# Supplementary Material: A Minimax Approach to Supervised Learning

Farzan Farnia, David Tse Department of Electrical Engineering Stanford University {farnia,dntse}@stanford.edu

# 1 Proof of Theorem 1

# 1.a Weak Version

First, we list the assumptions of the weak version of Theorem 1:

- $\Gamma$  is convex and closed,
- Loss function  $L$  is bounded by a constant  $C$ ,
- $\mathcal{X}, \mathcal{Y}$  are finite,
- Risk set  $S = \{ [L(y, a)]_{y \in \mathcal{Y}} : a \in \mathcal{A} \}$  is closed.

Given these assumptions, Sion's minimax theorem [\[1\]](#page-7-0) implies that the minimax problem has a finite answer  $H^*$ ,

$$
H^* := \sup_{P \in \Gamma} \inf_{\psi \in \Psi} \mathbb{E}[L(Y, \psi(X))] = \inf_{\psi \in \Psi} \sup_{P \in \Gamma} \mathbb{E}[L(Y, \psi(X))]. \tag{1}
$$

Thus, there exists a sequence of decision rules  $(\psi_n)_{n=1}^{\infty}$  for which

<span id="page-0-1"></span>
$$
\lim_{n \to \infty} \sup_{P \in \Gamma} \mathbb{E}[L(Y, \psi_n(X))] = H^*.
$$
\n(2)

As we supposed, the risk set S is closed. Therefore, the randomized risk set<sup>[1](#page-0-0)</sup>  $S_r = \{ [L(y, \zeta)]_{y \in \mathcal{Y}} :$  $\zeta \in \mathcal{Z}$  } defined over the space of randomized acts  $\mathcal Z$  is also closed and, since L is bounded, is a compact subset of  $\mathbb{R}^{|\mathcal{Y}|}$ . Therefore, since X and Y are both finite, we can find a randomized decision rule  $\psi^*$  which on taking a subsequence  $(n_k)_{k=1}^{\infty}$  satisfies

$$
\forall x \in \mathcal{X}, y \in \mathcal{Y}: L(y, \psi^*(x)) = \lim_{k \to \infty} L(y, \psi_{n_k}(x)).
$$
 (3)

Then  $\psi^*$  is a robust Bayes decision rule against  $\Gamma$ , because

$$
\sup_{P \in \Gamma} \mathbb{E}\left[L(Y, \psi^*(X))\right] = \sup_{P \in \Gamma} \lim_{k \to \infty} \mathbb{E}\left[L(Y, \psi_{n_k}(X))\right] \le \lim_{k \to \infty} \sup_{P \in \Gamma} \mathbb{E}[L(Y, \psi_{n_k}(X))] = H^*.
$$
 (4)

Moreover, since  $\Gamma$  is assumed to be convex and closed (hence compact),  $H(Y|X)$  achieves its supremum over  $\Gamma$  at some distribution  $P^*$ . By the definition of conditional entropy, [\(4\)](#page-0-1) implies that

$$
E_{P^*}[L(Y, \psi^*(X))] \le \sup_{P \in \Gamma} \mathbb{E}\left[L(Y, \psi^*(X))\right] \le H^* = H_{P^*}(Y|X),\tag{5}
$$

which shows that  $\psi^*$  is a Bayes decision rule for  $P^*$  as well. This completes the proof.

<span id="page-0-0"></span> ${}^{1}L(y,\zeta)$  is a short-form for  $E[L(y,A)]$  where  $A \in \mathcal{A}$  is a random action distributed according to  $\zeta$ .

<sup>30</sup>th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

#### 1.b Strong Version

Let's recall the assumptions of the strong version of Theorem 1:

- $\bullet$   $\Gamma$  is convex.
- For any distribution  $P \in \Gamma$ , there exists a Bayes decision rule.
- We assume continuity in Bayes decision rules over  $\Gamma$ , i.e., if a sequence of distributions  $(Q_n)_{n=1}^{\infty} \in \Gamma$  with the corresponding Bayes decision rules  $(\psi_n)_{n=1}^{\infty}$  converges to Q with a Bayes decision rule  $\psi$ , then under any  $P \in \Gamma$ , the expected loss of  $\psi_n$  converges to the expected loss of  $\psi$ .
- $P^*$  maximizes the conditional entropy  $H(Y|X)$ .

**Note:** A particular structure used in our paper is given by fixing the marginal  $P_X$  across Γ. Under this structure, the condition of the continuity in Bayes decision rules reduces to the continuity in Bayes acts over  $P_Y$ 's in  $\Gamma_{Y|X}$ . It can be seen that while this condition holds for the logarithmic and quadratic loss functions, it does not hold for the 0-1 loss.

<span id="page-1-0"></span>Let  $\psi^*$  be a Bayes decision rule for  $P^*$ . We need to show that  $\psi^*$  is a robust Bayes decision rule against  $\Gamma$ . To show this, it suffices to show that  $(P^*, \psi^*)$  is a saddle point of the mentioned minimax problem, i.e.,

$$
\mathbb{E}_{P^*}[L(Y,\psi^*(X))] \le \mathbb{E}_{P^*}[L(Y,\psi(X))],\tag{6}
$$

<span id="page-1-1"></span>and

$$
\mathbb{E}_{P^*}[L(Y, \psi^*(X))] \ge \mathbb{E}_P[L(Y, \psi^*(X))].\tag{7}
$$

Clearly, inequality [\(6\)](#page-1-0) holds due to the definition of the Bayes decision rule. To show [\(7\)](#page-1-1), let us fix an arbitrary distribution  $P \in \Gamma$ . For any  $\lambda \in (0,1]$ , define  $P_{\lambda} = \lambda P + (1 - \lambda)P^*$ . Notice that  $P_{\lambda} \in \Gamma$ since Γ is convex. Let  $\psi_{\lambda}$  be a Bayes decision rule for  $P_{\lambda}$ . Due to the linearity of the expected loss in the probability distribution, we have

$$
\mathbb{E}_{P}[L(Y,\psi_{\lambda}(X))] - \mathbb{E}_{P^*}[L(Y,\psi_{\lambda}(X))] = \frac{\mathbb{E}_{P_{\lambda}}[L(Y,\psi_{\lambda}(X))] - \mathbb{E}_{P^*}[L(Y,\psi_{\lambda}(X))]}{\lambda} \leq \frac{H_{P_{\lambda}}(Y|X) - H_{P^*}(Y|X)}{\lambda} \leq 0,
$$

for any  $0 < \lambda < 1$ . Here the first inequality is due to the definition of the conditional entropy and the last inequality holds since  $P^*$  maximizes the conditional entropy over  $\Gamma$ . Applying the assumption of the continuity in Bayes decision rules, we have

$$
\mathbb{E}_{P}[L(Y,\psi^{*}(X))] - \mathbb{E}_{P^{*}}[L(Y,\psi^{*}(X))] = \lim_{\lambda \to 0} \mathbb{E}_{P}[L(Y,\psi_{\lambda}(X))] - \mathbb{E}_{P^{*}}[L(Y,\psi_{\lambda}(X))] \le 0, (8)
$$

which makes the proof complete.

# 2 Proof of Theorem 2

Let us recall the definition of the set  $\Gamma(Q)$ :

$$
\Gamma(Q) = \{ P_{\mathbf{X}, Y} : P_{\mathbf{X}} = Q_{\mathbf{X}}, \forall 1 \le i \le t : \| \mathbb{E}_P [\theta_i(Y) \mathbf{X}] - \mathbb{E}_Q [\theta_i(Y) \mathbf{X}] \| \le \epsilon_i \}.
$$
\n(9)

Defining  $\tilde{\mathbf{E}}_i \triangleq \mathbb{E}_Q \left[ \theta_i(Y) \mathbf{X} \right]$  and  $C_i \triangleq \{ \mathbf{u} : ||\mathbf{u} - \tilde{\mathbf{E}}_i|| \leq \epsilon_i \}$ , we have

$$
\max_{P \in \Gamma(Q)} H(Y|\mathbf{X}) = \max_{P, \mathbf{w}: \forall i: \mathbf{w}_i = \mathbb{E}_P[\theta_i(Y)\mathbf{X}]} \mathbb{E}_{Q_{\mathbf{X}}}[H_P(Y|\mathbf{X}=\mathbf{x})] + \sum_{i=1}^t I_{C_i}(\mathbf{w}_i)
$$
(10)

where  $I_C$  is the indicator function for the set C defined as

<span id="page-1-2"></span>
$$
I_C(x) = \begin{cases} 0 & \text{if } x \in C, \\ -\infty & \text{Otherwise.} \end{cases}
$$
 (11)

First of all, the law of iterated expectations implies that  $\mathbb{E}_P\left[\theta_i(Y) \mathbf{X}\right] = \mathbb{E}_{Q_{\mathbf{X}}}\bigg[\, \mathbf{X}\, \mathbb{E}[\theta_i(Y)|\mathbf{X}=\mathbf{x}]\, \bigg].$ Furthermore, [\(10\)](#page-1-2) is equivalent to a convex optimization problem where it is not hard to check that

the Slater condition is satisfied. Hence strong duality holds and we can write the dual problem as

$$
\min_{\mathbf{A}} \sup_{P_{Y|\mathbf{X}}, \mathbf{w}} \mathbb{E}_{Q_{\mathbf{X}}} \left[ H_P(Y|\mathbf{X}=\mathbf{x}) + \sum_{i=1}^t \mathbb{E}[\theta_i(Y)|\mathbf{X}=\mathbf{x}]\mathbf{A}_i\mathbf{X} \right] + \sum_{i=1}^t \left[ I_{C_i}(\mathbf{w}_i) - \mathbf{A}_i\mathbf{w}_i \right],\tag{12}
$$

where the rows of matrix A, denoted by  $A_i$ , are the Lagrange multipliers for the constraints of  $\mathbf{w}_i = \mathbb{E}_P [\theta_i(Y) \mathbf{X}]$ . Notice that the above problem decomposes across  $P_{Y | \mathbf{X} = \mathbf{x}}$ 's and  $\mathbf{w}_i$ 's. Hence, the dual problem can be rewritten as

$$
\min_{\mathbf{A}} \left[ \mathbb{E}_{Q_{\mathbf{X}}} \left[ \sup_{P_{Y|{\mathbf{X}}=\mathbf{x}}} H_P(Y|{\mathbf{X}}=\mathbf{x}) + \sum_{i=1}^t \mathbb{E}[\theta_i(Y)|{\mathbf{X}}=\mathbf{x}] \mathbf{A}_i {\mathbf{X}} \right] + \sum_{i=1}^t \sup_{\mathbf{w}_i} \left[ I_{C_i}(\mathbf{w}_i) - \mathbf{A}_i \mathbf{w}_i \right] \right] \tag{13}
$$

Furthermore, according to the definition of  $F_{\theta}$ , we have

<span id="page-2-2"></span><span id="page-2-1"></span><span id="page-2-0"></span>
$$
F_{\theta}(\mathbf{A}\mathbf{x}) = \sup_{P_{Y|\mathbf{X}=\mathbf{x}}} H(Y|\mathbf{X}=\mathbf{x}) + \mathbb{E}[\theta(Y)|\mathbf{X}=\mathbf{x}]^T \mathbf{A}\mathbf{x}.
$$
 (14)

Moreover, the definition of the dual norm  $\|\cdot\|_*$  implies

$$
\sup_{\mathbf{w}_i} I_{C_i}(\mathbf{w}_i) - \mathbf{A}_i \mathbf{w}_i = \max_{\mathbf{u} \in C_i} -\mathbf{A}_i \mathbf{u} = -\mathbf{A}_i \tilde{\mathbf{E}}_i + \epsilon_i ||\mathbf{A}_i||_*.
$$
 (15)

Plugging [\(14\)](#page-2-0) and [\(15\)](#page-2-1) in [\(13\)](#page-2-2), the dual problem can be simplified to

$$
\min_{\mathbf{A}} \mathbb{E}_{Q_{\mathbf{X}}} \left[ F_{\theta}(\mathbf{A}\mathbf{X}) - \sum_{i=1}^{t} \mathbf{A}_{i} \tilde{\mathbf{E}}_{i} \right] + \sum_{i=1}^{t} \epsilon_{i} \|\mathbf{A}_{i}\|_{*}
$$
\n
$$
= \min_{\mathbf{A}} \mathbb{E}_{Q} \left[ F_{\theta}(\mathbf{A}\mathbf{X}) - \theta(Y)^{T} \mathbf{A}\mathbf{X} \right] + \sum_{i=1}^{t} \epsilon_{i} \|\mathbf{A}_{i}\|_{*}, \tag{16}
$$

which is equal to the primal problem [\(10\)](#page-1-2) since the strong duality holds. Furthermore, note that we can rewrite the definition given for  $F_{\theta}$  as

$$
F_{\theta}(\mathbf{z}) = \max_{\mathbf{E} \in \mathbb{R}^t} G(\mathbf{E}) + \mathbf{E}^T \mathbf{z},
$$
 (17)

where we define

$$
G(\mathbf{E}) = \begin{cases} \max_{P \in \mathcal{P}_{\mathcal{Y}} : \mathbb{E}[\theta(Y)] = \mathbf{E}} H(Y) & \text{if } \{P \in \mathcal{P}_{\mathcal{Y}} : \mathbb{E}[\theta(Y)] = \mathbf{E}\} \neq \emptyset \\ -\infty & \text{Otherwise.} \end{cases}
$$
(18)

Observe that  $F_{\theta}$  is the convex conjugate of the convex  $-G$ . Therefore, applying the derivative property of convex conjugates [\[2\]](#page-7-1) to [\(14\)](#page-2-0),

<span id="page-2-3"></span>
$$
\mathbb{E}_{P^*}[\boldsymbol{\theta}(Y) | \mathbf{X} = \mathbf{x}] \in \partial F_{\boldsymbol{\theta}}(\mathbf{A}^* \mathbf{x}).
$$
\n(19)

Here,  $\partial F_{\theta}$  denotes the subgradient of  $F_{\theta}$ . Assuming  $F_{\theta}$  is differentiable at  $\mathbf{A}^*$ x, [\(19\)](#page-2-3) implies that

$$
\mathbb{E}_{P^*}[\theta(Y) | \mathbf{X} = \mathbf{x}] = \nabla F_{\theta}(\mathbf{A}^* \mathbf{x}).
$$
\n(20)

# 2.a A generalization of Theorem 2

It can be seen that the above proof can be slightly generalized to prove the following generalization of Theorem 2.

Theorem. *Given a conjugate pair of convex functions* g, g<sup>∗</sup> *, the following duality holds*

$$
\max_{P:\;P_X=Q_X} H(Y|\mathbf{X}) - \sum_{i=1}^t g\bigg(\mathbb{E}_P[\theta_i(Y)\mathbf{X}] - \mathbb{E}_Q[\theta_i(Y)\mathbf{X}]\bigg) = \tag{21}
$$

$$
\min_{\mathbf{A}\in\mathbb{R}^{t\times d}}\mathbb{E}_{Q}\left[F_{\theta}(\mathbf{A}\mathbf{X})-\theta(Y)^{T}\mathbf{A}\mathbf{X}\right]+\sum_{i=1}^{t}g^{*}(\mathbf{A}_{i}),
$$
\n(22)

where  $A_i$  denotes the ith row of  $A$ . In addition, for the optimal  $P^*$  and  $A^*$ 

$$
\mathbb{E}_{P^*}[\boldsymbol{\theta}(Y) | \mathbf{X} = \mathbf{x}] = \nabla F_{\boldsymbol{\theta}}(\mathbf{A}^* \mathbf{x}).
$$
 (23)

**Corollary.** *Consider a pair of dual norms*  $\|\cdot\|$ ,  $\|\cdot\|_*$ *. Then, the following duality holds* 

$$
\max_{P:\;P_X=Q_X} H(Y|\mathbf{X}) - \sum_{i=1}^t \frac{1}{2\lambda_i} \left\| \mathbb{E}_P[\theta_i(Y)\mathbf{X}] - \mathbb{E}_Q[\theta_i(Y)\mathbf{X}] \right\|^2 = \tag{24}
$$

$$
\min_{\mathbf{A}\in\mathbb{R}^{t\times d}}\mathbb{E}_{Q}\left[F_{\theta}(\mathbf{A}\mathbf{X})-\theta(Y)^{T}\mathbf{A}\mathbf{X}\right]+\sum_{i=1}^{t}\frac{\lambda_{i}}{2}\left\|\mathbf{A}_{i}\right\|_{*}^{2},\tag{25}
$$

where  $\lambda_i$ 's are positive real numbers and  $\mathbf{A}_i$  denotes the ith row of  $\mathbf{A}$ *. Moreover, for the optimal*  $P^*$ *and* A<sup>∗</sup>

$$
\mathbb{E}_{P^*}[\boldsymbol{\theta}(Y) | \mathbf{X} = \mathbf{x}] = \nabla F_{\boldsymbol{\theta}}(\mathbf{A}^* \mathbf{x}).
$$
 (26)

# 3 Proof of Theorem 3

First, we aim to show that

$$
\max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X}))] \leq \mathbb{E}_{\tilde{P}}\left[F_{\theta}(\hat{\mathbf{A}}_n\mathbf{X}) - \theta(Y)^T\hat{\mathbf{A}}_n\mathbf{X}\right] + \sum_{i=1}^t \epsilon_i \|\hat{\mathbf{A}}_{n_i}\|_*\tag{27}
$$

where  $\hat{A}_n$  denotes the solution to the RHS of the duality equation in Theorem 2 for the empirical distribution  $\hat{P}_n$ . Similar to the duality proven in Theorem 2, we can show that

 $\mathbf{r}$ 

$$
\max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X}))] = \min_{\mathbf{A}} \mathbb{E}_{\tilde{P}_X} \bigg[ \sup_{P_{Y|X} \in \mathcal{P}_{\mathcal{Y}}} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X})) | \mathbf{X} = \mathbf{x}] + \mathbb{E}[\boldsymbol{\theta}(Y) | \mathbf{X} = \mathbf{x}]^T \mathbf{A} \mathbf{X} \bigg]
$$
  
\n
$$
- \mathbb{E}_{\tilde{P}}[\boldsymbol{\theta}(Y)^T \mathbf{A} \mathbf{X}] + \sum_{i=1}^t \epsilon_i ||\mathbf{A}_i||_*
$$
  
\n
$$
\leq \mathbb{E}_{\tilde{P}_X} \bigg[ \sup_{P_{Y|X=x} \in \mathcal{P}_{\mathcal{Y}}} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X})) | \mathbf{X} = \mathbf{x}] + \mathbb{E}[\boldsymbol{\theta}(Y) | \mathbf{X}]^T \hat{\mathbf{A}}_n \mathbf{X} \bigg]
$$
  
\n
$$
- \mathbb{E}_{\tilde{P}}[\boldsymbol{\theta}(Y)^T \hat{\mathbf{A}}_n \mathbf{X}] + \sum_{i=1}^t \epsilon_i ||\hat{\mathbf{A}}_{n_i}||_*
$$
  
\n
$$
= \mathbb{E}_{\tilde{P}} \bigg[ F_{\boldsymbol{\theta}}(\hat{\mathbf{A}}_n \mathbf{X}) - \boldsymbol{\theta}(Y)^T \hat{\mathbf{A}}_n \mathbf{X} \bigg] + \sum_{i=1}^t \epsilon_i ||\hat{\mathbf{A}}_{n_i}||_*.
$$

Here we first upper bound the minimum by taking the specific  $A = \hat{A}_n$ . Then the equality holds because  $\hat{\psi}_n$  is a robust Bayes decision rule against  $\Gamma(\hat{P}_n)$  and therefore adding the second term based on  $\hat{\mathbf{A}}_n\mathbf{x}$ ,  $\hat{\psi}_n(\mathbf{x})$  results in a saddle point for the following problem

<span id="page-3-0"></span>
$$
F_{\theta}(\hat{\mathbf{A}}_{n}\mathbf{X}) = \sup_{P \in \mathcal{P}_{\mathcal{Y}}} H(Y) + \mathbb{E}[\theta(Y)]^{T} \hat{\mathbf{A}}_{n}\mathbf{X}
$$
  
= 
$$
\sup_{P \in \mathcal{P}_{\mathcal{Y}}} \inf_{\zeta \in \mathcal{Z}} \mathbb{E}[L(Y,\zeta)] + \mathbb{E}[\theta(Y)]^{T} \hat{\mathbf{A}}_{n}\mathbf{X}
$$
  
= 
$$
\sup_{P \in \mathcal{P}_{\mathcal{Y}}} \mathbb{E}[L(Y,\hat{\psi}_{n}(\mathbf{X}))] + \mathbb{E}[\theta(Y)]^{T} \hat{\mathbf{A}}_{n}\mathbf{X}.
$$

Therefore, by Theorem 2 we have

$$
\max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X}))] - \max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \tilde{\psi}(\mathbf{X}))] \leq
$$
\n
$$
\mathbb{E}_{\tilde{P}}[F_{\theta}(\hat{\mathbf{A}}_n \mathbf{X}) - \theta(Y)^T \hat{\mathbf{A}}_n \mathbf{X}] + \sum_{i=1}^t \epsilon_i ||\hat{\mathbf{A}}_{n_i}||_* - \mathbb{E}_{\tilde{P}}[F_{\theta}(\tilde{\mathbf{A}} \mathbf{X}) - \theta(Y)^T \tilde{\mathbf{A}} \mathbf{X}] - \sum_{i=1}^t \epsilon_i ||\tilde{\mathbf{A}}_i||_*.
$$
\n(28)

As a result, we only need to bound the uniform convergence rate in the other side of the duality. Note that by the definition of  $F_{\theta}$ ,

$$
\forall P \in \mathcal{P}_{\mathcal{Y}}, \mathbf{z} \in \mathbb{R}^t : \quad F_{\theta}(\mathbf{z}) - \mathbb{E}_P[\theta(Y)]^T \mathbf{z} \ge H_P(Y) \ge 0. \tag{29}
$$

Hence,  $\forall$  **A** :  $F_{\theta}$ (**AX**) –  $\mathbb{E}[\theta(Y)]^T A X \ge 0$  and comparing the optimal solution to the RHS of the duality equation in Theorem 2 to the case  $A = 0$  implies that for any possible solution  $A^*$ 

$$
\forall 1 \leq i \leq t: \quad \epsilon_i \| \mathbf{A}_i^* \|_q \leq \sum_{j=1}^t \epsilon_j \| \mathbf{A}_j^* \|_q \leq F_{\theta}(\mathbf{0}) = \max_{P \in \mathcal{P}_{\mathcal{Y}}} H(Y) = M. \tag{30}
$$

Hence, since  $1 \le q \le 2$ , we only need to bound the uniform convergence rate in a bounded space where  $\forall 1 \leq i \leq t : ||A_i||_2 \leq ||A_i||_q \leq \frac{M}{\epsilon_i}$ . Also, applying the derivative property of the conjugate relationship indicates that  $\partial F_{\theta}(\mathbf{z})$  is a subset of the convex hull of  $\{\mathbb{E}[\theta(Y)] : P \in \mathcal{P}_{\mathcal{Y}}\}$ . Therefore, when  $\theta(Y)$  includes only one variable, for any  $u \in \partial F_{\theta}(z)$  we have  $|u| \leq L$ , and  $F_{\theta}(z) - \theta(Y)z$ is 2L-Lipschitz in z. As a result, since  $||\mathbf{X}||_2 \leq B$  and  $|\theta(Y)| \leq L$  for any  $\alpha_1, \alpha_2 \in \mathbb{R}^d$  such that  $\|\boldsymbol{\alpha}_i\|_2 \leq \frac{M}{\epsilon},$ 

$$
\forall \mathbf{x}_1, \mathbf{x}_2, y_1, y_2: [F_{\theta}(\boldsymbol{\alpha}_1^T \mathbf{x}_1) - \theta(y_1)\boldsymbol{\alpha}_1^T \mathbf{x}_1] - [F_{\theta}(\boldsymbol{\alpha}_2^T \mathbf{x}_2) - \theta(y_2)\boldsymbol{\alpha}_2^T \mathbf{x}_2] \le \frac{4BML}{\epsilon}
$$
(31)

Consequently, we can apply standard uniform convergence results given convexity-Lipschitzness-boundedness [\[3\]](#page-7-2) to show that for any  $\delta > 0$  with a probability at least  $1 - \delta$ 

$$
\forall \alpha \in \mathbb{R}^d, \|\alpha\|_2 \le \frac{M}{\epsilon}:
$$
\n
$$
\mathbb{E}_{\tilde{P}}\left[F_{\theta}(\alpha^T \mathbf{X}) - \theta(Y)\alpha^T \mathbf{X}\right] - \mathbb{E}_{\hat{P}_n}\left[F_{\theta}(\alpha^T \mathbf{X}) - \theta(Y)\alpha^T \mathbf{X}\right] \le \frac{4BLM}{\epsilon \sqrt{n}} \left(1 + \sqrt{\frac{\log(2/\delta)}{2}}\right).
$$
\n(32)

Therefore, considering  $\hat{\alpha}_n$  and  $\tilde{\alpha}$  as the solution to the dual problems corresponding to the empirical and underlying cases, for any  $\delta > 0$  with a probability at least  $1 - \delta/2$ 

$$
\mathbb{E}_{\tilde{P}}\left[F_{\theta}(\hat{\alpha}_n^T \mathbf{X}) - \theta(Y)\hat{\alpha}_n^T \mathbf{X}\right] + \epsilon \|\hat{\alpha}_n\|_q
$$
\n
$$
-\mathbb{E}_{\hat{P}_n}\left[F_{\theta}(\hat{\alpha}_n^T \mathbf{X}) - \theta(Y)\hat{\alpha}_n^T \mathbf{X}\right] - \epsilon \|\hat{\alpha}_n\|_q \le \frac{4BLM}{\epsilon \sqrt{n}} \left(1 + \sqrt{\frac{\log(4/\delta)}{2}}\right).
$$
\n(33)

Since  $\hat{\alpha}_n$  is minimizing the objective for  $Q = \hat{P}_n$ ,

<span id="page-4-3"></span><span id="page-4-2"></span><span id="page-4-1"></span><span id="page-4-0"></span>
$$
\mathbb{E}_{\hat{P}_n} \left[ F_{\theta}(\hat{\alpha}_n^T \mathbf{X}) - \theta(Y) \hat{\alpha}_n^T \mathbf{X} \right] + \epsilon ||\hat{\alpha}_n||_q
$$
\n
$$
- \mathbb{E}_{\hat{P}_n} \left[ F_{\theta}(\tilde{\alpha}^T \mathbf{X}) - \theta(Y) \tilde{\alpha}^T \mathbf{X} \right] - \epsilon ||\tilde{\alpha}||_q \le 0.
$$
\n(34)

Also, since  $\tilde{\alpha}$  does not depend on the samples, the Hoeffding's inequality implies that with a probability at least  $1 - \delta/2$ 

$$
\mathbb{E}_{\hat{P}_n} \left[ F_{\theta}(\tilde{\boldsymbol{\alpha}}^T \mathbf{X}) - \theta(Y) \tilde{\boldsymbol{\alpha}}^T \mathbf{X} \right] + \epsilon \| \tilde{\boldsymbol{\alpha}} \|_q
$$
\n
$$
- \mathbb{E}_{\tilde{P}} \left[ F_{\theta}(\tilde{\boldsymbol{\alpha}}^T \mathbf{X}) - \theta(Y) \tilde{\boldsymbol{\alpha}}^T \mathbf{X} \right] - \epsilon \| \tilde{\boldsymbol{\alpha}} \|_q \le \frac{2BML}{\epsilon} \sqrt{\frac{\log(4/\delta)}{2n}}.
$$
\n(35)

Applying the union bound, combining [\(33\)](#page-4-0), [\(34\)](#page-4-1), [\(35\)](#page-4-2) shows that with a probability at least  $1 - \delta$ , we have

$$
\mathbb{E}_{\hat{P}_n} \left[ F_{\theta}(\hat{\alpha}_n^T \mathbf{X}) - \theta(Y) \hat{\alpha}_n^T \mathbf{X} \right] + \epsilon ||\hat{\alpha}_n||_q
$$
\n
$$
- \mathbb{E}_{\tilde{P}} \left[ F_{\theta}(\tilde{\alpha}^T \mathbf{X}) - \theta(Y) \tilde{\alpha}^T \mathbf{X} \right] - \epsilon ||\tilde{\alpha}||_q \le \frac{4BLM}{\epsilon \sqrt{n}} \left( 1 + \frac{3}{2} \sqrt{\frac{\log(4/\delta)}{2}} \right).
$$
\n(36)

Given  $(28)$  and  $(36)$ , the proof is complete.

Note that we can improve the result in the case  $q = 1$  by using the same proof and plugging in the Rademacher complexity of the  $\ell_1$ -bounded linear functions. Here, we replace the assumption that  $\|\mathbf{X}\|_2 \leq B$  with  $\|\mathbf{X}\|_{\infty} \leq B$  which can be much weaker for high-dimensional X's.

**Theorem.** *Consider a loss function L* with the entropy *H* and suppose  $\theta(Y)$  *includes only one element. Let*  $M = \max_{P \in \mathcal{P}_{\mathcal{Y}}} H(Y)$  *be the maximum entropy value over*  $\mathcal{P}_{\mathcal{Y}}$ *. Also, take*  $\| \cdot \|/\| \cdot \|_*$ *to be the*  $\ell_{\infty}/\ell_1$  *pair . Given that*  $\bf{X}$  *is a d-dimensional vector with*  $\|\bf{X}\|_{\infty} \leq B$ *, and*  $|\theta(Y)| \leq L$ *, for any*  $\delta > 0$  *with probability at least*  $1 - \delta$ 

$$
\max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \hat{\psi}_n(\mathbf{X}))] - \max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \tilde{\psi}(\mathbf{X}))] \le \frac{4BLM}{\epsilon \sqrt{n}} \left(\sqrt{2\log(2d)} + \sqrt{\frac{9\log(4/\delta)}{8}}\right). \tag{37}
$$

# 4 0-1 Loss: minimax SVM

#### **4.a**  $F_{\theta}$  derivation

Given the defined one-hot encoding  $\theta$  we define  $\tilde{\mathbf{z}} = (\mathbf{z}, 0)$  and represent each randomized decision rule  $\zeta$  with its corresponding loss vector  $\mathbf{L} \in \mathbb{R}^{t+1}$  such that  $L_i = L_{0-1}(i, \zeta)$  denotes the 0-1 loss suffered by  $\zeta$  when  $Y=i.$  It can be seen that  ${\bf L}$  is a feasible loss vector if and only if  $\forall\,i:0\le L_i\le 1$ and  $\sum_{i=1}^{t+1} L_i = t$ . Then,

$$
F_{\theta}(\mathbf{z}) = \max_{\substack{\mathbf{p} \in \mathbb{R}^{t+1}: \mathbf{1}^T \mathbf{p} = 1, \ \mathbf{L} \in \mathbb{R}^{t+1}: \mathbf{1}^T \mathbf{L} = t, \\ \forall i: 0 \le p_i}} \sum_{\substack{\mathbf{v} \in \mathbb{R}^{t+1}: \mathbf{1}^T \mathbf{L} = t, \\ \forall i: 0 \le L_i \le 1}} \sum_{i=1}^{t+1} p_i(\tilde{z}_i + L_i). \tag{38}
$$

Hence, Sion's minimax theorem implies that the above minimax problem has a saddle point. Thus,

$$
F_{\theta}(\mathbf{z}) = \min_{\substack{\mathbf{L} \in \mathbb{R}^{t+1}: \mathbf{1}^T \mathbf{L} = t, \ 1 \le i \le t+1 \\ \forall i : 0 \le L_i \le 1}} \max_{1 \le i \le t+1} \{ \tilde{z}_i + L_i \}. \tag{39}
$$

Consider  $\sigma$  as the permutation sorting  $\tilde{z}$  in a descending order and for simplicity let  $\tilde{z}_{(i)} = \tilde{z}_{\sigma(i)}$ . Then,

$$
\forall 1 \le k \le t+1: \quad \max_{1 \le i \le t+1} \{ \tilde{z}_i + L_i \} \ge \frac{1}{k} \sum_{i=1}^k [\tilde{z}_{\sigma(i)} + L_{\sigma(i)}] \ge \frac{k-1 + \sum_{i=1}^k \tilde{z}_{(i)}}{k},\tag{40}
$$

which is independent of the value of  $L_i$ 's. Therefore,

<span id="page-5-0"></span>
$$
\max_{1 \le k \le t+1} \frac{k-1 + \sum_{i=1}^k \tilde{z}_{(i)}}{k} \le F_{\theta}(\mathbf{z}).\tag{41}
$$

On the other hand, if we let  $k_{\max}$  be the largest index satisfying  $\sum_{i=1}^{k_{\max}} [\tilde{z}_{(i)} - \tilde{z}_{(k_{\max})}] < 1$  and define

$$
\forall 1 \le j \le t+1: \quad L^*_{\sigma(j)} = \begin{cases} \frac{k_{\max} - 1 + \sum_{i=1}^{k_{\max}} \tilde{z}_{(i)}}{k_{\max}} - \tilde{z}_{(j)} & \text{if } \sigma(j) \le k_{\max} \\ 1 & \text{if } \sigma(j) > k_{\max}, \end{cases} \tag{42}
$$

we can simply check that  $\mathbf{L}^*$  is a feasible point since  $\sum_{i=1}^{t+1} L_i^* = t$  and  $L_{\sigma(k_{\text{max}})}^* \leq 1$  so for all *i*'s  $L^*_{\sigma(i)} \leq 1$ . Also,  $L^*_{\sigma(1)} \geq 0$  because  $\tilde{z}_{(1)} - \tilde{z}_{(j)} < 1$  for any  $j \leq k_{\text{max}}$ , so for all  $i$ 's  $L^*_{\sigma(i)} \geq 0$ . Then for this  $L^*$  we have

$$
F_{\theta}(\mathbf{z}) \le \max_{1 \le i \le t+1} \{ \tilde{z}_i + L_i^* \} = \frac{k_{\text{max}} - 1 + \sum_{i=1}^{k_{\text{max}}} \tilde{z}_{(i)}}{k_{\text{max}}}.
$$
(43)

Therefore, [\(41\)](#page-5-0) holds with equality and achieves its maximum at  $k = k_{\text{max}}$ ,

$$
F_{\theta}(\mathbf{z}) = \max_{1 \le k \le t+1} \frac{k-1+\sum_{i=1}^{k} \tilde{z}_{(i)}}{k} = \frac{k_{\max} - 1 + \sum_{i=1}^{k_{\max}} \tilde{z}_{(i)}}{k_{\max}}.
$$
(44)

Moreover,  $L^*$  corresponds to a randomized robust Bayes act, where we select label  $i$  according to the probability vector  $\overline{p^*} = 1 - \overline{L^*}$  that is

$$
\forall 1 \le j \le t+1: \quad p^*_{\sigma(j)} = \begin{cases} \frac{1 - \sum_{i=1}^{k_{\text{max}}} \tilde{z}_{(i)}}{k_{\text{max}}} + \tilde{z}_{(j)} & \text{if } \sigma(j) \le k_{\text{max}}\\ 0 & \text{if } \sigma(j) > k_{\text{max}}. \end{cases} \tag{45}
$$

Given  $F_{\theta}$  we can simply derive the gradient  $\nabla F_{\theta}$  to find the entropy maximizing distribution. Here if the inequality  $\sum_{i=1}^{k_{\text{max}}} [\tilde{\mathbf{z}}_{\sigma(i)} - \tilde{\mathbf{z}}_{(k_{\text{max}}+1)}] \ge 1$  holds strictly, which is true almost everywhere on  $\mathbb{R}^t$ ,

$$
\forall 1 \le i \le t: \left(\nabla F_{\theta}(\mathbf{z})\right)_i = \begin{cases} 1/k_{\text{max}} & \text{if } \sigma(i) \le k_{\text{max}}, \\ 0 & \text{Otherwise.} \end{cases} \tag{46}
$$

If the inequality does not strictly hold,  $F_{\theta}$  is not differentiable at z; however, the above vector still lies in the subgradient  $\partial F_{\theta}(\mathbf{z})$ .

#### 4.b Sufficient Conditions for Applying Theorem 1.a

As supposed in Theorem 1.a, the space  $X$  should be finite in order to apply that result. Here, we show for the proposed structure on  $\Gamma(Q)$  one can relax this condition while Theorem 1.a still holds. It is because, as shown in the proofs of Theorems 2 and 3, we have

$$
\inf_{\psi \in \Psi} \max_{P \in \Gamma(\tilde{P})} \mathbb{E}[L(Y, \psi(\mathbf{X}))] = \inf_{\psi \in \Psi} \min_{\mathbf{A}} \mathbb{E}_{\tilde{P}_{X}} \Big[ \sup_{P_{Y|X} \in \mathcal{P}_{Y}} \mathbb{E}[L(Y, \psi(\mathbf{X})) | \mathbf{X} = \mathbf{x}] \n+ \mathbb{E}[\theta(Y) | \mathbf{X} = \mathbf{x}]^{T} \mathbf{A} \mathbf{X} \Big] - \mathbb{E}_{\tilde{P}}[\theta(Y)^{T} \mathbf{A} \mathbf{X}] + \sum_{i=1}^{t} \epsilon_{i} ||\mathbf{A}_{i}||_{*} \n= \min_{\mathbf{A}} \mathbb{E}_{\tilde{P}_{X}} \Big[ \inf_{\psi(\mathbf{x}) \in \mathcal{Z}} \sup_{P_{Y|X} \in \mathcal{P}_{Y}} \mathbb{E}[L(Y, \psi(\mathbf{x})) | \mathbf{X} = \mathbf{x}] \n+ \mathbb{E}[\theta(Y) | \mathbf{X} = \mathbf{x}]^{T} \mathbf{A} \mathbf{X} \Big] - \mathbb{E}_{\tilde{P}}[\theta(Y)^{T} \mathbf{A} \mathbf{X}] + \sum_{i=1}^{t} \epsilon_{i} ||\mathbf{A}_{i}||_{*}.
$$

Therefore, given this structure the minimax problem decouples across different x's. Hence, the assumption of finite X is no longer needed, because as long as  $\theta$  is a bounded function (which is true for the one-hot encoding  $\theta$ ), the rest of assumptions suffice to guarantee the existence of a saddle point given  $X = x$  for any x.

## 5 Quadratic Loss: Linear Regression

#### 5.a  $F_{\theta}$  derivation

Here, we find  $F_{\theta}(\mathbf{z}) = \max_{P \in \mathcal{P}_{\mathcal{Y}}} H(Y) + \mathbb{E}[\theta(Y)]^T \mathbf{z}$  for  $\theta(Y) = Y$  and  $\mathcal{P}_{\mathcal{Y}} = \{P_Y : \mathbb{E}[Y^2] \leq \theta(Y)\}$  $\rho^2$ . Since for quadratic loss  $H(Y) = \text{Var}(Y) = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$ , the problem is equivalent to

$$
F_{\boldsymbol{\theta}}(z) = \max_{\mathbb{E}[Y^2] \le \rho^2} \mathbb{E}[Y^2] - \mathbb{E}[Y]^2 + z \mathbb{E}[Y] \tag{47}
$$

As  $\mathbb{E}[Y]^2 \leq \mathbb{E}[Y^2]$ , it can be seen for the solution  $\mathbb{E}_{P^*}[Y^2] = \rho^2$  and therefore we equivalently solve

$$
F_{\theta}(z) = \max_{|\mathbb{E}[Y]| \le \rho} \rho^2 - \mathbb{E}[Y]^2 + z\mathbb{E}[Y] = \begin{cases} \rho^2 + z^2/4 & \text{if } |z/2| \le \rho \\ \rho |z| & \text{if } |z/2| > \rho. \end{cases}
$$
(48)

#### 5.b Applying Theorem 2 while restricting  $P_y$

For the quadratic loss, we first change  $\mathcal{P}_{\mathcal{Y}} = \{P_Y : \mathbb{E}[Y^2] \leq \rho^2\}$  and then apply Theorem 2. Note that by modifying  $F_{\theta}$  based on the new  $\mathcal{P}_{\mathcal{Y}}$  we also solve a modified version of the maximum conditional entropy problem

$$
\max_{\substack{P:\ P_{\mathbf{X},Y}\in\Gamma(Q) \\ \forall \mathbf{x}:\ P_{Y|\mathbf{X}=\mathbf{x}}\in\mathcal{P}_{\mathcal{Y}}}} H(Y|\mathbf{X})\tag{49}
$$

In the case  $P_{\mathcal{Y}} = \{P_Y : \mathbb{E}[Y^2] \leq \rho^2\}$  Theorem 2 remains valid given the above modification in the maximum conditional entropy problem. This is because the inequality constraint  $\mathbb{E}[Y^2|\mathbf{X}=\mathbf{x}] \leq \rho^2$ is linear in  $P_{Y|X=x}$ , and thus the problem is still convex and strong duality holds as well. Also, when we move the constraints of  $w_i = \mathbb{E}_P [\theta_i(Y)X]$  to the objective function, we get a similar dual problem

$$
\min_{\mathbf{A}} \sup_{P_{Y|\mathbf{X},\mathbf{w}:}} \mathbb{E}_{Q_{\mathbf{X}}}\left[H_{P}(Y|\mathbf{X}=\mathbf{x}) + \sum_{i=1}^{t} \mathbb{E}[\theta_{i}(Y)|\mathbf{X}=\mathbf{x}]\mathbf{A}_{i}\mathbf{X}\right] + \sum_{i=1}^{t} [I_{C_{i}}(\mathbf{w}_{i}) - \mathbf{A}_{i}\mathbf{w}_{i}]
$$
\n
$$
\forall \mathbf{x}: P_{Y|\mathbf{X}=\mathbf{x}} \in \mathcal{P}_{Y}
$$
\n(50)

Following the next steps of the proof of Theorem 2, we complete the proof assuming the modification on  $F_{\theta}$  and the maximum conditional entropy problem.

#### 5.c Derivation of group lasso

To derive the group lasso problem, we slightly change the structure of  $\Gamma(Q)$ . First assume the subsets  $I_1, \ldots, I_k$  are disjoint. Consider a set of distributions  $\Gamma_{GL}(Q)$  with the following structure

$$
\Gamma_{GL}(Q) = \{ P_{\mathbf{X},Y} : P_{\mathbf{X}} = Q_{\mathbf{X}},
$$
  
\n
$$
\forall 1 \le j \le k : \|\mathbb{E}_P \left[ Y \mathbf{X}_{I_j} \right] - \mathbb{E}_Q \left[ Y \mathbf{X}_{I_j} \right] \| \le \epsilon_j \}.
$$
\n(51)

Now we prove a modified version of Theorem 2,

<span id="page-7-3"></span>
$$
\max_{P \in \Gamma_{GL}(Q)} H(Y|\mathbf{X}) = \min_{\mathbf{\alpha}} \mathbb{E}_Q \left[ F_{\theta}(\mathbf{\alpha}^T \mathbf{X}) - Y \mathbf{\alpha}^T \mathbf{X} \right] + \sum_{j=1}^k \epsilon_j \|\mathbf{\alpha}_{I_j}\|_*.
$$
 (52)

To prove this identity, we can use the same proof provided for Theorem 2. We only need to redefine  $\tilde{\mathbf{E}}_j = \mathbb{E}_Q \left[ Y \mathbf{X}_{I_j} \right]$  and  $C_j = \{ \mathbf{u} : ||\mathbf{u} - \tilde{\mathbf{E}}_j|| \leq \epsilon_j \}$  for  $1 \leq j \leq k$ . Notice that here  $t = 1$ . Using the same technique in that proof, the dual problem can be formulated as

$$
\min_{\mathbf{\alpha}} \sup_{P_{Y|\mathbf{X}}, \mathbf{w}} \mathbb{E}_{Q_{\mathbf{X}}} \left[ H_P(Y|\mathbf{X}=\mathbf{x}) + \mathbb{E}[Y|\mathbf{X}=\mathbf{x}]\boldsymbol{\alpha}^T\mathbf{X} \right] + \sum_{j=1}^k \left[ I_{C_j}(\mathbf{w}_{I_j}) - \boldsymbol{\alpha}_{I_j}\mathbf{w}_{I_j} \right].
$$
 (53)

Similarly, we can decouple and simplify the above problem to derive the RHS of [\(52\)](#page-7-3). Then, if we let  $\|\cdot\|$  be the  $\ell_q$ -norm, we will get the group lasso problem with the  $\ell_{1,p}$  regularizer.

If the subsets are not disjoint, we can create new copies of each feature corresponding to a repeated index, such that there will be no repeated indices after adding the new features. Note that since  $P_X$  has been fixed over  $\Gamma_{GL}(Q)$  adding the extra copies of original features does not change the maximum-conditional entropy problem. Hence, we can use the result proven for the disjoint case and derive the overlapping group lasso problem.

# References

- <span id="page-7-0"></span>[1] Maurice Sion. On general minimax theorems. *Pacific J. Math*, 8(1):171–176, 1958.
- <span id="page-7-1"></span>[2] Ralph Rockafellar. Characterization of the subdifferentials of convex functions. *Pacific Journal of Mathematics*, 17(3):497–510, 1966.
- <span id="page-7-2"></span>[3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.