### ① Problem 1.

↳ Theorem 1.1: Let $X$ be the space of features and $Y = \{0,1\}$ be the labels space for a binary classification problem. Define the Bayes Classifier $h^*$ as followed:

$$h^*(X) = \begin{cases} 1 & \text{if } p(X) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

Where $p(x) = P(Y = 1 \mid X)$. Then, we have the following properties of $h^*$:

(i) $R(h^*) = R^* \Rightarrow h^*$ is the Bayes Classifier.

(ii) $\forall h: X \to Y$:

$$R(h) - R^* = 2\mathbb{E}_X\left[\left|p(x) - \frac{1}{2}\right| \mathbb{1}(h(x) \neq h^*(x))\right]$$

(iii) $R^* = \mathbb{E}_X\left[\min(p(x), 1 - p(x))\right]$.

↳ Problem: Extend the above theorem to multi-class classification problem. Meaning $Y = \{1, 2, \ldots, M\}$.

→ Solution: Re-define $h^*$ as followed

$$p(x) = \begin{bmatrix} p_1(x) \\ \vdots \\ p_M(x) \end{bmatrix} \quad \text{where } p_c(x) = P(Y = c \mid X = x).$$

and let $h^*(x) = \underset{c}{\arg\max}\{p_c(x)\}$.

$\Rightarrow$ We have that: $\sum_{y \in Y} p_y(x) = 1, \quad \forall x \in X.$

(i) Given an arbitrary classifier $h: X \to Y$, we have:

$$R(h) = \mathbb{E}_X\left[\mathbb{E}_{Y|X}\left[\mathbb{1}(h(x) \neq Y)\right]\right]$$

$$= \mathbb{E}_{x \sim X}\left[\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, P(Y=y \mid X=x)\right] = \mathbb{E}_{x \sim X}\left[\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x)\right]$$

$\Rightarrow$ Let $\hat{y} = h(x)$ for a given $x \in \mathcal{X}$. We have:

$$R(h) = \mathbb{E}_{x \sim X}\left[\sum_{y \in \mathcal{Y};\, y \neq \hat{y}} p_y(x)\right] = \mathbb{E}_{x \sim X}\left[1 - p_{\hat{y}}(x)\right]$$

$\Rightarrow$ To minimize $R(h)$, we need to maximize $p_{\hat{y}}(x)$ for all $x \in \mathcal{X}$. Hence, we have:

$$\hat{y} = h^*(x) = \operatorname{argmax}_c\{p_c(x)\} \Rightarrow h^* = \inf_h\{R(h)\} . \quad \square$$

(iii) From (i) we have:

$$R^* = \mathbb{E}_X\left[1 - \max(\{p_c(x) : c \in \mathcal{Y}\})\right] .$$

$\Rightarrow$ In other words, if we set $\overline{p_c(x)} = P(Y \neq c \mid X=x)$. We have:

$$R^* = \mathbb{E}_X\left[\min(\{\overline{p_c(x)} : c \in \mathcal{Y}\})\right] . \quad \square \qquad \text{(This is similar to the formula for binary classification case).}$$

① <u>Problem 1</u> ((Continued...)

(ii) Compute the excess risk $R(h) - R^*$ :

$$R(h) - R^* = \mathbb{E}_{x \sim X}\left[\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x)\right] - \mathbb{E}_{x \sim X}\left[1 - \max\left(\{p_c(x) : c \in \mathcal{Y}\}\right)\right]$$

$$= \mathbb{E}_{x \sim X}\left[\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x) + \max\left(\{p_c(x) : c \in \mathcal{Y}\}\right) - 1\right]$$

→ When $h(x) = h^*(x)$ ⟹ $\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x) + \max\left(\{p_c(x) : c \in \mathcal{Y}\}\right) = 1$

⟹ $\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x) + \max\left(\{p_c(x) : c \in \mathcal{Y}\}\right) - 1 = 0$

→ When $h(x) \neq h^*(x)$ : Let $\hat{y} = h(x)$ and $\hat{y}_* = h^*(x)$. We have

$$\sum_{y \in \mathcal{Y}} \mathbb{1}(h(x) \neq y)\, p_y(x) + \max\left(\{p_c(x) : c \in \mathcal{Y}\}\right) - 1$$

$$= \sum_{y \neq \hat{y}} p_y(x) + p_{\hat{y}_*}(x) - 1$$

$$= 2 p_{\hat{y}_*}(x) - 1 + \sum_{\substack{y \in \mathcal{Y}; y \neq \hat{y}; \\ y \neq \hat{y}_*}} p_y(x) = 2 p_{\hat{y}_*}(x) - \left(1 - \sum_{\substack{y \in \mathcal{Y}; y \neq \hat{y}; \\ y \neq \hat{y}_*}} p_y(x)\right)$$

$$= 2 p_{\hat{y}_*}(x) - P\left(Y \in \{\hat{y}, \hat{y}_*\} \mid X = x\right)$$

⟹ $R(h) - R^* = \mathbb{E}_{x \sim X}\left[\left(2 p_{\hat{y}_*}(x) - \underbrace{P\left(Y \in \{\hat{y}, \hat{y}_*\} \mid X = x\right)}_{= \;p_{\hat{y}}(x) \,+\, p_{\hat{y}_*}(x).}\right)\mathbb{1}(h(x) \neq h^*(x))\right]$

$$= \mathbb{E}_{x \sim X}\left[\left(p_{\hat{y}_*}(x) - p_{\hat{y}}(x)\right)\mathbb{1}(h(x) \neq h^*(x))\right] . \quad \square$$

## ② Problem 2.

↳ **Problem:** Let $\alpha \in (0,1)$. Define the $\alpha$-cost-sensitive risk of $h : \mathcal{X} \to \mathcal{Y}$ to be:

$$R_\alpha(h) = \mathbb{E}_{XY}\left[(1-\alpha)\,\mathbb{1}(y=1, h(x)=0) + \alpha\,\mathbb{1}(y=0, h(x)=1)\right]$$

→ Determine the Bayes Classifier and prove an analogue of Theorem 1.1 for this risk.

→ **Solution:**

(i) Determine the Bayes Classifier & Bayes Risk:

→ We have: $R_\alpha(h) = \mathbb{E}_{X \sim X}\left[p(x)(1-\alpha)\mathbb{1}(y=1, h(x)=0) + \alpha(1-p(x))\mathbb{1}(y=0, h(x)=1)\right]$

⇒ We have the following table of values for the integrand inside $\mathbb{E}_{X \sim X}[\cdots]$:

| $h(x)$ | | |
|---|---|---|
| 0 | 0 | $p(x)(1-\alpha)$ |
| 1 | $\alpha(1-p(x))$ | 0 |
| | 0 | 1 |

$y$

⇒

For $R_\alpha(h)$ to be minimize, we have:

$$h^*(x) = \begin{cases} 1 & \text{when } \alpha(1-p(x)) \leq p(x)(1-\alpha) \\ 0 & \text{Otherwise.} \end{cases}$$

⇒ The Bayes Classifier is then defined as:

$$h^*(x) = \begin{cases} 1 & \text{when } p(x) \geq \alpha. \\ 0 & \text{Otherwise.} \end{cases}$$

⇒ We have the following Bayes Risk:

$$R_\alpha^* = R_\alpha(h^*) = \mathbb{E}_X\left[\min\big(\alpha(1-p(x)), p(x)(1-\alpha)\big)\right] \quad \square .$$

② **Problem 2** (Continued...)

(ii) Compute the Excess Risk: For all $h: \mathcal{X} \to \mathcal{Y}$, we have

$$R_\alpha(h) = \mathbb{E}_{XY}\left[(1-\alpha)\mathbb{1}\left(Y=1, h(x)=0\right) + \alpha\,\mathbb{1}\left(Y=0, h(x)=1\right)\right]$$

$$= \mathbb{E}_{x \sim X}\left[p(x)(1-\alpha)\mathbb{1}\left(y=1, h(x)=0\right) + (1-p(x))\alpha\,\mathbb{1}\left(y=0, h(x)=1\right)\right]$$

$$= \mathbb{E}_{x \sim X}\left[p(x)(1-\alpha)\mathbb{1}\left(h(x)=0\right) + (1-p(x))\alpha\,\mathbb{1}\left(h(x)=1\right)\right]$$

$$\left(\text{Because } \mathbb{1}\left(h(x)=0, y=0\right) \text{ and } \mathbb{1}\left(h(x)=1, y=1\right) \text{ carries no risk}\right).$$

$$\Rightarrow R_\alpha(h) - R_\alpha^* = \mathbb{E}_{x \sim X}\left[p(x)(1-\alpha)\left[\mathbb{1}\left(h(x)=0\right) - \mathbb{1}\left(h^*(x)=0\right)\right]\right.$$

$$\left. + (1-p(x))\alpha\left[\mathbb{1}(h(x)=1) - \mathbb{1}(h^*(x)=1)\right]\right]$$

$$= \mathbb{E}_{x \sim X}\left[p(x)(1-\alpha)\left[\mathbb{1}\left(h(x)=0, h^*(x)=1\right) - \mathbb{1}\left(h(x)=1, h^*(x)=0\right)\right]\right.$$

$$\left. + (1-p(x))\alpha\left[\mathbb{1}\left(h(x)=1, h^*(x)=0\right) - \mathbb{1}\left(h(x)=0, h^*(x)=1\right)\right]\right]$$

$$= \mathbb{E}_{x \sim X}\left[\mathbb{1}\left(h(x)=0, h^*(x)=1\right)\left[p(x)(1-\alpha) - \alpha(1-p(x))\right]\right.$$

$$\left. + \mathbb{1}\left(h(x)=1, h^*(x)=0\right)\left[\alpha(1-p(x)) - p(x)(1-\alpha)\right]\right]$$

$$= \mathbb{E}_{x \sim X}\left[\mathbb{1}\left(h(x)=0, h^*(x)=1\right)\left(p(x)-\alpha\right)\right.$$

$$\left. + \mathbb{1}\left(h(x)=1, h^*(x)=0\right)\left(\alpha - p(x)\right)\right]$$

$$= \mathbb{E}_X\left[|p(x)-\alpha|\,\mathbb{1}\left(h(x) \neq h^*(x)\right)\right]. \quad \square.$$

③ Problem 3.

↳ Corollary 1.2: $\quad R(\hat{h}_n) - R^* \leq 2\mathbb{E}_x\left[|p(x) - \hat{p}_n(x)|\right]$

↳ Problem: Prove Corollary 1.2.

→ Solution: By Theorem 1.1, we have

$$R(\hat{h}_n) - R^* = 2\mathbb{E}_x\left[|p(x)-\tfrac{1}{2}|\,\mathbb{1}\left(\hat{h}_n(x) \neq h^*(x)\right)\right]$$

→ We have that $\hat{h}_n(x) \neq h^*(x)$ when:

(1) $\hat{p}_n(x) < \tfrac{1}{2}$ and $p(x) \geq \tfrac{1}{2}$

(2) or: $\hat{p}_n(x) \geq \tfrac{1}{2}$ and $p(x) < \tfrac{1}{2}$.

→ For (1) we have $p(x)-\tfrac{1}{2} \leq p(x) - \hat{p}_n(x)$. (Both sides positive)

→ For (2) we have $p(x)-\tfrac{1}{2} \geq p(x) - \hat{p}_n(x)$. (Both sides negative)

⇒ For both cases (1) and (2), we have $|p(x) - \hat{p}_n(x)| \geq |p(x) - \tfrac{1}{2}|$

⇒ We have: $2\mathbb{E}_x\left[|p(x)-\hat{p}_n(x)|\right]$

$$= 2\mathbb{E}_x\left[|p(x)-\hat{p}_n(x)|\left[\mathbb{1}\left(\hat{h}_n(x)\neq h^*(x)\right) + \mathbb{1}\left(\hat{h}_n(x)=h^*(x)\right)\right]\right]$$

$$\geq 2\mathbb{E}_x\left[|p(x)-\hat{p}_n(x)|\,\mathbb{1}\left(\hat{h}_n(x)\neq h^*(x)\right)\right]$$

$$\geq 2\mathbb{E}_x\left[|p(x)-\tfrac{1}{2}|\,\mathbb{1}\left(\hat{h}_n(x)\neq h^*(x)\right)\right] = R(\hat{h}_n) - R^*. \quad \square$$