

**Exercise 1.20.** Suppose that we are in the multi-class classification setting, with  $Y = \{1, \dots, K\}$ . Suppose that we posit the following probabilistic model:

$$p_{w_{1:K}}(y = k | x) \propto \exp(\langle w_k, x \rangle), k \in \{1, \dots, K\}.$$

Show that the corresponding empirical risk for the population risk

$$L[w] = \mathbb{E}[\text{KL}(p(y | x) \| p_{w_{1:K}}(y | x))]$$

is given by the following *cross-entropy loss*:

$$\hat{L}_n[w] = -\frac{1}{n} \sum_{i=1}^n \left[ \langle w_{y_i}, x_i \rangle - \log \left[ \sum_{j=1}^K \exp(\langle w_j, x_i \rangle) \right] \right].$$

That is, show that  $\mathbb{E}[\hat{L}_n[w]] = L[w] + C$  for every fixed  $w = \{w_k\}_{k=1}^K$ , where  $C$  is a constant that does not depend on  $w$ .

$$p_{w_{1:K}}(y = k | x) = \frac{\exp(\langle w_k, x \rangle)}{\sum_{j=1}^K \exp(\langle w_j, x \rangle)}$$

Take log on both sides, we obtain

$$\log p_{w_{1:K}}(y = k | x) = \langle w_k, x \rangle - \log \left( \sum_{j=1}^K \exp(\langle w_j, x \rangle) \right)$$

$$\Rightarrow \hat{L}_n[w] = -\frac{1}{n} \sum_{i=1}^n \log p_{w_{1:K}}(y = y_i | x_i)$$

Recall Kullback - Leibler Divergence:

Given discrete prob. distr.  $P$  and  $Q$  defined on the same sample space,  $X$ , the KL divergence of  $P$  from  $Q$  is defined to be

$$\text{KL}(P \| Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$$

$$\Rightarrow KL(p(y|x) || P_{w_{1:k}}(y|x))$$

$$L[w] = \mathbb{E}_{\mathbf{x}} \left[ \sum_{k=1}^K p(y=k|x) \log \frac{p(y=k|x)}{P_{w_{1:k}}(y=k|x)} \right]$$

$$= -C - \mathbb{E} \left[ \sum_{k=1}^K p(y=k|x) \log P_{w_{1:k}}(y=k|x) \right]$$

where  $C = -\mathbb{E} \left[ \sum_{k=1}^K p(y=k|x) \log p(y=k|x) \right]$   
 $\downarrow$  independent of  $w$ .

$$\begin{aligned} \mathbb{E}[\hat{L}_n[w]] &= \mathbb{E}_{(x_i, y_i)} \left[ -\frac{1}{n} \sum_{i=1}^n \log P_{w_{1:k}}(y=y_i|x_i) \right] \\ &= \mathbb{E}_{(x, y)} [\log P_{w_{1:k}}(y|x)] \\ &= -\mathbb{E}_x \left[ \sum_{k=1}^K p(y=k|x) \log P_{w_{1:k}}(y=k|x) \right] \\ &= L[w] + C \end{aligned}$$

$$w^* = \underset{w}{\operatorname{argmin}} L[w]$$

$$= \underset{w}{\operatorname{argmin}} L[w] + C$$

$$= \underset{w}{\operatorname{argmin}} \mathbb{E}[\hat{L}_n[w]]$$

**Proposition 1.23.** Suppose that  $\ell$  is the zero-one loss,  $\mathcal{F}$  is finite, and there exists  $f \in \mathcal{F}$  satisfying  $L[f] = 0$ . Then, the empirical risk minimizer  $\hat{f}_n$  satisfies, with probability at least  $1 - \delta$ ,

$$L[\hat{f}_n] \leq \frac{\log(|\mathcal{F}|/\delta)}{n} \quad \text{with prob. } > \delta$$

**Proof:** For any  $t > 0$ , want to look at  $\{L(\hat{f}_n) > t\}$

Define  $B(t) = \{f \in \mathcal{F} \mid L(f) > t\}$

$$\{L(\hat{f}_n) > t\} = \{\hat{f}_n \in B(t)\}$$

Since by assumption, we have an  $f \in \mathcal{F}$   
s.t.  $L(f) = 0$ , then the ERM  $\hat{f}_n$  will  
always achieve zero training error, i.e.,

$$\hat{L}_n[\hat{f}_n] = 0$$

$$\underset{\text{Term} \in \arg\min_{f \in \mathcal{F}}}{\hat{L}_n[f]} = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \neq 0$$

$$\underline{L[f]} = \underline{\mathbb{E}[\ell(f(x), y)]} = 0$$

If  $\hat{f}_n \in B(t)$ , then  $\exists f \in B(t)$  s.t.  $\hat{L}_n[f] = 0$

for a fixed  $f \in B(t)$

$$\Pr[\hat{L}_n[f] = 0] = \Pr\left(\bigcap_{i=1}^n \{\operatorname{sgn}(f(x_i)) = y_i\}\right)$$

$$= \prod_{i=1}^n \left(1 - \Pr[\operatorname{sgn}(f(x_i)) \neq y_i]\right)$$

$$= (1 - L(f))^n \quad L(f)$$

$$\leq (1-t)^n$$

$$1-x \leq \exp(-x) \quad \forall x \in \mathbb{R}$$

$$\leq \exp(-tn)$$

$$\begin{aligned}
 P\{\hat{f}_n \in B(H)\} &\leq P\{\exists f \in B(H), \hat{L}_n(f) = 0\} \\
 &\leq \sum_{f \in B(H)} P(\hat{L}_n(f) = 0) \\
 &\leq |B(H)| \exp(-tn) \\
 &\leq |\mathcal{F}| \exp(-tn) = \delta
 \end{aligned}$$

## • Rademacher Random Variable

$\varepsilon_i$  is a Rademacher r.v.

$$P(\varepsilon_i = \pm 1) = \frac{1}{2}$$

$$P(\varepsilon_i = k) = 0 \quad \forall k \neq \pm 1$$

## • Rademacher Complexity

$$R_n(\mathcal{F}) := E[\sup_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n \varepsilon_i f(x_i)]$$

We first look at R.C. of linear function classes  
 $(x \sim N(0, I))$

$$\mathcal{F}_p := \{x \mapsto \langle x, w \rangle \mid \|w\|_p \leq 1\}$$

Can use dual norm  $\|\cdot\|_*$

$$\|\cdot\| = \|\cdot\|_p$$

$$\|\cdot\|_* = \|\cdot\|_q$$

$$\begin{aligned}
 \|x\|_* &:= \sup_{\|z\| \leq 1} \langle z, x \rangle \quad \frac{1}{p} + \frac{1}{q} = 1 \\
 R_n(\mathcal{F}_p) &= E[\sup_{\|w\|_p \leq 1} n^{-1} \sum_{i=1}^n \varepsilon_i \langle x_i, w \rangle] \\
 &= n^{-1} \underbrace{\sup_{\|w\|_p \leq 1} \langle w, \sum_{i=1}^n \varepsilon_i x_i \rangle}_{\text{red line}}
 \end{aligned}$$

$$= n^{-1} \sum_{i=1}^n \varepsilon_i \|x_i\|_q$$

$$\text{where } \frac{1}{p} + \frac{1}{q} = 1$$