

## A Proof of Theorem 2

*Proof.* To bound the expectation of  $\Psi$ , we first derive an upper bound on  $\Psi$  that does not depend on  $\mathbf{w}$ . Let  $\hat{D}$  denote the empirical distribution related to the sample  $S$ . In the following, the expectations with respect to  $\hat{D}$  assume a fixed  $\mathbf{w}$ . We use directly the properties of  $\mathbf{w}$  and  $\mathbf{w}^*$  as minimizers of  $F_S$  and  $F^*$ . Writing  $\nabla F_S(\mathbf{w})=0$  and  $\nabla F^*(\mathbf{w}^*)=0$  and taking the difference yield immediately

$$\mathbf{w}^* - \mathbf{w} = -\frac{1}{2\lambda} \left[ \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w})] \right] \quad (16)$$

$$= -\frac{1}{2\lambda} \left[ \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] + \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w})] \right]. \quad (17)$$

Taking the inner product with  $(\mathbf{w}^* - \mathbf{w})$  and using the convexity of  $L_z$ , which implies  $(\mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \nabla L_z(\mathbf{w})) \cdot (\mathbf{w}^* - \mathbf{w}) \geq 0$ , lead to

$$\|\mathbf{w}^* - \mathbf{w}\|^2 \leq -\frac{1}{2\lambda} \left[ \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] \right] \cdot (\mathbf{w}^* - \mathbf{w}) \quad (18)$$

$$\leq \frac{1}{2\lambda} \left\| \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] - \mathbb{E}_{\hat{D}}[\nabla L_z(\mathbf{w}^*)] \right\| \|\mathbf{w}^* - \mathbf{w}\|. \quad (19)$$

Thus, we can write  $2\lambda\|\mathbf{w}^* - \mathbf{w}\| \leq \left\| \frac{1}{m} \sum_{i=1}^m \mathbf{Z}_i \right\|$ , where  $\mathbf{Z} = \nabla L_z(\mathbf{w}^*) - \mathbb{E}[\nabla L_z(\mathbf{w}^*)]$  and  $\mathbf{Z}_i = \nabla L_{z_i}(\mathbf{w}^*) - \mathbb{E}[\nabla L_{z_i}(\mathbf{w}^*)]$ , for all  $i \in [1, m]$ . Note that this upper bound does not depend on  $\mathbf{w}$ , which makes it easier to analyze its expectation with respect to the choice of  $S$ .

By Jensen's inequality,  $2\lambda \mathbb{E}[\Psi] \leq \mathbb{E} \left[ \left\| \frac{1}{m} \sum_{i=1}^m \mathbf{Z}_i \right\| \right] \leq \sqrt{\mathbb{E} \left[ \left\| \frac{1}{m} \sum_{i=1}^m \mathbf{Z}_i \right\|^2 \right]}$ . Using the fact that the variables  $\mathbf{Z}_i$ s are i.i.d. with  $\mathbb{E}[\mathbf{Z}_i] = 0$ , we obtain

$$\mathbb{E} \left[ \left\| \frac{1}{m} \sum_{i=1}^m \mathbf{Z}_i \right\|^2 \right] = \frac{1}{m^2} \left[ \sum_{i=1}^m \mathbb{E}[\|\mathbf{Z}_i\|^2] + \sum_{i \neq j} \mathbb{E}[\mathbf{Z}_i] \cdot \mathbb{E}[\mathbf{Z}_j] \right] = \frac{1}{m} \mathbb{E}[\|\mathbf{Z}_1\|^2] = \frac{1}{m} \text{Var}(\mathbf{Z}_1).$$

Using the expression of  $\nabla L_z(\mathbf{w}^*)$  already derived in the proof of Theorem 1 and the elementary fact that if  $\mathbf{Z}_1$  and  $\mathbf{Z}_2$  are independent and identically distributed, then  $\text{Var}(\mathbf{Z}_1) = 1/2 \mathbb{E}[(\mathbf{Z}_1 - \mathbf{Z}_2)^2]$ , this shows that  $\mathbb{E}[\Psi] \leq \frac{1}{2\lambda} \sqrt{\frac{1}{2} \frac{(4R)^2}{m}} = \frac{2R}{\lambda\sqrt{2m}}$ .  $\square$