

Supplementary Information - Proofs of Theorems

Appendix A The adversarial distribution in RSVM₂(σ)

In general there is no probability distribution $p(\bar{\mathbf{x}}|\mathbf{x})$ that attains the optimal value of the adversarial maximization problem (Equation 6 in the paper). However, the optimum is achieved as a limit over a sequence of distributions as described next. We shall assume *wlog* that the smallest ball that contains the set of vectors \mathbf{w} is centered at zero and has radius $r = \max_{\bar{y}} \|\mathbf{w}_{\bar{y}}\|$ (Every set \mathbf{w} can be shifted by translation to such a set, and the construction can be extended accordingly). In other words, we can assume $0 \in \text{conv}\{\mathbf{w}_{\bar{y}} : \|\mathbf{w}_{\bar{y}}\|_2 = r\}$, with $r = \max_{\bar{y}} \|\mathbf{w}_{\bar{y}}\|$. There are $\lambda_{\bar{y}}$ that satisfies.

$$\sum \lambda_{\bar{y}} \mathbf{w}_{\bar{y}} = 0$$

with $\sum \lambda_{\bar{y}} = 1$ and $\lambda_{\bar{y}} \geq 0$. Additionally $\lambda_{\bar{y}} \neq 0$ iff $\|\mathbf{w}_{\bar{y}}\|_2 = r$. We proceed to define a set of distributions $p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})$, each parameterized by a value $\gamma > 0$. The distribution $p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})$ has non zero mass only on $L + 1$ points, and is defined as follows:

$$\begin{aligned} p_{\gamma}(\mathbf{x}|\mathbf{x}) &= 1 - \gamma \\ p_{\gamma}\left(\sigma \frac{\mathbf{w}_{\bar{y}}}{r\gamma} + \mathbf{x}|\mathbf{x}\right) &= \lambda_{\bar{y}}\gamma, \quad \forall \bar{y} \end{aligned}$$

It is easy to see that $p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})$ is a valid distribution. To see that $p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})$ satisfies the constraints in \mathcal{S}_{ℓ_2} note that:

$$E_{p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})}(\bar{\mathbf{x}}) = (1 - \gamma)\mathbf{x} + \sum_{\lambda_{\bar{y}} \neq 0} \sigma \lambda_{\bar{y}} \frac{\mathbf{w}_{\bar{y}}}{r} + \lambda_{\bar{y}}\gamma \mathbf{x} = \mathbf{x}$$

$$E_{p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})}(\|\mathbf{x} - \bar{\mathbf{x}}\|_2) = \sum_{\lambda_{\bar{y}} \neq 0} \lambda_{\bar{y}} \sigma \gamma \frac{\|\mathbf{w}_{\bar{y}}\|_2}{r\gamma} = \sigma$$

Finally, we want to show that as $\gamma \rightarrow 0$ we obtain the optimal value of the adversarial problem. To show this, note that when γ is sufficiently small the loss $\ell\left(\frac{\sigma \mathbf{w}_{\bar{y}}}{r\gamma} + \mathbf{x}, y; \mathbf{w}\right)$ is given by

$$\ell\left(\frac{\mathbf{w}_{\bar{y}}}{r\gamma} + \mathbf{x}, y; \mathbf{w}\right) = e_{\bar{y}, y} + \Delta \mathbf{w}_{\bar{y}}^T \left(\sigma \frac{\mathbf{w}_{\bar{y}}}{r\gamma} + \mathbf{x}\right)$$

We can now write the loss corresponding to $p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})$ and take the limit $\gamma \rightarrow 0$.

$$E_{p_{\gamma}(\bar{\mathbf{x}}|\mathbf{x})}(\ell(\bar{\mathbf{x}}, \bar{y}; \mathbf{w})) =$$

$$\begin{aligned} &(1-\gamma)\ell(\mathbf{x}, \bar{y}; \mathbf{w}) + \gamma \sum_{\lambda_{\bar{y}} \neq 0} \lambda_{\bar{y}} \left(e_{\bar{y}, y} + \Delta \mathbf{w}_{\bar{y}}^T \left(\sigma \frac{\mathbf{w}_{\bar{y}}}{r\gamma} + \mathbf{x} \right) \right) \\ &\xrightarrow{\gamma \rightarrow 0} \ell(\bar{\mathbf{x}}, \bar{y}; \mathbf{w}) + \sigma \sum_{\lambda_{\bar{y}} \neq 0} \lambda_{\bar{y}} \Delta \mathbf{w}_{\bar{y}}^T \frac{\mathbf{w}_{\bar{y}}}{r} = \ell(\mathbf{x}, \bar{y}; \mathbf{w}) + \sigma \|\mathbf{w}_y\|_2 \end{aligned}$$

We obtained the optimum value of the adversarial maximization problem (see Theorem 3.1), and thus the limit of p_{γ} corresponds to the optimal adversary.

In case \mathbf{w} is not centered around zero, note that by our proof the optimal adversary is given by

$$E_{p(\bar{\mathbf{x}}|\mathbf{x})}(\ell(\mathbf{x}, y; \mathbf{w})) = \sigma \min_{\beta} \max_{\bar{y}} \|\mathbf{w}_{\bar{y}} - \beta\| + \ell(\mathbf{x}, y; \mathbf{w})$$

The β that solves this optimization problem, will be the center of the smallest ball containing the set of vectors $\mathbf{w}_{\bar{y}}$. The construction of p_{γ} is similar. Note also that the optimization problem we are considering (Eq. 13) will converge to a set \mathbf{w} , centered at zero.

Appendix B

Appendix B.1 Proof of Theorem 3.5

Theorem. $RSVM_2^2(\sigma)$ is equivalent to the problem

$$\begin{aligned} &\min_{\mathbf{w}_y; \alpha_i; \beta_i; \gamma_i} \frac{1}{n} \sum_i \alpha_i \sigma + \alpha_i \|\mathbf{x}_i\|^2 + \mathbf{x}_i^T \beta_i + \gamma_i \quad \text{s.t.} \\ &\forall \bar{y} \left[\begin{array}{cc} \alpha Id & \frac{1}{2} (\beta_i - \Delta^i \mathbf{w}_{\bar{y}}) \\ \frac{1}{2} (\beta_i - \Delta^i \mathbf{w}_{\bar{y}})^T & \gamma_i - e_{y_i, \bar{y}} \end{array} \right] \succeq 0. \end{aligned}$$

Proof. Our starting point is,

$$\begin{aligned} &\max_{p \in \mathcal{P}} E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\ell(\bar{\mathbf{x}}; y; \mathbf{w})] \\ &\text{s.t.} \quad E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\bar{\mathbf{x}}] = \mathbf{x}, \quad E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\|\bar{\mathbf{x}} - \mathbf{x}\|_2^2] = \sigma \end{aligned}$$

which is equivalent to,

$$\begin{aligned} &\max_{p \in \mathcal{P}} E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\ell(\bar{\mathbf{x}}; y; \mathbf{w})] \\ &\text{s.t.} \quad E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\bar{\mathbf{x}}] = \mathbf{x}, \quad E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\|\bar{\mathbf{x}}\|_2^2] = \sigma + \|\mathbf{x}\|_2^2 \end{aligned}$$

As before, given a labeled example $(\mathbf{x}; y)$, we define $\Delta \mathbf{w}_{\bar{y}} = \mathbf{w}_{\bar{y}} - \mathbf{w}_y$. The dual of the last problem is

$$\begin{aligned} &\min. \quad \alpha \sigma + \alpha \|\mathbf{x}\|_2^2 + \beta^T \mathbf{x} + \gamma \\ &\text{s.t.} \quad \alpha \bar{\mathbf{x}}^T \bar{\mathbf{x}} + \beta^T \bar{\mathbf{x}} + \gamma \geq e_{\bar{y}, y} + \Delta \mathbf{w}_{\bar{y}}^T \bar{\mathbf{x}} \quad \forall \bar{y} \forall \bar{\mathbf{x}} \end{aligned}$$

In the above, each constraint is quadratic in $\bar{\mathbf{x}}$, where α is the coefficient of the quadratic term. We note that $\alpha > 0$, since otherwise, the constraints will be violated. Hence we replace the infinitely many constraints with a constraint on the point that achieves the minimum value and get the equivalent problem:

$$\begin{aligned} &\min. \quad \alpha \sigma + \alpha \|\mathbf{x}\|_2^2 + \beta^T \mathbf{x} + \gamma \\ &\text{s.t.} \quad \gamma - e_{y, \bar{y}} - \frac{(\beta - \Delta \mathbf{w}_{\bar{y}})^T (\beta - \Delta \mathbf{w}_{\bar{y}})}{4\alpha} \geq 0, \quad \forall \bar{y} \end{aligned}$$

Moving to the Schur complement we obtain the following problem,

$$\begin{aligned} \min. \quad & \alpha\sigma + \alpha\|\mathbf{x}\|^2 + \boldsymbol{\beta}^T \mathbf{x} + \gamma \\ \text{s.t.} \quad & \begin{bmatrix} \alpha Id & \frac{1}{2}(\boldsymbol{\beta} - \Delta \mathbf{w}_{\bar{y}}) \\ \frac{1}{2}(\boldsymbol{\beta} - \Delta \mathbf{w}_{\bar{y}})^T & \gamma - e_{y,\bar{y}} \end{bmatrix} \succeq 0, \forall \bar{y}. \end{aligned} \quad (1)$$

and the result is immediate. \square

Note that the dual of Eq. 1 is

$$\begin{aligned} \max. \quad & \sum_{\bar{y}} \Delta \mathbf{w}_{\bar{y}} b_{\bar{y}} + \sum_{y \neq \bar{y}} c_{\bar{y}} \\ \text{s.t.} \quad & \sum_y \begin{bmatrix} a_{\bar{y}} Id & b_{\bar{y}} \\ b_{\bar{y}}^T & c_{\bar{y}} \end{bmatrix} = \begin{bmatrix} \sigma + \|\mathbf{x}\|^2 & \mathbf{x} \\ \mathbf{x}^T & 1 \end{bmatrix} \\ & \begin{bmatrix} a_{\bar{y}} Id & b_{\bar{y}} \\ b_{\bar{y}}^T & c_{\bar{y}} \end{bmatrix} \succeq 0 \end{aligned}$$

Conceptually, we construct a probability distribution over the labels $p(\cdot|\mathbf{x})$ such that for each point $\frac{1}{c_y} b_y$ we give probability c_y . (The positivity constraint implies that $c_y = 0$ will entail $b_y = 0$ hence there is no problem in dividing by c_y). The constraints then can be interpreted as

$$E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\bar{\mathbf{x}}] = \mathbf{x} \text{ and } E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\|\bar{\mathbf{x}} - \mathbf{x}\|^2] \leq \sigma.$$

The expected loss $E_{p(\bar{\mathbf{x}}|\mathbf{x})}[\ell(\bar{\mathbf{x}}; \mathbf{y}; \mathbf{x})]$ has exactly the optimal value of the problem. Furthermore, the optimal value is an upper bound on the expected loss for a probability that satisfies these constraints. Hence p is the desired probability. At first sight, it seems that this optimization problem requires us to solve for each example point \mathbf{x} an SDP with complexity that scales with the dimension of \mathbf{x} . In fact, the complexity of each SDP problem (i.e. minimization of $\{\alpha_i, \beta_i, \gamma_i\}$ for a given \mathbf{w}) can be reduced to scale with the number of classes. Intuitively, this follows from the fact that there is no point in putting adversarial noise on the space orthogonal to the space spanned by \mathbf{w} , hence the adversarial problem can be solved in that space.

Appendix B.2 Proof of Theorem 3.6

Theorem. *If $y \in \{1, -1\}$ is binary, $RSVM_2^2(\sigma)$ is equivalent to the problem*

$$\min_{\mathbf{w}} \frac{1}{n} \sum_i \frac{\sqrt{\sigma \|\mathbf{w}\|^2 + (1 - y\mathbf{w}^T \mathbf{x})^2} + (1 - y\mathbf{w}^T \mathbf{x})}{2} \quad (2)$$

Proof. For simplicity we let $y \in \{1, -1\}$, and as noted above, we also have $\frac{1}{2}\mathbf{w}_1 = -\frac{1}{2}\mathbf{w}_{-1} = \mathbf{w}$. The problem

of Eq. 1 becomes,

$$\begin{aligned} \min. \quad & \alpha\sigma + \alpha\|\mathbf{x}\|^2 + \boldsymbol{\beta}^T \mathbf{x} + \gamma \\ \text{s.t.} \quad & \begin{bmatrix} \alpha Id & \frac{1}{2}(\boldsymbol{\beta} + y\mathbf{w}) \\ \frac{1}{2}(\boldsymbol{\beta} + y\mathbf{w})^T & \gamma - 1 \end{bmatrix} \succeq 0 \\ & \begin{bmatrix} \alpha Id & \frac{1}{2}\boldsymbol{\beta} \\ \frac{1}{2}\boldsymbol{\beta}^T & \gamma \end{bmatrix} \succeq 0 \end{aligned}$$

As noted above, the variables of the dual problem belong to the simplex (and thus define a probability distribution). Since we assume $\sigma > 0$ this probability measure can not be degenerate, hence the dual variables are not zero. By complementary slackness, at the optimal values, the matrices used in the constraints of the last problem cannot be of a full rank. Thus, using the Schur complement again we obtain,

$$\begin{aligned} \gamma - \frac{1}{4\alpha} \boldsymbol{\beta}^T \boldsymbol{\beta} &= 0 \\ \gamma - 1 - \frac{1}{4\alpha} (\boldsymbol{\beta} + y\mathbf{w})^T (\boldsymbol{\beta} + y\mathbf{w}) &= 0. \end{aligned}$$

We rewrite the problem and get,

$$\begin{aligned} \min. \quad & \alpha \left(\sigma + \|\mathbf{x} + \frac{1}{2\alpha} \boldsymbol{\beta}\|^2 \right) \\ \text{s.t.} \quad & \alpha = -\frac{2y\boldsymbol{\beta}^T \mathbf{w} + \|\mathbf{w}\|^2}{4}, \end{aligned}$$

which is equivalent to the problem

$$\begin{aligned} \min. \quad & \alpha \left(\sigma + \frac{(w^T \mathbf{x} + \frac{1}{2\alpha} \boldsymbol{\beta}^T \mathbf{w})^2}{\|\mathbf{w}\|^2} \right) \\ \text{s.t.} \quad & y\boldsymbol{\beta}^T \mathbf{w} = -2\alpha + \frac{1}{2} \|\mathbf{w}\|^2 \end{aligned}$$

We plug the value of $y\boldsymbol{\beta}^T \mathbf{w}$ into the objective and get,

$$\min. \frac{\alpha}{\|\mathbf{w}\|^2} \left(\|\mathbf{w}\|^2 \sigma + \left(y\mathbf{w}^T \mathbf{x} - 1 + \frac{1}{4\alpha} \|\mathbf{w}\|^2 \right)^2 \right). \quad (3)$$

Setting to zero the derivative of the last problem with respect to α we get,

$$\|\mathbf{w}\|^2 \sigma + (y\mathbf{w}^T \mathbf{x} - 1)^2 - \left(\frac{1}{4\alpha} \|\mathbf{w}\|^2 \right)^2 = 0$$

Yielding,

$$\frac{\|\mathbf{w}\|^2}{\alpha} = 4\sqrt{\|\mathbf{w}\|^2 \sigma + (1 - y\mathbf{w}^T \mathbf{x})^2}.$$

Plugging the last result back into Eq. 3 yields the desired result. \square