


سوال ۱

الف) درست، اضافه کردن ترم منظم ساز باعث ساده تر شدن مدل می شود، لذا خطای بایس احتراس می باید از طرفی، مدل های ساده تر دارای خطای واریانس کمتری هستند.

ب) غلط، ریس کاهش نرادان در صورت همگرایی، به مینیمم محلی می رسد، لذا اگر تابع convex نباشد لزوماً به global min نمی رسد.

ج) غلط، میزان حتم در ساختار داده ای KD-tree بستگی به پراکندگی داده ها دارد، لذا هزینه جستجو ثابت نیست.

د) درست، استفاده از پارامترهای مشترک باعث کاهش تعداد کل پارامترها می شود.

ه) درست، با استفاده از شبکه RBF می توان تابعی به شکل  را تقریب زد و لذا از آن رهم در آموزش نوروهای RBF می توان حر تابعی را تقریب زد. به عبارت دیگر، شبکه های RBF جزء $\text{universal func. approx}$ هستند.

و) غلط، از آنجایی که شبکه ها عصبی $\text{universal func. approx}$ هستند، خطای بایس کمی دارند.

سوال ۲
مرز دو کلاس در تقاطعی است که احتمال پسین دو کلاس سادی باشد:

$$P(y=1|x) = P(y=2|x)$$

$$\Rightarrow \frac{P(x|y=1) \cancel{P(y=1)}}{\cancel{P(x)}} = \frac{P(x|y=2) \cancel{P(y=2)}}{\cancel{P(x)}} \Rightarrow P(x|y=1) = P(x|y=2)$$

$$\Rightarrow \frac{1}{2\pi \sqrt{\frac{1}{9 \times 8}}} \exp \left\{ -\frac{1}{2} [x_1 \ x_2] \begin{bmatrix} 9 & 0 \\ 0 & 8 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\} = \frac{1}{2\pi \sqrt{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} [x_1 \ x_2] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\}$$

$$\xRightarrow{\log} +\log \sqrt{72} - \frac{9}{2} x_1^2 - 4 x_2^2 = +\log \sqrt{2} - x_1^2 - \frac{1}{2} x_2^2 \Rightarrow \frac{7}{2} x_1^2 + \frac{7}{2} x_2^2 = \log \sqrt{\frac{72}{2}}$$

$$\Rightarrow x_1^2 + x_2^2 = \frac{2}{7} \log 6 \rightarrow \text{مرز جدا کننده یک دایره به شعاع } \sqrt{\frac{2}{7} \log 6}$$

الف) $\hat{\theta}_{ML} = \arg \max_{\theta} \log P(D|\theta) = \arg \max_{\theta} \log P(x_1, \dots, x_n | \theta) \stackrel{i.i.d}{=} \arg \max_{\theta} \log \prod_{i=1}^n P(x_i | \theta)$ سوال ۳

$$= \arg \max_{\theta} \sum_{i=1}^n \log P(x_i | \theta) = \arg \max_{\theta} \sum_{i=1}^n (\log \theta + \theta \log b - (\theta + 1) \log x_i)$$

$$= \arg \max_{\theta} \underbrace{n \log \theta + n \theta \log b - (\theta + 1) \sum_{i=1}^n \log x_i}_{L(\theta)}$$

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{n}{\theta} + n \log b - \sum_{i=1}^n \log x_i = 0 \Rightarrow \frac{n}{\theta} = \sum_{i=1}^n \log x_i - n \log b$$

$$\Rightarrow \boxed{\hat{\theta}_{ML} = \frac{n}{\sum_{i=1}^n \log x_i - n \log b}}$$

4)

$$P(\theta|D) \propto P(D|\theta) P(\theta) \stackrel{i.i.d.}{=} \prod_{i=1}^n P(x_i|\theta) P(\theta) = \prod_{i=1}^n \left(\frac{\theta b^\theta}{x_i^{\theta+1}} \right) c \theta^{\alpha-1} e^{-\beta\theta}$$

$$= \frac{\theta^n b^{n\theta}}{\left(\prod_{i=1}^n x_i \right)^{\theta+1}} c \theta^{\alpha-1} e^{-\beta\theta} = c \frac{\theta^{n+\alpha-1} e^{n\theta \ln b} e^{-\beta\theta}}{e^{(\theta+1) \ln \prod_{i=1}^n x_i}}$$

$$= \frac{c}{\underbrace{e^{\sum \ln x_i}}_{c'}} \theta^{n+\alpha-1} e^{-\theta(\beta - n \ln b + \sum_{i=1}^n \ln x_i)}$$

$$= \text{Gamma}(\theta \mid \underbrace{n+\alpha}_{\alpha_{\text{new}}}, \underbrace{\beta - n \ln b + \sum_{i=1}^n \ln x_i}_{\beta_{\text{new}}})$$

ج) بلہ، همان طور کہ در سمت قبل می بینید، توزیع احتمال پس از همان نوع توزیع امکان پسین سده است.

(د)

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta | D) = \frac{\alpha_{new} - 1}{\beta_{new}} = \frac{n + \alpha - 1}{\beta - n \log b + \sum_{i=1}^n \log x_i}$$

ه) بلہ، زیرا:

$$\lim_{n \rightarrow \infty} \hat{\theta}_{MAP} = \lim_{n \rightarrow \infty} \frac{n + \alpha - 1}{\beta - n \log b + \sum_{i=1}^n \log x_i} = \lim_{n \rightarrow \infty} \frac{n}{-n \log b + \sum_{i=1}^n \log x_i} = \lim_{n \rightarrow \infty} \hat{\theta}_{ML}$$

الف) $y_i = w^T x_i + \epsilon_i \Rightarrow P(y_i | x_i) = \text{laplace}(y_i | w^T x_i, 1) = \frac{1}{2} e^{-|y_i - w^T x_i|}$ سوال ٤

$$\hat{w}_{ML} = \arg \max_w \log P(D|w) = \arg \max_w \log \prod_{i=1}^n P(y_i | x_i) = \arg \max_w \sum_{i=1}^n \log P(y_i | x_i)$$

$$= \arg \max_w \sum_{i=1}^n \left(\underbrace{-\log 2}_{\text{const}} - |y_i - w^T x_i| \right) = \arg \max_w - \sum_{i=1}^n |y_i - w^T x_i|$$

$$= \arg \min_w \sum_{i=1}^n |y_i - w^T x_i|$$

ب)

$$\hat{w}_{MAP} = \arg \max_w \overset{\log}{P(w|D)} = \arg \max_w \overset{\log}{P(D|w) P(w)} = \arg \max_w \underbrace{\log P(D|w)}_{\text{رسمت قبل می باشد}} + \log P(w)$$

$$= \arg \max_w - \sum_{i=1}^n \underbrace{(\log 2 - |y_i - w^T x_i|)}_{\text{const with respect to } w} - \underbrace{d \log 2 b}_{\text{const with respect to } w} - \frac{\|w\|_1}{b}$$

$$= \arg \min_w \sum_{i=1}^n |y_i - w^T x_i| + \frac{\|w\|_1}{b}$$

لذا ۱ و ۲ رابطه عکس دارند.

1. [4 points] What is the entropy $H(\text{Passed})$?

★ ANSWER:

$$H(\text{Passed}) = -\left(\frac{2}{6} \log_2 \frac{2}{6} + \frac{4}{6} \log_2 \frac{4}{6}\right)$$

$$H(\text{Passed}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right)$$

$$H(\text{Passed}) = \boxed{\log_2 3 - \frac{2}{3} \approx 0.92}$$

2. [4 points] What is the entropy $H(\text{Passed} \mid \text{GPA})$?

★ ANSWER:

$$H(\text{Passed} \mid \text{GPA}) = -\frac{1}{3} \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) - \frac{1}{3} \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) - \frac{1}{3} (1 \log_2 1)$$

$$H(\text{Passed} \mid \text{GPA}) = \frac{1}{3}(1) + \frac{1}{3}(1) + \frac{1}{3}(0)$$

$$H(\text{Passed} \mid \text{GPA}) = \boxed{\frac{2}{3} \approx 0.66}$$

★ ANSWER:

$$H(Passed|Studied) = -\frac{1}{2}(\frac{1}{3}\log_2 \frac{1}{3} + \frac{2}{3}\log_2 \frac{2}{3}) - \frac{1}{2}(1\log_2 1)$$

$$H(Passed|Studied) = \frac{1}{2}(\log_2 3 - \frac{2}{3})$$

$$H(Passed|Studied) = \frac{1}{2}\log_2 3 - \frac{1}{3} \approx 0.46$$

4. [4 points] Draw the full decision tree that would be learned for this dataset. You do not need to show any calculations.

★ ANSWER: We want to split first on the variable which maximizes the information gain $H(Passed) - H(Passed|A)$. This is equivalent to minimizing $H(Passed|A)$, so we should split on "Studied?" first.

