

۱-

(الف)

$$H(y) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

(ب)

$$P(x_1=T) = \frac{6}{10}, \quad P(x_1=I) = \frac{4}{10}, \quad P(y=w|x_1=T) = \frac{1}{3}, \quad P(y=w|x_1=I) = \frac{3}{4}$$

$$IG(y, x_1) = 1 - \left(\frac{6}{10} \times -\left(\frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3} \right) + \frac{4}{10} \times -\left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right) \right) \approx 0,12$$

$$P(x_2=M) = \frac{4}{10}, \quad P(x_2=P) = \frac{6}{10}, \quad P(y=w|x_2=M) = \frac{3}{4}, \quad P(y=w|x_2=P) = \frac{1}{3}$$

$$IG(y, x_2) = 1 - \left(\frac{6}{10} \times -\left(\frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3} \right) + \frac{4}{10} \times -\left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right) \right) \approx 0,12$$

$$P(x_3=S) = \frac{5}{10}, \quad P(x_3=C) = \frac{5}{10}, \quad P(y=w|x_3=S) = \frac{4}{5}, \quad P(y=w|x_3=C) = \frac{1}{5}$$

$$IG(y, x_3) = 1 - \left(\frac{5}{10} \times -\left(\frac{4}{5} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5} \right) + \frac{5}{10} \times -\left(\frac{1}{5} \log \frac{1}{5} + \frac{4}{5} \log \frac{4}{5} \right) \right) \approx 0,28$$

این معیار باید ویژگی x_3 انتخاب شود چون Information Gain بیشتری دارد.

(الف)

$$MSE = \frac{1}{n} \sum_{i=1}^n (M_T(x_i) - y_i)^2 \xrightarrow{M_T(x_i) = \sum_{t=1}^T \beta_t G_t(x_i)} MSE = \frac{1}{n} \sum_{i=1}^n \left(\sum_{t=1}^T \beta_t G_t(x_i) - y_i \right)^2$$

(ب)

$$MSE = \frac{1}{n} \sum_{i=1}^n \left(\beta_T G_T(x_i) + \sum_{t=1}^{T-1} \beta_t G_t(x_i) - y_i \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\beta_T G_T(x_i) + M_{T-1}(x_i) - y_i \right)^2$$

$$\frac{\partial MSE}{\partial \beta_T} = \frac{1}{n} \sum_{i=1}^n 2 \left(\beta_T G_T(x_i) + M_{T-1}(x_i) - y_i \right) \times G_T(x_i) = 0$$

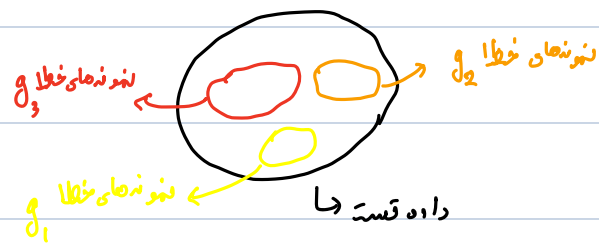
$$\Rightarrow \beta_T \sum_{i=1}^n G_T^2(x_i) = \sum_{i=1}^n (y_i G_T(x_i) - G_T(x_i) M_{T-1}(x_i))$$

$$\Rightarrow \beta_T = \frac{\sum_{i=1}^n (y_i G_T(x_i) - G_T(x_i) M_{T-1}(x_i))}{\sum_{i=1}^n G_T^2(x_i)} = \frac{\sum_{i=1}^n G_T(x_i) (y_i - M_{T-1}(x_i))}{\sum_{i=1}^n G_T^2(x_i)}$$

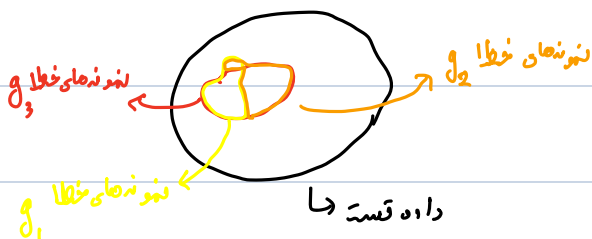
$$P^N = \left(1 - \frac{1}{N}\right)^{PN} = \text{احتمال انتخاب نشدن یک نمونه خاص}$$

$$\Rightarrow \text{چون } N \text{ خیلی بزرگ است} \quad N \times \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^{PN} = N \times e^{-P}$$

ب) ابتدا *lower bound* را بدست می آوریم. برای این کار می توان در نظر گرفت هر کدام از درخت های تقسیم ما ادوی داده های متناوبی اشتباه دسته بندی کنند، به عبارت دیگر هرگاه نمونه ای توسط یکی از درخت ها اشتباه دسته بندی می شود، توسط دو درخت دیگر به درستی دسته بندی می شود. چون جمع خطای 3 درخت از 1 کمتر است ($0.15 + 0.25 + 0.35 = 0.75$)، این حالت ممکن است و در این حالت خطای مدل کلی ما یعنی G بر روی داده تست 0 می شود.



حال *upper bound* را بدست می آوریم. برای این کار باید تعداد و بخشی از داده تست که حداقل توسط 2 تا از درخت ها به اشتباه دسته بندی شوند را پیشینه کنیم. برای این کار به درختی نگاه می کنیم که بیشترین خطا را دارد یعنی g_3 که این میزان 0.35 است. حال این 0.35 را به در قسمت با اندازه های 2.25، 2.25 تقسیم می کنیم. حال اگر فرض کنیم تمام قسمت اول توسط g_2 نیز اشتباه دسته بندی می شوند و تمام قسمت دوم نیز توسط g_1 اشتباه دسته بندی می شوند، همچنین 0.25 از هر کدام از g_1 و g_2 باقی مانده این بخش از این دو نیز با یکدیگر *overlap* دارند. مدل G ما بعد این 0.375 از داده تست را اشتباه دسته بندی می کند بنابراین در این حالت خطای مدل 0.375 خواهد بود.



$$0 \leq E_{out}(G) \leq 0.375$$

الف) درست است. در bagging، ما داده‌های آموزشی خود را از داده‌های آموزشی اصلی sample می‌کنیم، بنابراین بعد از این مرحله sample کردن و بدست آوردن چند دسته داده آموزشی sample شده، یادگیرنده‌ها را می‌توان به شکل جداگانه و موازی بر روی هر کدام از این دسته‌ها آموزش داد.

ب) نادرست است. در روش boosting، یادگیرنده‌ها از هم مستقل نیستند و آموزش هر یادگیرنده به یادگیرنده‌های قبلی بستگی دارد، بنابراین نمی‌توان آن‌ها را به شکل موازی آموزش داد.

ج) نادرست است. در روش bagging، ما برای هر یادگیرنده داده‌های آموزشی را از داده‌های آموزشی اصلی sample می‌کنیم (با replacement). بنابراین ممکن است از کل داده آموزشی اصلی استفاده نکنیم.

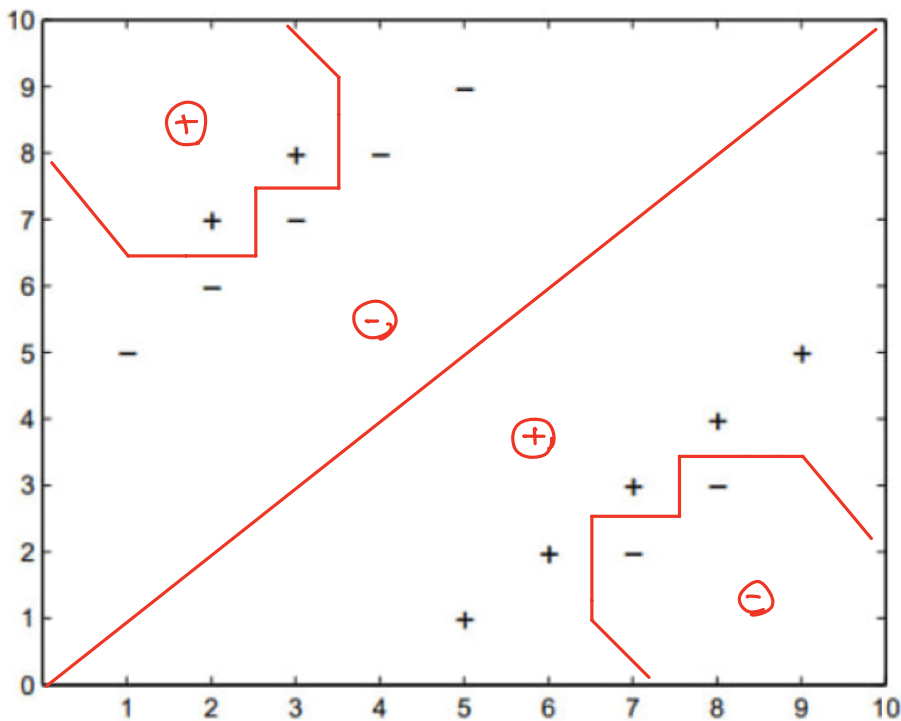
د) درست است. در روش boosting، هر یادگیرنده پس از قبلی، بر روی کل داده‌های آموزشی، آموزش داده می‌شود.

5-

الف) چندی هر کدام از این نمونه ها در داده آموزش نیز وجود دارند و می توانند نزدیکترین همسایه مفرد باشند، اگر $k=1$ باشد، error کمینه شده برابر با صفر می شود.

ب) اگر k خیلی بزرگ باشد، Resolution مدل ماکم شود و بعداً نظر که مشخص است، اگر $k \leq 13$ باشد، همه نمونه ها اشتباه دسته بندی می شوند (با استفاده از روش leave one out cross validation). اما اگر k خیلی کوچک باشد منجر به overfit شدن مدل می شود.

ج) اگر مقدار k را برابر با 5 یا 7 قرار دهیم میزان error مدل کمینه می شود و مقدار error نیز برابر با $\frac{4}{14}$ می شود.



(>

$$H(T) = - \sum_{i=1}^k P(c_i) \log P(c_i) = - \sum_{i=1}^k \frac{1}{k} \log \frac{1}{k} = \frac{1}{k} \sum_{i=1}^k \log k = \log k$$

چون ویژگی A و T را به m_A دسته disjoint تقسیم می‌کند و توزیع کلاس‌ها در دسته به شکل Uniform است:

$$H(T_i^A) = H(T) = \log k, \quad P(A^i) = \frac{|T_i^A|}{|T|}$$

$$\Rightarrow H(T|A) = \sum_{i=1}^{m_A} P(A^i) \underbrace{H(T_i^A)}_{\log k} = \log k \underbrace{\sum_{i=1}^{m_A} P(A^i)}_1 = \log k$$

$$IG(T, A) = H(T) - H(T|A) = \log k - \log k = 0$$

به عبارت دیگر آر ویژگی A را بدانی، اطلاعاتی به ما در مورد کلاس نمونه‌های T نمی‌دهد و عدم قطعیت ما را کم نمی‌کند.

یا می‌توان گفت کلاس نمونه‌های موجود در T مستقل از ویژگی A هستند.