

Theory Questions

1. (10 points) **Suboptimality of ID3.** Solve exercise 2 in chapter 18 in the course book: Understanding Machine Learning: From Theory to Algorithms.
2. **(Suboptimality of ID3)**

Consider the following training set, where $\mathcal{X} = \{0, 1\}^3$ and $\mathcal{Y} = \{0, 1\}$:

1. $((1, 1, 1), 1)$
2. $((1, 0, 0), 1)$
3. $((1, 1, 0), 0)$
4. $((0, 0, 1), 0)$

Suppose we wish to use this training set in order to build a decision tree of depth 2 (i.e., for each input we are allowed to ask two questions of the form $(x_i = 0?)$ before deciding on the label).

1. Suppose we run the ID3 algorithm up to depth 2 (namely, we pick the root node and its children according to the algorithm, but instead of keeping on with the recursion, we stop and pick leaves according to the majority label in each subtree). Assume that the subroutine used to measure the quality of each feature is based on the entropy function (so we measure the *information gain*), and that if two features get the same score, one of them is picked arbitrarily. Show that the training error of the resulting decision tree is at least $1/4$.

לפ' עליה גains '32 gain > ak eopne (גיאנס ak eanf, nlelo)

$$\begin{aligned} Gain(S, i) &= \underbrace{C(\mathbb{P}_S [Y = 1])}_{\text{pre-split error}} - \\ &\quad \underbrace{(\mathbb{P}_S[X_i = 1] C(\mathbb{P}_S[Y = 1|X_i = 1]) + \mathbb{P}_S[X_i = 0] C(\mathbb{P}_S[Y = 1|X_i = 0])))}_{\text{post-split error}} \end{aligned}$$

Entropy, $C(\alpha) = -\frac{1}{2} (\alpha \log \alpha + (1 - \alpha) \log (1 - \alpha))$

$$gain(S, 1) = C\left(\frac{1}{2}\right) - \left[\frac{3}{4} \cdot C\left(\frac{2}{3}\right) + \frac{1}{4} \cdot C\left(\frac{1}{3}\right) \right]$$

$$= -\frac{1}{2} \cdot \underbrace{\left[\frac{1}{2} \cdot \log\left(\frac{1}{2}\right) + \frac{1}{2} \cdot \log\left(\frac{1}{2}\right) \right]}_{-1} - \frac{3}{4} \cdot \left(-\frac{1}{2} \right) \cdot \underbrace{\left[\frac{2}{3} \log\left(\frac{2}{3}\right) + \frac{1}{3} \cdot \log\left(\frac{1}{3}\right) \right]}_{-0.78} \approx 0.21$$

$$= -\frac{1}{2} \cdot \underbrace{\left[\frac{1}{2} \cdot \log(1/2) + \frac{1}{2} \cdot \log(1/2) \right]}_{-1} - \frac{3}{4} \cdot (-\frac{1}{2}) \cdot \underbrace{\left[\frac{4}{3} \log(\frac{2}{3}) + \frac{1}{3} \cdot \log(\frac{1}{3}) \right]}_{-0.78} \approx 0.21$$

$$\text{gain}(S, 2) = \underbrace{C(1/2)}_1 - \left[\frac{1}{2} \cdot \underbrace{C(1/2)}_1 + \frac{1}{2} \cdot \underbrace{C(1/2)}_1 \right] = 0$$

$$\text{gain}(S, 3) = C(1/2) - \left[\frac{1}{3} \cdot C(1/2) + \frac{2}{3} \cdot C(1/3) \right] = 0$$

לפנינו מושג אחד (המוגדר בתרגיל 3, פה מושג שני)
 שיפרנו את המושג הראשון, אך שיפרנו אותו לא יותר מאשר 3, פה מושג שני

בנוסף לכך, ניתן לשים לב כי המושג השני מושג נסיעה מושג אחד.

$\therefore (X_2 = 1 ?)$ מושג אחד (1)

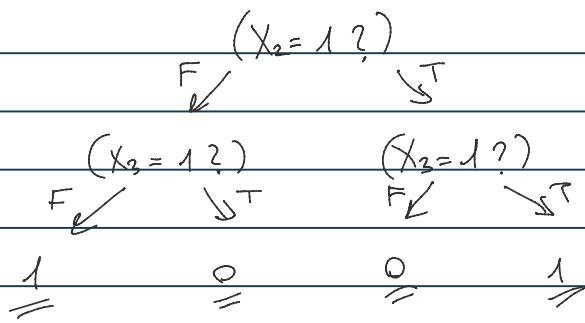
ונובע מכך שמשתנה $y_3 = 1$ או $y_3 = 0$ מושג נסיעה מושג אחד.

$\therefore (X_3 = 1 ?)$ מושג אחד (2)

ונובע מכך שמשתנה $y_2 = 1$ או $y_2 = 0$ מושג נסיעה מושג אחד.

לפנינו מושג אחד וzero, והו מושג שמייד יתאפשר לרשום:

2. Find a decision tree of depth 2 that attains zero training error.



3. (20 points) AdaBoost. Let $x_1, \dots, x_m \in \mathbb{R}^d$ and $y_1, \dots, y_m \in \{-1, 1\}$ its labels. We run the AdaBoost algorithm as given in the lecture, and we are in iteration t . Assume that $\epsilon_t > 0$.

- (a) Show that the error of the current hypothesis relative to the new distribution is exactly $1/2$, that is:

$$\Pr_{x \sim D_{t+1}} [h_t(x) \neq y] = \frac{1}{2}.$$

$$\text{Defining } p_0 \quad \sum_{i: h_t(x_i) \neq y_i} D_{t+1}(i) = \frac{1}{2} \quad \text{as a rough definition}$$

$$\Pr_{\substack{x \sim D_{t+1}}} [h_t(x) \neq y] = \mathbb{E}_{D_{t+1}} \left[\mathbb{I}\{h_t(x) \neq y\} \right] = \sum_{i=1}^m D_{t+1}(i) \cdot \mathbb{I}\{h_t(x_i) \neq y\} = \sum_{i: h_t(x_i) \neq y} D_{t+1}(i)$$

$$= \frac{e^{w_t}}{z_t} \cdot \sum_{i: h_t(x_i) \neq y_i} D_b(i) = e^{w_t} \cdot \underbrace{\frac{(1-\epsilon_t)}{\epsilon_t} \cdot \frac{\epsilon_t}{2 \cdot \sqrt{\epsilon_t(1-\epsilon_t)}}}_{\text{je, je n. g. Koeffizient}} = \frac{(1-\epsilon_t)^{0.5}}{2 \cdot \epsilon_t^{0.5} \cdot (1-\epsilon_t)^{0.5}} = \frac{1}{2}$$

- (b) Show that AdaBoost will not pick the same hypothesis twice consecutively; that is $h_{t+1} \neq h_t$.

ה_t = h_{t+1} מינימיזירן את ה-errort באלגוריתם AdaBoost והיינו מילאנו את הדרישה
 של $\sum_{i=1}^n \hat{y}_i^2 \leq \frac{1}{2}$ ו- $E_{t+1} < \frac{1}{2} \wedge E_t < \frac{1}{2}$ ו- $\hat{y}_i^2 \leq 1$ ו- $\hat{y}_i \in \{-1, 1\}$
 מנו. $\frac{1}{2}$ ל- D_{t+1} מינימיזירן ht ו- \hat{y}_i מ- y_i ו- $\hat{y}_i^2 \leq 1$ ו- $\hat{y}_i \in \{-1, 1\}$

4. (20 points) Sufficient Condition for Weak Learnability. Let $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be a training set and let \mathcal{H} be a hypothesis class. Assume that there exists $\gamma > 0$, hypotheses $h_1, \dots, h_k \in \mathcal{H}$ and coefficients $a_1, \dots, a_k \geq 0$, $\sum_{i=1}^k a_i = 1$ for which the following holds:

$$y_i \sum_{j=1}^k a_j h_j(x_i) \geq \gamma \quad (1)$$

for all $(x_i, y_i) \in S$.

- (a) Show that for any distribution D over S there exists $1 \leq j \leq k$ such that

$$\Pr_{i \sim D} [h_j(x_i) \neq y_i] \leq \frac{1}{2} - \frac{\gamma}{2}.$$

(Hint: Take expectation of both sides of inequality (1) with respect to D .)

Remark: Note that the condition above is sufficient for *empirical* weak learnability, the condition defined in lecture #9 for the Adaboost analysis.

$$\Pr_{i \sim D} [h_j(x_i) \neq y_i] > \frac{1 - \gamma}{2} \quad \text{: e.g. } 1 \leq j \leq k \text{ für } \forall i$$

$$\gamma = E[\gamma] \leq E \left[y_i \cdot \sum_{j=1}^k a_j h_j(x_i) \right] = \sum_{i=1}^n \Pr_{X, Y \sim D} [X = x_i, Y = y_i] \cdot y_i \cdot \sum_{j=1}^k a_j h_j(x_i)$$

$$\gamma = \sum_{i=1}^n D(i) y_i \sum_{j=1}^k a_j h_j(x_i) = \sum_{j=1}^k a_j \sum_{i=1}^n D(i) y_i h_j(x_i) = *$$

$$\text{: if } -1 = y_i \cdot h_j(x_i) \text{ e.g. } h_j(x_i) \neq y_i \text{ relabel } 1 = y_i \cdot h_j(x_i) \text{ e.g. } h_j(x_i) = y_i \text{ relabel } y_i$$

$$* = \sum_{j=1}^k a_j \left[\sum_{i: h_j(x_i) = y_i} D(i) - \sum_{i: h_j(x_i) \neq y_i} D(i) \right] = \text{relabel}$$

$$= \sum_{j=1}^k a_j \left(\Pr_{i \sim D} [h_j(x_i) = y_i] - \Pr_{i \sim D} [h_j(x_i) \neq y_i] \right) =$$

$$= \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) = y_i] - \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] =$$

$$= \sum_{j=1}^k a_j (1 - \Pr_{i \sim D} [h_j(x_i) \neq y_i]) - \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] =$$

$$= \sum_{j=1}^k a_j - 2 \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] = 1 - 2 \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i]$$

$$\gamma \leq 1 - 2 \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] \quad \text{: if } \gamma > 0$$

$$\Leftrightarrow \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] \leq \frac{1 - \gamma}{2}$$

$$\frac{1 - \gamma}{2} = \sum_{j=1}^k a_j \left(\frac{1 - \gamma}{2} \right) \leq \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] \text{ after expand (1)}$$

$$\frac{1-\gamma}{2} = \sum_{j=1}^k a_j \cdot \left(\frac{1-\gamma}{2}\right) \leq \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] \text{ if } \gamma < \frac{1}{2}$$

$$\text{then } \frac{1-\gamma}{2} \leq \sum_{j=1}^k a_j \Pr_{i \sim D} [h_j(x_i) \neq y_i] \leq \frac{1-\gamma}{2} \quad \text{by def}$$

- (b) Let $S = \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathbb{R}^d \times \{-1, 1\}$ be a training set that is realized by a d -dimensional hyper-rectangle classifier, i.e., there exists a d dimensional hyper-rectangle $[b_1, c_1] \times \dots \times [b_d, c_d]$ which contains all of the positive points in S and doesn't contain the negative points in S . Let \mathcal{H} be the class of decision stumps of the form

$$h(x) = \begin{cases} 1 & x_j \leq \theta \\ -1 & x_j > \theta \end{cases}, \quad h(x) = \begin{cases} 1 & x_j \geq \theta \\ -1 & x_j < \theta \end{cases},$$

for $1 \leq j \leq d$ and $\theta \in \mathbb{R} \cup \{\infty, -\infty\}$ (for $\theta \in \{\infty, -\infty\}$ we get constant hypotheses which predict always 1 or always -1). Show that there exist $\gamma > 0$, $k > 0$, hypotheses $h_1, \dots, h_k \in \mathcal{H}$ and $a_1, \dots, a_k \geq 0$ with $\sum_{i=1}^k a_i = 1$, such that the condition in inequality (1) holds for the training set S and hypothesis class \mathcal{H} . This implies that \mathcal{H} is empirically weak learnable w.r.t. data realizable by a d -dimensional hyper-rectangle.

(Hint: Set $k = 4d - 1$, $a_i = \frac{1}{4d-1}$ and let $2d - 1$ of the hypotheses be constant.)

$$\forall i \in [k]: a_i = \frac{1}{4d-1}, \quad k = 4d-1 \quad (\text{why})$$

$$h_{b_j}(x) = \begin{cases} 1 & x_j \geq b_j \\ -1 & \text{else} \end{cases}, \quad h_{c_j}(x) = \begin{cases} 1 & x_j \leq c_j \\ -1 & \text{else} \end{cases}$$

↑
($j \in [d]$)'וְנִזְמָנָה
↑
 d

היפוך גיאומטריה של -1 בפונקציית $y = 1/x$ (בזווית 90°)

$$\forall j \in [d]: b_j \leq x_j \leq c_j \quad \text{and} \quad \text{היפוך גיאומטריה של } y = 1/x \text{ בזווית } 90^\circ$$

$$y \cdot \sum_{j=1}^k a_j h_j(x) = \frac{1}{4d-1} \cdot \sum_{j=1}^k h_j(x) = \frac{1}{4d-1} \cdot (2d - (2d-1)) = \frac{1}{4d-1}$$

$\left[\begin{array}{l} \text{היפוך גיאומטריה של } y = 1/x \text{ בזווית } 90^\circ \\ \text{היפוך גיאומטריה של } y = 1/x \text{ בזווית } 90^\circ \end{array} \right] \left[\begin{array}{l} \text{היפוך גיאומטריה של } y = 1/x \text{ בזווית } 90^\circ \\ \text{היפוך גיאומטריה של } y = 1/x \text{ בזווית } 90^\circ \end{array} \right]$

לעומת נסחאות אחרות מתקיימת:

$\exists j \in [d]: x_j < b_j \vee x_j > c_j$ ו- $y = -1$ פולש.

$$y \cdot \sum_{j=1}^k a_j h_j(x) = -\frac{1}{4d-1} \cdot \sum_{j=1}^k h_j(x) = -\frac{1}{4d-1} \left(-(2d-1) + \sum_{j=1}^d h_{b_j}(x) + \sum_{j=1}^d h_{c_j}(x) \right)$$

$\underbrace{\sum_{j=1}^d h_{b_j}(x) + \sum_{j=1}^d h_{c_j}(x)}_{\text{הסכום של הנקודות שפוגעות}}$

$$= \frac{2d-1}{4d-1} - \frac{1}{4d-1} \cdot \left(\sum_{j=1}^d h_{b_j}(x) + \sum_{j=1}^d h_{c_j}(x) \right)$$

בנוסף ל- $x_j < b_j$ או $x_j > c_j$ מתקיימת לפחות אחת נוספת. לכן $x_j < b_1, \dots, b_d, x_j > c_1, \dots, c_d$.

$$\leq 2d-1 - 1 = 2d-2$$

$$\frac{2d-1}{4d-1} - \frac{1}{4d-1} \left(\sum_{j=1}^d h_{b_j}(x) + \sum_{j=1}^d h_{c_j}(x) \right) \geq \frac{2d-1}{4d-1} - \frac{1}{4d-1} \cdot (2d-2) = \frac{1}{4d-1}$$

$$y \cdot \sum_{j=1}^k a_j h_j(x) = \frac{1}{4d-1} : y = 1 \text{ או } y = -1 \text{ יתגלו}$$

$$y \cdot \sum_{j=1}^k a_j h_j(x) \geq \frac{1}{4d-1} : y = -1 \text{ או } y = 1$$

■ $\therefore \text{השאלה שפוגעת ב-} \frac{1}{4d-1} \text{ היא שפוגעת ב-} \frac{1}{4d-1}$

2. (20 points) Properties of KL divergence. Recall the definition of the KL-divergence from slide 15 in recitation 8.

- (a) Show that the KL-divergence is always non-negative.

(Hint: Consider the convex function $f(y) = y \log y$, and use Jensen's inequality which states that for any distribution q , $\mathbb{E}_{z \sim q}[f(z)] \geq f(\mathbb{E}_{z \sim q}[z])$.

- (b) Let p_1, p_2, q_1, q_2 be distributions over \mathcal{X} such that p_1 is independent of p_2 and q_1 is independent of q_2 . Denote the product distributions $p = p_1 \times p_2$ and $q = q_1 \times q_2$ over \mathcal{X}^2 (i.e. $p(x_1, x_2) = p_1(x_1)p_2(x_2)$ for any $x_1, x_2 \in \mathcal{X}$ and similarly for q). Prove the following:

$$D_{KL}(p, q) = D_{KL}(p_1, q_1) + D_{KL}(p_2, q_2).$$

直观上讲，如果 $f(z) = -\log z$ 是一个凸函数， $f'(z) = -\frac{1}{z} < 0$

$$\begin{aligned} D_{KL}(p, q) &:= \sum_{x \in \mathcal{X}} p(x) \cdot \log \left(\frac{p(x)}{q(x)} \right) = \sum_{x \in \mathcal{X}} p(x) \cdot \left[-\log \left(\frac{q(x)}{p(x)} \right) \right] \\ &\stackrel{\text{jensen}}{=} \mathbb{E}_p \left[-\log \left(\frac{q(x)}{p(x)} \right) \right] = -\log \left(\mathbb{E}_p \left[\frac{q(x)}{p(x)} \right] \right) = -\log \sum_{x \in \mathcal{X}} p(x) \cdot \frac{q(x)}{p(x)} \\ &= -\log \sum_{x \in \mathcal{X}} q(x) = -\log 1 = 0 \end{aligned}$$

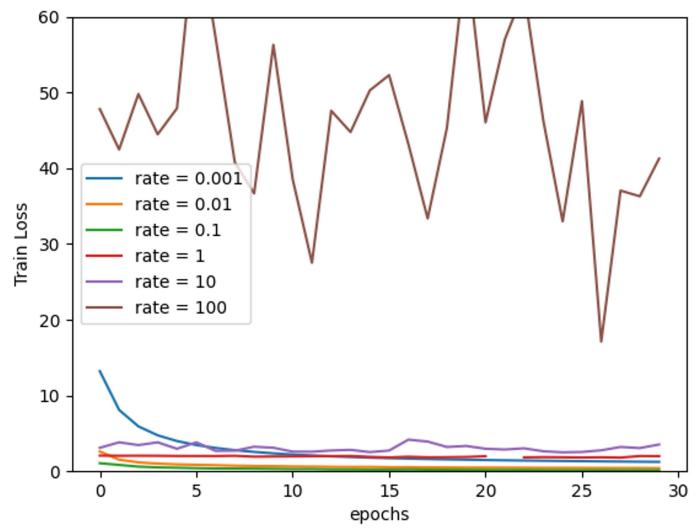
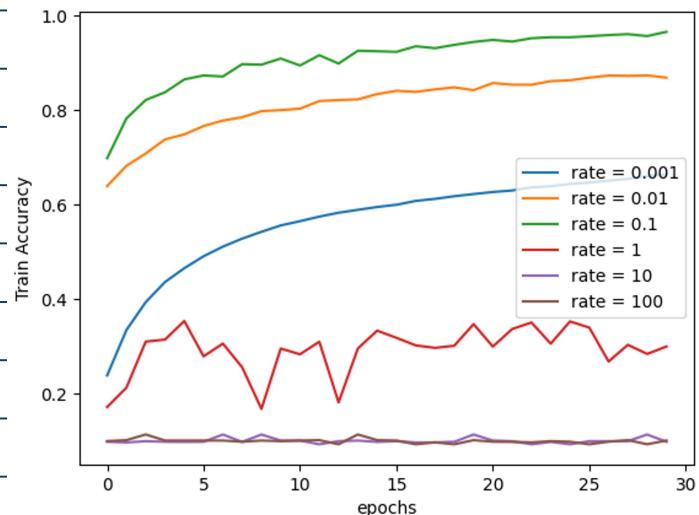
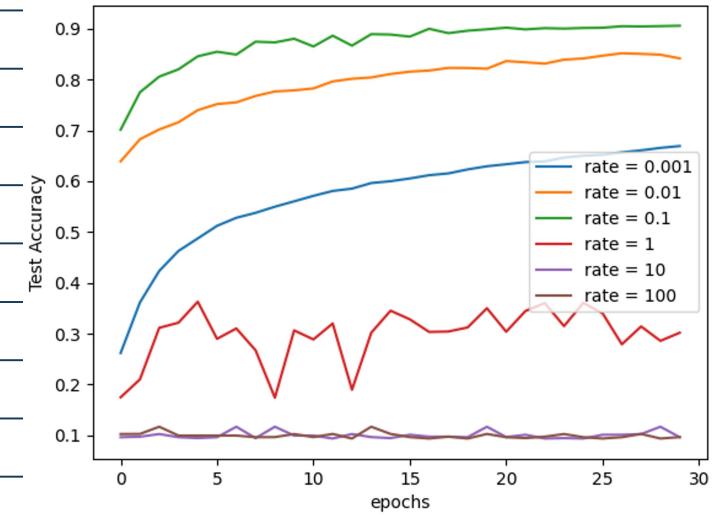
■

$$\begin{aligned} D_{KL}(p, q) &= \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} p(x_1, x_2) \cdot \log \left(\frac{p(x_1, x_2)}{q(x_1, x_2)} \right) = \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} p_1(x_1) \cdot p_2(x_2) \cdot \log \left(\frac{p_1(x_1) \cdot p_2(x_2)}{q_1(x_1) \cdot q_2(x_2)} \right) \quad (b) \\ &= \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} p_1(x_1) \cdot p_2(x_2) \cdot \left[\log \left(\frac{p_1(x_1)}{q_1(x_1)} \right) + \log \left(\frac{p_2(x_2)}{q_2(x_2)} \right) \right] \\ &= \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} p_1(x_1) \cdot p_2(x_2) \cdot \log \left(\frac{p_1(x_1)}{q_1(x_1)} \right) + \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} p_1(x_1) \cdot p_2(x_2) \cdot \log \left(\frac{p_2(x_2)}{q_2(x_2)} \right) \\ &= \sum_{x_1 \in \mathcal{X}} p_1(x_1) \cdot \underbrace{\log \left(\frac{p_1(x_1)}{q_1(x_1)} \right)}_{\text{只计算 } q_1 \text{ 的项}} \cdot \sum_{x_2 \in \mathcal{X}} p_2(x_2) + \sum_{x_2 \in \mathcal{X}} p_2(x_2) \cdot \underbrace{\log \left(\frac{p_2(x_2)}{q_2(x_2)} \right)}_{\text{只计算 } q_2 \text{ 的项}} \cdot \sum_{x_1 \in \mathcal{X}} p_1(x_1) \\ &= D_{KL}(p_1, q_1) + D_{KL}(p_2, q_2) \end{aligned}$$

■

نیکوں کی

3480



finished SGD
Epoch 29: test accuracy: 0.906

lossin "for if we'll be 201-138. If we have 0.2, 0.4, 0.6, 0.8, 1 = 23% > 32%.
So if we have 0.1, it's 100 = 23% > 32%.
SGD in 13/33, 3N, 1/N, 1/2, 1/3, 0.9, 0.1, 0.6, 0.4, 0.2.
The SGD is 13/33, 3N, 1/N, 1/2, 1/3, 0.9, 0.1, 0.6, 0.4, 0.2.

8/10

Epoch 27 test accuracy: 0.9463
Epoch 28 test accuracy: 0.9413
Epoch 29 test accuracy: 0.9433
Finished SGD - test accuracy: 0.9433