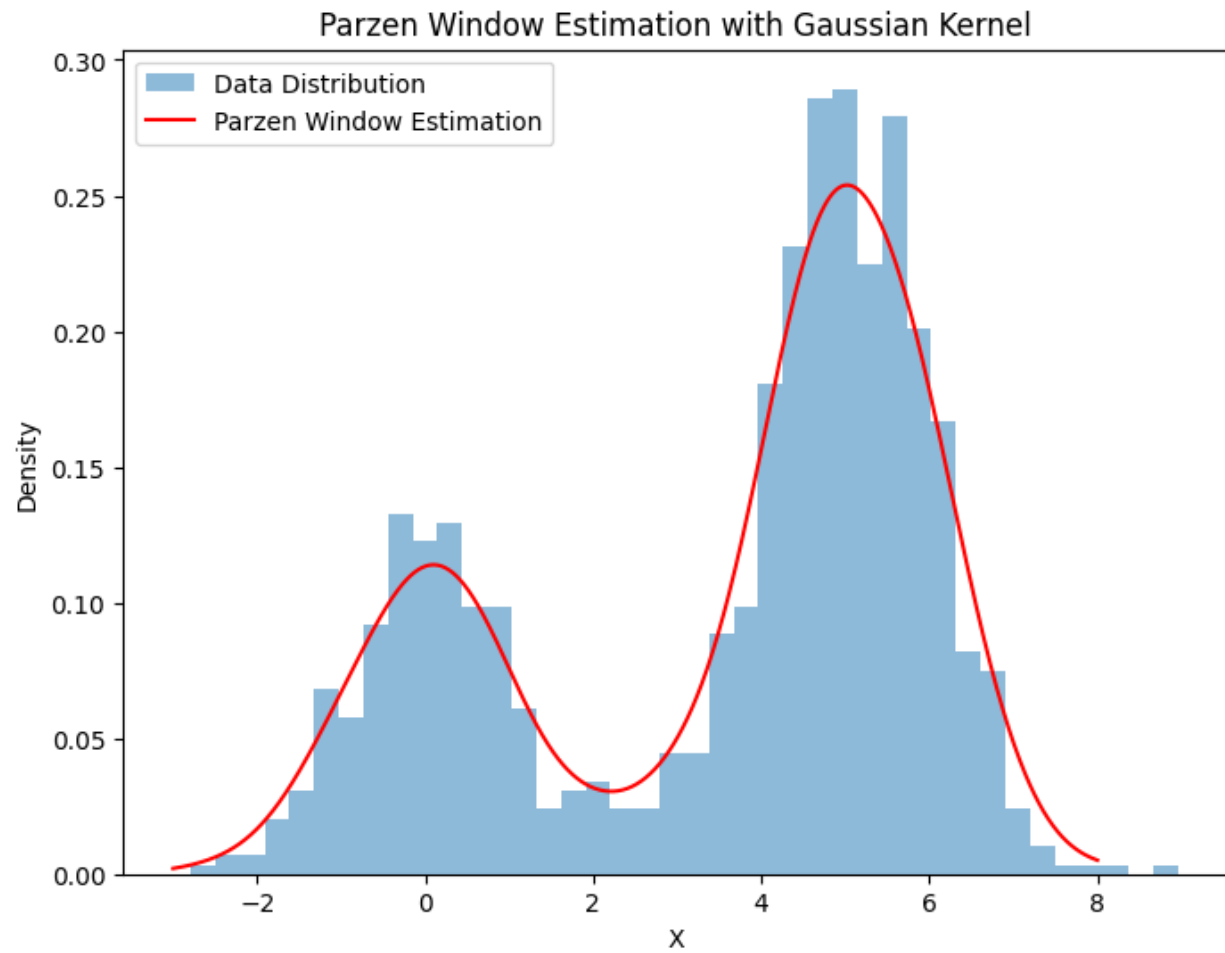
۷-۳. تاثیر اندازه پنجره پارزن را روی توزیع تخمین زده شده بررسی کنید (حداقل ۳ اندازه مختلف مثلاً: ۱۰ و

۱ و ۰.۱)

بعد از ساخت مجموعه داده اول آن را نمایش می‌دهیم

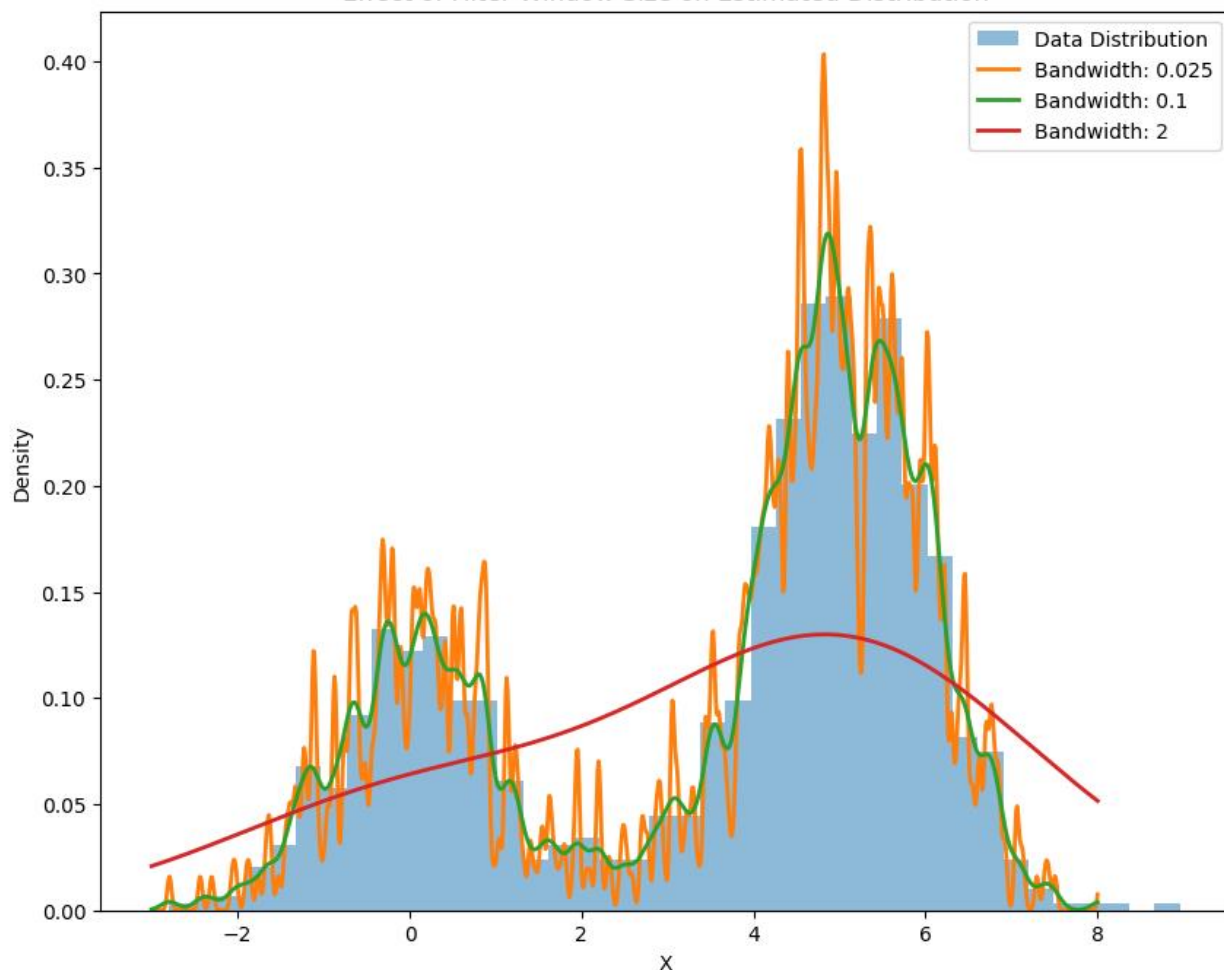


بعد از نمایش آن ما با استفاده از پنجره پارزن توزیع گوسی آن را نمایش می‌دهیم



و در نهایت با h های 2, 0.1, 0.025 ان را رسم میکنیم

Effect of Filter Window Size on Estimated Distribution



در زمینه تجزیه و تحلیل داده‌ها و آمار، اندازه پنجره پارزن به پارامتر هموارسازی مربوط می‌شود، که به عنوان پهنای باند در تخمین چگالی هسته‌ای (KDE) شناخته می‌شود. KDE یک روش غیر پارامتریک برای تخمین تابع چگالی احتمال یک متغیر تصادفی است.

نمودار شامل یک توزیع داده و سه توزیع تخمینی با استفاده از پهنای باندهای مختلف است. پهنای باند کوچکتر منجر به یک خط حساس‌تر می‌شود که داده‌ها را به دقت دنبال می‌کند، که می‌تواند منجر به بیش برآزش و واریانس بالا شود. در مقابل، پهنای باند بزرگتر تخمین را صاف می‌کند، که ممکن است منجر به کم برآزش داده‌ها و نادیده گرفتن ساختار ظریف‌تر آن‌ها شود.

هدف از تنظیم اندازه پنجره فیلتر یافتن تعادلی است که ویژگی‌های مهم توزیع داده‌ها را بدون بیش برآزش یا کم برآزش بگیرد. از نمودار، به نظر می‌رسد که پهنای باند ۰.۱ شکل توزیع داده‌ها را نزدیک‌تر از بقیه دنبال می‌کند، در حالی که پهنای باند ۲ بیش از حد بزرگ است و باعث هموار شدن زیادی می‌شود و پهنای باند ۰.۰۲۵ بیش از حد کوچک است و نویز زیادی را ثبت می‌کند.



این جنبه‌ای حیاتی از تجزیه و تحلیل داده‌ها است، زیرا انتخاب پهنای باند بر نتایجی که می‌توان از داده‌ها استخراج کرد تأثیر می‌گذارد.

پاسخ ۸.

۱-۸ ابتدا دیتاست زیر را با استفاده از قطعه کد زیر ایجاد کنید.

```
from sklearn import cluster, datasets, mixture
noisy_moons=datasets.make_moons(n_samples=500, noise=0.11)
```

۲-۸ یک بار هر کلاس را با توزیع نرمال تقریب بزنید و پارامترهای آن را به دست آورده و کانتورهای مربوطه را رسم نمایید.

۳-۸ این بار از روش GMM استفاده کنید. روش GMM را با تعداد مولفه های ۱ تا ۱۶ تست کنید و شکل داده ها و کانتورها را برای تعداد مولفه برابر با ۳ و ۸ و ۱۶ بدست بیاورید.

۴-۸ تعداد مولفه های بهینه را با توجه به متریک های AIC و BIC به دست بیاورید.

بعد از ساختن دیتای مورد نظر با استفاده از کد زیر توزیع نرمال دیتاست زیر را به دست می آوریم

```
# Calculate mean and covariance for each class
mean_0 = np.mean(class_0, axis=0)
cov_0 = np.cov(class_0, rowvar=False)

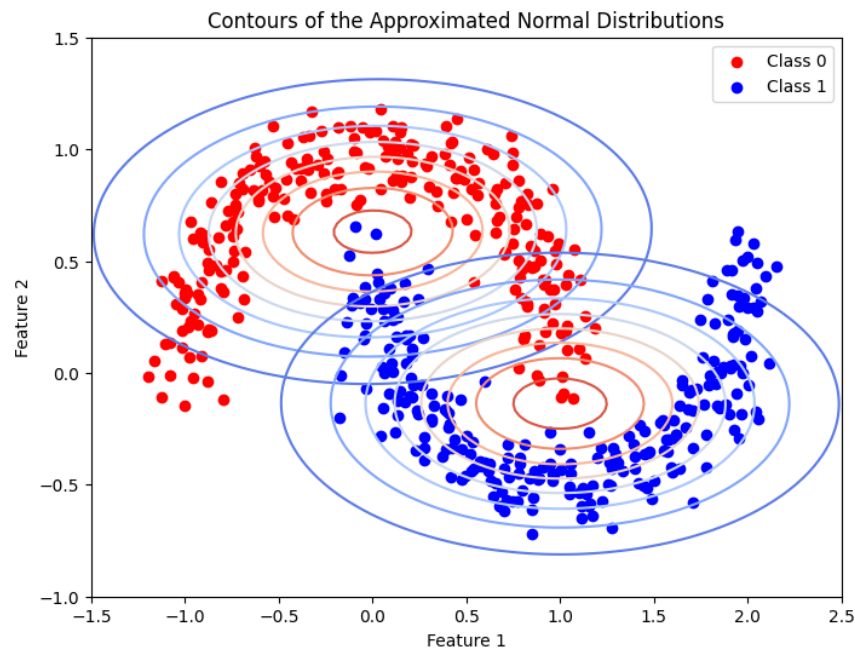
mean_1 = np.mean(class_1, axis=0)
cov_1 = np.cov(class_1, rowvar=False)

# Function to calculate the Gaussian density
def gaussian_density(x, mean, cov):
    n = mean.shape[0]
    diff = (x - mean).reshape(-1, 1)
    return np.exp(-0.5 * np.dot(np.dot(diff.T, np.linalg.inv(cov)), diff)) / \
        ((2 * np.pi)**(n/2) * np.sqrt(np.linalg.det(cov)))

# Creating a grid of points to plot
x, y = np.linspace(-1.5, 2.5, 100), np.linspace(-1, 1.5, 100)
X, Y = np.meshgrid(x, y)
Z0 = np.zeros(X.shape)
Z1 = np.zeros(X.shape)

# Calculating the density for each point in the grid
for i in range(X.shape[0]):
    for j in range(X.shape[1]):
        Z0[i, j] = gaussian_density(np.array([X[i, j], Y[i, j]]), mean_0, cov_0)
        Z1[i, j] = gaussian_density(np.array([X[i, j], Y[i, j]]), mean_1, cov_1)
```

Contour plot هر کلاس که به شکل جدا جدا رسم شده است به صورت زیر می باشد:



در محله ی بعدی از GMM استفاده میشود

مدل مخلوط گاوسی (GMM) یک روش آماری برای مدل سازی توزیع های احتمالاتی پیچیده است که می تواند به عنوان ترکیبی از چندین توزیع نرمال ساده تر در نظر گرفته شوند. در اینجا کلیدی ترین مفاهیم و کاربردهای GMM را توضیح می دهیم:

۱. ترکیب چند گانه: GMM فرض می کند که داده ها از ترکیبی از چندین توزیع نرمال (گاوسی) تولید شده اند. هر یک از این توزیع های نرمال را می توان به عنوان یک "جزء" یا "گروه" در نظر گرفت.

۲. پارامترها: هر جزء دارای سه پارامتر اصلی است: میانگین (که مرکز توزیع را مشخص می کند)، کوواریانس (که شکل و گستردگی توزیع را تعیین می کند)، و وزن (که نشان دهنده احتمال وجود داده ها در آن جزء است).

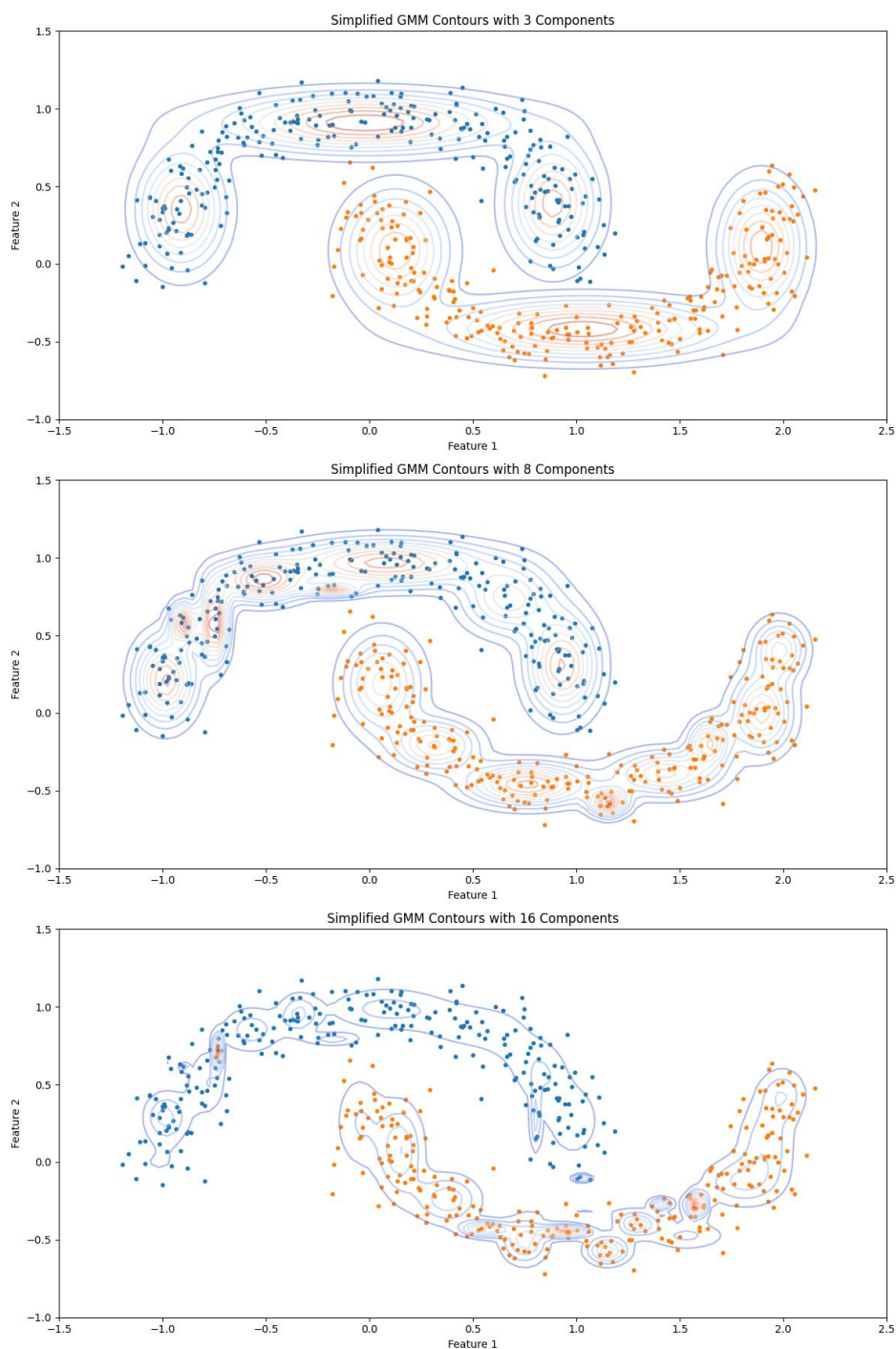
۳. یادگیری پارامترها: برای یادگیری پارامترهای جزءها، معمولاً از الگوریتم Expectation-Maximization (EM) استفاده می شود. این الگوریتم به صورت تکراری احتمال تعلق هر نقطه داده به هر جزء را برآورد می کند و سپس پارامترهای جزءها را براساس این احتمالات به روزرسانی می کند.

۴. کاربردها: GMM در بسیاری از زمینه هایی که نیاز به مدل سازی توزیع های پیچیده داده ها وجود دارد، کاربرد دارد. این شامل شناسایی الگو، تجزیه و تحلیل خوشه ای، کاهش بُعد و سایر کاربردهای مرتبط با یادگیری ماشین است.



۵. مزایا و معایب: مزیت اصلی GMM در انعطاف‌پذیری آن در مدل‌سازی توزیع‌های پیچیده است. با این حال، انتخاب تعداد صحیح جزءها و تعیین پارامترهای اولیه می‌تواند چالش‌برانگیز باشد. همچنین، GMM ممکن است در مواجهه با داده‌های دارای ابعاد بالا یا داده‌هایی که به خوبی از هم جدا نشده‌اند، با مشکلاتی روبرو شود.

توجه: باتوجه به این که در قسمت قبل برای هر کلاس باید بکشیم و این قسمت هم گفته شده که حال این کار را GMM انجام بدید پس ما برای هر کلاس انجام میدیم GMM را.





و در نهایت با استفاده از متود های BIC و AIC تعداد بهینه را مشخص میکنیم

AIC (Akaike Information Criterion)

AIC یک معیار برای ارزیابی کیفیت مدل های آماری است که هم کیفیت برازش داده ها و هم تعداد پارامترهای مدل را در نظر می گیرد. AIC بر اساس فرمول زیر محاسبه می شود:

$$AIC = 2k - 2 \ln(L)$$

که در آن k تعداد پارامترهای مدل است و L احتمال بیشینه (maximum likelihood) است. AIC سعی دارد تعادلی بین پیچیدگی مدل (تعداد پارامترها) و کیفیت برازش مدل به داده ها ایجاد کند.

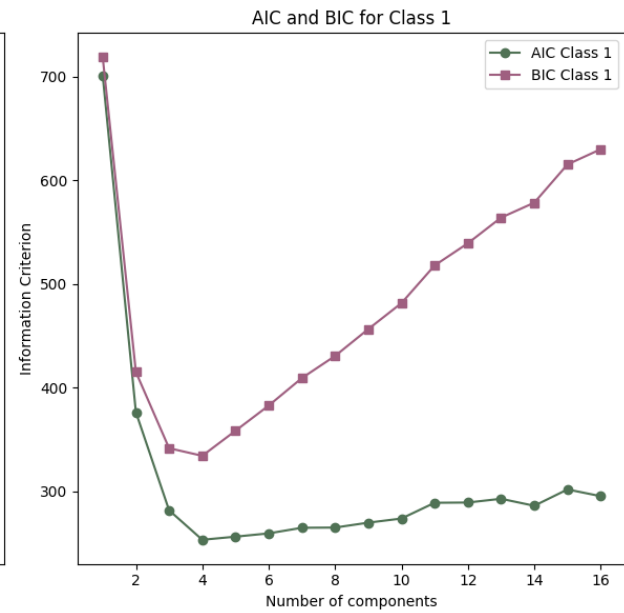
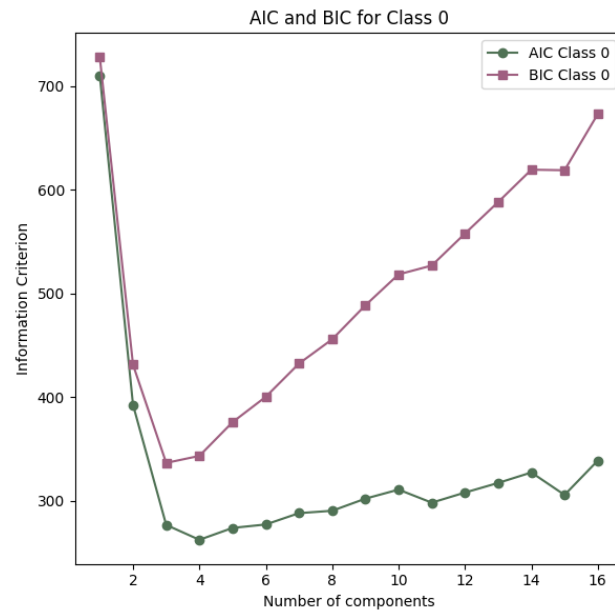
BIC (Bayesian Information Criterion)

BIC، که گاهی به عنوان معیار اطلاعات شوارتز (Schwarz Criterion) نیز شناخته می شود، یک شاخص برای انتخاب مدل است که هم برازش داده ها و هم اندازه نمونه و تعداد پارامترهای مدل را در نظر می گیرد. BIC بر اساس فرمول زیر محاسبه می شود:

$$BIC = \ln(n)k - 2 \ln(L)$$

که در آن n اندازه نمونه، k تعداد پارامترهای مدل و L احتمال بیشینه است.

مانند AIC، BIC نیز به دنبال یافتن تعادل بین پیچیدگی مدل و برازش داده ها است. با این حال، BIC به اندازه نمونه نیز توجه دارد و بنابراین در موقعیتهایی با نمونه های بزرگ، به شدت به پیچیدگی مدل حساس تر است.



تعداد بهینه کلاس ها هم به شرح زیر میباشد

```
Optimal number of components based on AIC for Class zero: 4
Optimal number of components based on BIC for Class zero: 3
Optimal number of components based on AIC for Class one: 4
Optimal number of components based on BIC for Class one: 4
```

پی نوشت : کد هر سه سوال اخر ضمیمه شده است.