

A. Proofs

A.1. Proof of Theorem 1

If the composite loss $\tilde{\ell}(z)$ is convex, it is linear.

Proof: The composite loss is an odd function:

$$\tilde{\ell}(-z) = \ell(-z) - \ell(z) = -\tilde{\ell}(z),$$

Therefore, $\frac{d^2}{dz^2}\tilde{\ell}(z) = -\frac{d^2}{dz^2}\tilde{\ell}(-z)$. If the composite loss $\tilde{\ell}(z)$ is convex, $\frac{d^2}{dz^2}\tilde{\ell}(z) \geq 0$ holds for all z . Since the convexity of $\tilde{\ell}(z)$ implies the convexity of $\tilde{\ell}(-z)$, $\frac{d^2}{dz^2}\tilde{\ell}(-z) \geq 0$ should also hold for all z . However, if $\frac{d^2}{dz^2}\tilde{\ell}(z) > 0$, then $\frac{d^2}{dz^2}\tilde{\ell}(-z) < 0$ holds, which is contradictory to the convexity of $\tilde{\ell}(-z)$. Therefore, $\frac{d^2}{dz^2}\tilde{\ell}(z) = 0$ should hold, which is satisfied only when $\tilde{\ell}(z)$ is linear. \square

A.2. Proof of Lemma 2

$J_S(\alpha)$ is strongly convex in α with parameter at least λ , and thus

$$\begin{aligned} J_S(\alpha) &\geq J_S(\alpha_S^*) + \nabla J_S(\alpha_S^*)^\top (\alpha - \alpha_S^*) + \lambda \|\alpha - \alpha_S^*\|_2^2 \\ &\geq J_S(\alpha_S^*) + \lambda \|\alpha - \alpha_S^*\|_2^2, \end{aligned}$$

where we use the optimality condition $\nabla J_S(\alpha_S^*) = \mathbf{0}$. Similarly, we can prove the other two inequalities. \square

A.3. Proof of Lemma 3

The difference function can be written as

$$J_S(\alpha, \mathbf{u}) - J_S(\alpha) = \frac{1}{4} \alpha^\top \mathbf{u}_1 \alpha + \frac{1}{2} \mathbf{u}_2^\top \alpha - \pi \mathbf{u}_3^\top \alpha,$$

with a partial gradient

$$\frac{\partial}{\partial \alpha} (J_S(\alpha, \mathbf{u}) - J_S(\alpha)) = \frac{1}{2} \mathbf{u}_1 \alpha + \frac{1}{2} \mathbf{u}_2 - \pi \mathbf{u}_3.$$

Given the δ -ball of α_S^* , i.e., $B_\delta(\alpha_S^*) = \{\alpha \mid \|\alpha - \alpha_S^*\|_2 \leq \delta\}$, it is easy to see that for any $\alpha \in B_\delta(\alpha_S^*)$,

$$\|\alpha\|_2 \leq \|\alpha - \alpha_S^*\|_2 + \|\alpha_S^*\|_2 \leq 1 + M_\alpha,$$

and then

$$\left\| \frac{\partial}{\partial \alpha} (J_S(\alpha, \mathbf{u}) - J_S(\alpha)) \right\|_2 \leq \frac{1}{2} (1 + M_\alpha) \|\mathbf{u}_1\|_{\text{Fro}} + \frac{1}{2} \|\mathbf{u}_2\|_2 + \pi \|\mathbf{u}_3\|_2.$$

This means that $J_S(\cdot, \mathbf{u}) - J_S(\cdot)$ is Lipschitz continuous on $B_\delta(\alpha_S^*)$ with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}_1\|_{\text{Fro}} + \|\mathbf{u}_2\|_2 + \|\mathbf{u}_3\|_2)$. \square

A.4. Proof of Lemma 5

The difference function can be written as

$$J_{LL}(\alpha, \mathbf{u}) - J_{LL}(\alpha) = -\pi \mathbf{u}_3^\top \alpha + u_4(\alpha).$$

Given $\alpha \in B_\delta(\alpha_{LL}^*)$, we have known that $-\pi \mathbf{u}_3^\top \alpha$ is Lipschitz continuous with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}_3\|_2)$ in the proof of Lemma 3. Consequently, $J_{LL}(\cdot, \mathbf{u}) - J_{LL}(\cdot)$ is Lipschitz continuous on $B_\delta(\alpha_{LL}^*)$ with a Lipschitz constant of order $\mathcal{O}(\|\mathbf{u}_3\|_2 + \text{Lip}(u_4))$. \square

A.5. Proof of Lemma 7

Same as the proof of Lemma 5. \square

A.6. Proof of Theorem 4

Let \mathbf{u}_1 , \mathbf{u}_2 and \mathbf{u}_3 be defined as in Eq. (13). According to the *central limit theorem*,

$$\|\mathbf{u}_1\|_{\text{Fro}} = \mathcal{O}_p(n'^{-1/2}), \quad \|\mathbf{u}_2\|_2 = \mathcal{O}_p(n'^{-1/2}), \quad \|\mathbf{u}_3\|_2 = \mathcal{O}_p(n^{-1/2}),$$

as $n, n' \rightarrow \infty$. Thus, we have

$$\begin{aligned} \|\widehat{\boldsymbol{\alpha}}_S - \boldsymbol{\alpha}_S^*\|_2 &\leq \lambda^{-1}\omega(\mathbf{u}) \\ &= \mathcal{O}(\|\mathbf{u}_1\|_{\text{Fro}} + \|\mathbf{u}_2\|_2 + \|\mathbf{u}_3\|_2) \\ &= \mathcal{O}_p(n^{-1/2} + n'^{-1/2}) \end{aligned}$$

by Lemma 2, Lemma 3, and Proposition 6.1 in Bonnans & Shapiro (1998, p. 19).

On the other hand,

$$|\widehat{J}_S(\widehat{\boldsymbol{\alpha}}_S) - J_S(\boldsymbol{\alpha}_S^*)| \leq |\widehat{J}_S(\widehat{\boldsymbol{\alpha}}_S) - \widehat{J}_S(\boldsymbol{\alpha}_S^*)| + |\widehat{J}_S(\boldsymbol{\alpha}_S^*) - J_S(\boldsymbol{\alpha}_S^*)|,$$

in which

$$\begin{aligned} \widehat{J}_S(\widehat{\boldsymbol{\alpha}}_S) - \widehat{J}_S(\boldsymbol{\alpha}_S^*) &= (\widehat{\boldsymbol{\alpha}}_S + \boldsymbol{\alpha}_S^*)^\top \left(\frac{1}{4n'} \sum_{i=1}^{n'} \boldsymbol{\varphi}(\mathbf{x}'_i) \boldsymbol{\varphi}(\mathbf{x}'_i)^\top + \frac{\lambda}{2} I_m \right) (\widehat{\boldsymbol{\alpha}}_S - \boldsymbol{\alpha}_S^*) \\ &\quad + \left(\frac{1}{2n'} \sum_{i=1}^{n'} \boldsymbol{\varphi}(\mathbf{x}'_i) \right)^\top (\widehat{\boldsymbol{\alpha}}_S - \boldsymbol{\alpha}_S^*) - \pi \left(\frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}(\mathbf{x}_i) \right)^\top (\widehat{\boldsymbol{\alpha}}_S - \boldsymbol{\alpha}_S^*), \\ \widehat{J}_S(\boldsymbol{\alpha}_S^*) - J_S(\boldsymbol{\alpha}_S^*) &= \frac{1}{4} \boldsymbol{\alpha}_S^{*\top} \mathbf{u}_1 \boldsymbol{\alpha}_S^* + \frac{1}{2} \mathbf{u}_2 \boldsymbol{\alpha}_S^* - \pi \mathbf{u}_3 \boldsymbol{\alpha}_S^*. \end{aligned}$$

Since $0 \leq \varphi_j(\mathbf{x}) \leq 1$, $\|\boldsymbol{\alpha}_S^*\|_2 \leq M_\alpha$ and $\|\widehat{\boldsymbol{\alpha}}_S\|_2 \leq M_\alpha$,

$$\begin{aligned} |\widehat{J}_S(\widehat{\boldsymbol{\alpha}}_S) - J_S(\boldsymbol{\alpha}_S^*)| &\leq |\widehat{J}_S(\widehat{\boldsymbol{\alpha}}_S) - \widehat{J}_S(\boldsymbol{\alpha}_S^*)| + |\widehat{J}_S(\boldsymbol{\alpha}_S^*) - J_S(\boldsymbol{\alpha}_S^*)| \\ &\leq \mathcal{O}_p(\|\widehat{\boldsymbol{\alpha}}_S - \boldsymbol{\alpha}_S^*\|_2) + \mathcal{O}_p(\|\mathbf{u}_1\|_{\text{Fro}} + \|\mathbf{u}_2\|_2 + \|\mathbf{u}_3\|_2) \\ &= \mathcal{O}_p(n^{-1/2} + n'^{-1/2}), \end{aligned}$$

which completes the proof. □

A.7. Proof of Theorem 6

Let \mathbf{u}_3 and $u_4(\boldsymbol{\alpha})$ be defined as in Eq. (14). The gradient of u_4 is given by

$$\nabla u_4(\boldsymbol{\alpha}) = \frac{1}{n'} \sum_{i=1}^{n'} \frac{\boldsymbol{\varphi}(\mathbf{x}'_i)}{1 + \exp(-\boldsymbol{\varphi}(\mathbf{x}'_i)^\top \boldsymbol{\alpha})} - \int \frac{\boldsymbol{\varphi}(\mathbf{x})}{1 + \exp(-\boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha})} p(\mathbf{x}) d\mathbf{x}.$$

According to the central limit theorem,

$$\|\mathbf{u}_3\|_2 = \mathcal{O}_p(n^{-1/2}), \quad \text{Lip}(u_4) = \mathcal{O}_p(n'^{-1/2}),$$

as $n, n' \rightarrow \infty$, since $\text{Lip}(u_4) = \sup_{\boldsymbol{\alpha}} \|\nabla u_4(\boldsymbol{\alpha})\|_2$ and

$$\sup_{\boldsymbol{\alpha} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^d} \left\| \frac{\boldsymbol{\varphi}(\mathbf{x})}{1 + \exp(-\boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha})} \right\|_2 \leq m^{1/2} < \infty.$$

Thus, we have

$$\begin{aligned} \|\widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\alpha}_{\text{LL}}^*\|_2 &\leq \lambda^{-1}\omega(\mathbf{u}) \\ &= \mathcal{O}(\|\mathbf{u}_3\|_2 + \text{Lip}(u_4)) \end{aligned}$$

$$= \mathcal{O}_p(n^{-1/2} + n'^{-1/2})$$

by Lemma 2, Lemma 5, and Proposition 6.1 in Bonnans & Shapiro (1998, p. 19).

On the other hand,

$$|\widehat{J}_{\text{LL}}(\widehat{\boldsymbol{\alpha}}_{\text{LL}}) - J_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)| \leq |\widehat{J}_{\text{LL}}(\widehat{\boldsymbol{\alpha}}_{\text{LL}}) - \widehat{J}_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)| + |\widehat{J}_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*) - J_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)|.$$

For the second term,

$$\begin{aligned} |\widehat{J}_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*) - J_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)| &= |-\pi \mathbf{u}_3^\top \boldsymbol{\alpha}_{\text{LL}}^* + u_4(\boldsymbol{\alpha}_{\text{LL}}^*)| \\ &\leq \pi M_\alpha \|\mathbf{u}_3\|_2 + |u_4(\boldsymbol{\alpha}_{\text{LL}}^*)| \\ &= \mathcal{O}_p(n^{-1/2} + n'^{-1/2}) \end{aligned}$$

according to the central limit theorem. For the first term, it is a bit more complex:

$$\begin{aligned} |\widehat{J}_{\text{LL}}(\widehat{\boldsymbol{\alpha}}_{\text{LL}}) - \widehat{J}_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)| &\leq \left| \frac{\lambda}{2} (\widehat{\boldsymbol{\alpha}}_{\text{LL}} + \boldsymbol{\alpha}_{\text{LL}}^*)^\top (\widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\alpha}_{\text{LL}}^*) \right| + \left| \pi \left(\frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}(\mathbf{x}_i) \right)^\top (\widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\alpha}_{\text{LL}}^*) \right| \\ &\quad + \frac{1}{n'} \sum_{i=1}^{n'} |\ln(1 + \exp(\boldsymbol{\varphi}(\mathbf{x}'_i)^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}})) - \ln(1 + \exp(\boldsymbol{\varphi}(\mathbf{x}'_i)^\top \boldsymbol{\alpha}_{\text{LL}}^*))|. \end{aligned}$$

Let $f(z, t) = \ln(1 + \exp(z + t))$, then $\lim_{t \rightarrow 0} f(z, t) = f(z, 0)$ and

$$\lim_{t \rightarrow 0} \frac{f(z, t) - f(z, 0)}{t} = \lim_{t \rightarrow 0} \frac{\partial}{\partial t} f(z, t) = \frac{1}{1 + \exp(-z - t)} < \infty,$$

where we use *L'Hôpital's rule*. In other words, $f(z, t)$ approaches $f(z, 0)$ in $\mathcal{O}(t)$ as $t \rightarrow 0$. Subsequently, for any $\mathbf{x} \in \mathbb{R}^d$, by $z = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*$ and $t = \boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*$ we can obtain

$$\begin{aligned} |\ln(1 + \exp(\boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}})) - \ln(1 + \exp(\boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*))| &= \mathcal{O}(|\boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*|) \\ &= \mathcal{O}(m^{1/2} \|\widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\alpha}_{\text{LL}}^*\|_2), \end{aligned}$$

which results in $|\widehat{J}_{\text{LL}}(\widehat{\boldsymbol{\alpha}}_{\text{LL}}) - \widehat{J}_{\text{LL}}(\boldsymbol{\alpha}_{\text{LL}}^*)| = \mathcal{O}_p(n^{-1/2} + n'^{-1/2})$. \square

A.8. Proof of Theorem 8

The proof goes along the same line as that of Theorem 6. Let \mathbf{u}_3 and $u_5(\boldsymbol{\alpha})$ be defined as in Eq. (15). Note that the function $\max\{0, (1+z)/2, z\}$ is piecewise linear in z , differentiable almost everywhere, and $0 \leq (d/dz) \max\{0, (1+z)/2, z\} \leq 1$. As a result,

$$\|\mathbf{u}_3\|_2 = \mathcal{O}_p(n^{-1/2}), \quad \text{Lip}(u_5) = \mathcal{O}_p(n'^{-1/2}),$$

as $n, n' \rightarrow \infty$, and

$$\begin{aligned} \|\widehat{\boldsymbol{\alpha}}_{\text{DH}} - \boldsymbol{\alpha}_{\text{DH}}^*\|_2 &\leq \lambda^{-1} \omega(\mathbf{u}) \\ &= \mathcal{O}(\|\mathbf{u}_3\|_2 + \text{Lip}(u_5)) \\ &= \mathcal{O}_p(n^{-1/2} + n'^{-1/2}) \end{aligned}$$

by Lemma 2, Lemma 7, and Proposition 6.1 in Bonnans & Shapiro (1998, p. 19).

On the other hand,

$$\begin{aligned} |\widehat{J}_{\text{DH}}(\widehat{\boldsymbol{\alpha}}_{\text{DH}}) - J_{\text{DH}}(\boldsymbol{\alpha}_{\text{DH}}^*)| &\leq |\widehat{J}_{\text{DH}}(\widehat{\boldsymbol{\alpha}}_{\text{DH}}) - \widehat{J}_{\text{DH}}(\boldsymbol{\alpha}_{\text{DH}}^*)| + |\widehat{J}_{\text{DH}}(\boldsymbol{\alpha}_{\text{DH}}^*) - J_{\text{DH}}(\boldsymbol{\alpha}_{\text{DH}}^*)| \\ &\leq \frac{1}{n'} \sum_{i=1}^{n'} |\max\{0, (1 + \boldsymbol{\varphi}(\mathbf{x}'_i)^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}})/2, \boldsymbol{\varphi}(\mathbf{x}'_i)^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}}\}| \end{aligned}$$

$$- \max\{0, (1 + \boldsymbol{\varphi}(\mathbf{x}'_i)^\top \boldsymbol{\alpha}_{\text{LL}}^*)/2, \boldsymbol{\varphi}(\mathbf{x}'_i)^\top \boldsymbol{\alpha}_{\text{LL}}^*\} + \mathcal{O}_p(n^{-1/2} + n'^{-1/2}).$$

Let $f(z, t) = \max\{0, (1 + z + t)/2, z + t\}$, then $\lim_{t \rightarrow 0} f(z, t) = f(z, 0)$ and for $z \in \mathbb{R} \setminus \{0, 1\}$,

$$\lim_{t \rightarrow 0} \frac{f(z, t) - f(z, 0)}{t} = \lim_{t \rightarrow 0} \frac{\partial}{\partial t} f(z, t) \in \left\{0, \frac{1}{2}, 1\right\}.$$

In other words, $f(z, t)$ approaches $f(z, 0)$ in $\mathcal{O}(t)$ as $t \rightarrow 0$ almost surely. Subsequently, for any $\mathbf{x} \in \mathbb{R}^d$, by $z = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{DH}}^*$ and $t = \boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{DH}} - \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{DH}}^*$ we can obtain

$$\begin{aligned} |\max\{0, (1 + \boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}})/2, \boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}}\} - \max\{0, (1 + \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*)/2, \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*\}| &= \mathcal{O}(|\boldsymbol{\varphi}(\mathbf{x})^\top \widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\alpha}_{\text{LL}}^*|) \\ &= \mathcal{O}(m^{1/2} \|\widehat{\boldsymbol{\alpha}}_{\text{LL}} - \boldsymbol{\alpha}_{\text{LL}}^*\|_2) \\ &= \mathcal{O}_p(n^{-1/2} + n'^{-1/2}), \end{aligned}$$

which completes the proof. \square

B. Optimization problems

In this section, we give exact optimization problems for the optimization methods presented in the paper. The logistic regression and logistic loss method is solved with a quasi-Newton method, and therefore we provide the derivatives in Sec. B.1.

The Hinge loss and Double Hinge loss result in quadratic problems. The ramp-loss is solved via a sequence of quadratic problems. All quadratic problems are expressed in the form

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \boldsymbol{\alpha}^\top H \boldsymbol{\alpha} + \mathbf{f}^\top \boldsymbol{\alpha} \\ \text{s.t.} \quad & L \boldsymbol{\alpha} \preceq \mathbf{k} \\ & \mathbf{l} \preceq \boldsymbol{\alpha} \end{aligned}$$

This standard form can then just be plugged into an off-the-shelf optimization package such as Gurobi, IBM CPLEX or MATLAB's internal 'quadprog' function.

B.1. Logistic loss

The gradient for the objective function in Eq. (8) is

$$\begin{aligned} \frac{\partial \widehat{J}_{\text{LL}}(\boldsymbol{\alpha}, b)}{\partial \boldsymbol{\alpha}} &= -\frac{\pi}{n} \Phi_{\text{p}}^\top \mathbf{1} + \lambda \boldsymbol{\alpha} \\ &\quad - \frac{1}{n'} \sum_{j=1}^{n'} \ell'_{\text{LL}}(-\boldsymbol{\alpha}^\top \boldsymbol{\varphi}(\mathbf{x}'_j) - b) \boldsymbol{\varphi}(\mathbf{x}'_j), \end{aligned}$$

where $\ell'_{\text{LL}}(z)$ is the derivative of $\ell_{\text{LL}}(z)$:

$$\ell'_{\text{LL}}(z) = -\frac{\exp(-z)}{1 + \exp(-z)}.$$

The derivative with respect to the unregularized constant b is

$$\frac{\partial \widehat{J}_{\text{LL}}(\boldsymbol{\alpha}, b)}{\partial b} = -\pi - \frac{1}{n'} \sum_{j=1}^{n'} \ell'_{\text{LL}}(-\boldsymbol{\alpha}^\top \boldsymbol{\varphi}(\mathbf{x}'_j) - b).$$

B.2. Double Hinge Loss - PU Learning

The objective function can be expressed as

$$-\frac{\pi}{n} \sum_{i=1}^n g(\mathbf{x}_i) + \frac{1}{n'} \sum_{j=1}^{n'} \max\left(0, \max\left(g(\mathbf{x}'_j), \frac{1}{2} + \frac{1}{2}g(\mathbf{x}'_j)\right)\right) + \frac{\lambda}{2} \|g\|_2^2$$

$$= -\frac{\pi}{n} \sum_{i=1}^n \left(\sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}_i) + b \right) + \frac{1}{n'} \sum_{j=1}^{n'} \max \left(0, \max \left(\sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}'_j) + b, \frac{1}{2} + \frac{1}{2} \left(\sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}'_j) + b \right) \right) \right) + \frac{\lambda}{2} \sum_{\ell=1}^m \alpha_{\ell}^2$$

The objective function can then be expressed as

$$\begin{aligned} \min_{\alpha, b, \xi} \quad & -\frac{\pi}{n} \mathbf{1}^{\top} \Phi_{\text{P}} \alpha - \pi b + \frac{1}{n'} \mathbf{1}^{\top} \xi + \frac{\lambda}{2} \alpha^{\top} \alpha \\ \text{s.t.} \quad & \xi \succeq \mathbf{0}, \\ & \xi \succeq \frac{1}{2} \mathbf{1} + \frac{1}{2} \Phi_{\text{U}} \alpha + \frac{1}{2} b \mathbf{1}, \\ & \xi \succeq \Phi_{\text{U}} \alpha + b \mathbf{1}, \end{aligned}$$

Let

$$\gamma = \begin{bmatrix} \alpha_{b \times 1} \\ b \\ \xi_{n' \times 1} \end{bmatrix}.$$

Then H is defined as

$$H = \begin{bmatrix} \lambda I_{m \times m} & O_{m \times 1} & O_{m' \times n'} \\ O_{1 \times m} & 0 & O_{1 \times n'} \\ O_{n' \times m} & O_{n' \times 1} & O_{n' \times n'} \end{bmatrix},$$

where $O_{n \times m}$ is a zero matrix of n rows and m columns. The linear part of the objective is

$$\mathbf{f} = \begin{bmatrix} -\frac{\pi}{n} \Phi_{\text{P}}^{\top} \mathbf{1} \\ -\pi \\ \frac{1}{n'} \mathbf{1}_{n' \times 1} \end{bmatrix}$$

The lower-bound is

$$\mathbf{l} = \begin{bmatrix} -\infty_{m \times 1} \\ -\infty \\ \mathbf{0}_{n' \times 1} \end{bmatrix}.$$

The first linear constraint is

$$\begin{aligned} \xi &\succeq \frac{1}{2} \mathbf{1} + \frac{1}{2} \Phi_{\text{U}} \alpha + \frac{1}{2} b \mathbf{1} \\ \frac{1}{2} \Phi_{\text{U}} \alpha + \frac{1}{2} b \mathbf{1} - \xi &\preceq -\frac{1}{2} \mathbf{1} \\ \left[\frac{1}{2} \Phi_{\text{U}} \quad \frac{1}{2} \mathbf{1}_{n' \times 1} \quad -I_{n' \times n'} \right] \begin{bmatrix} \alpha \\ u \\ \xi \end{bmatrix} &\preceq -\frac{1}{2} \mathbf{1}_{n' \times 1}. \end{aligned}$$

The second linear constraint is

$$\begin{aligned} \xi &\succeq \Phi_{\text{U}} \alpha + b \mathbf{1} \\ \Phi_{\text{U}} \alpha + b \mathbf{1} - \xi &\preceq \mathbf{0}_{n' \times 1} \\ \left[\Phi_{\text{U}} \quad \mathbf{1}_{n' \times 1} \quad -I_{n' \times n'} \right] \begin{bmatrix} \alpha \\ b \\ \xi \end{bmatrix} &\preceq \mathbf{0}_{n' \times 1}. \end{aligned}$$

Combining the two sets of inequalities, we get

$$L = \begin{bmatrix} \frac{1}{2} \Phi_{\text{U}} & \frac{1}{2} \mathbf{1}_{n' \times 1} & -I_{n' \times n'} \\ \Phi_{\text{U}} & \mathbf{1}_{n' \times 1} & -I_{n' \times n'} \end{bmatrix},$$

and

$$\mathbf{k} = \begin{bmatrix} -\frac{1}{2} \mathbf{1}_{n' \times 1} \\ \mathbf{0}_{n' \times 1} \end{bmatrix}.$$

B.3. Weighted hinge loss classifier

We want a cost-sensitive classifier with a per-sample weighting. Using the model

$$g(\mathbf{x}) = \sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}) + b,$$

where

$$\{\mathbf{c}_1, \dots, \mathbf{c}_m\} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\},$$

we wish to minimize

$$\begin{aligned} J(g) &= \frac{1}{n} \sum_{i=1}^n w_i \ell_H \left(y_i \sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}_i) + b \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha}, \\ &= \frac{1}{2n} \sum_{i=1}^n w_i \max \left(0, 1 - y_i \sum_{\ell=1}^m \alpha_{\ell} \varphi_{\ell}(\mathbf{x}_i) + b \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} \boldsymbol{\alpha}. \end{aligned}$$

This gives a QP of

$$\begin{aligned} \min_{\boldsymbol{\alpha}, b, \boldsymbol{\xi}} \quad & \frac{1}{2n} \mathbf{w}^{\top} \boldsymbol{\xi} + \frac{\lambda}{2} \boldsymbol{\alpha}^{\top} R \boldsymbol{\alpha} \\ \text{s.t.} \quad & \xi_i \geq 0, \quad \forall i = 1, \dots, n \\ & \xi_i \geq 1 - y_i \sum_{\ell=1}^m \alpha_{\ell} k(\mathbf{x}_i, \mathbf{c}_{\ell}) + u \quad \forall i = 1, \dots, n. \end{aligned}$$

We then set

$$\boldsymbol{\gamma} = \begin{bmatrix} \boldsymbol{\alpha} \\ b \\ \boldsymbol{\xi} \end{bmatrix}.$$

H is then

$$H = \begin{bmatrix} \lambda I & O_{m \times 1} & O_{m \times n} \\ O_{1 \times n} & 0 & O_{1 \times n} \\ O_{n \times n} & O_{n \times 1} & O_{n \times n} \end{bmatrix}.$$

The linear term is

$$\mathbf{f} = \begin{bmatrix} \mathbf{0}_{m \times 1} \\ 0 \\ \frac{1}{2n} \mathbf{w} \end{bmatrix}$$

The lower bound is

$$\mathbf{l} = \begin{bmatrix} -\infty_{m \times 1} \\ -\infty \\ \mathbf{0}_{n \times 1} \end{bmatrix}$$

Define $\bar{\Phi}$ as

$$\bar{\Phi}_{i\ell} = y_i \varphi_{\ell}(\mathbf{x}_i).$$

The constraint can be written in matrix form as

$$\begin{aligned} \boldsymbol{\xi} &\succeq \mathbf{1}_{n \times 1} - (\bar{\Phi} \boldsymbol{\alpha} + b \mathbf{y}) \\ -\bar{\Phi} \boldsymbol{\alpha} - b \mathbf{y} - \boldsymbol{\xi} &\preceq -\mathbf{1}_{n \times 1} \end{aligned}$$

The matrix is then

$$L = \begin{bmatrix} -\bar{\Phi} & -\mathbf{y} & -I_{n \times n} \end{bmatrix},$$

and \mathbf{k} is

$$\mathbf{k} = [-\mathbf{1}_{n \times 1}].$$

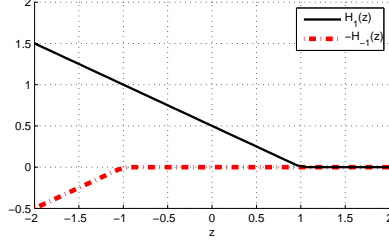


Figure 6. Decomposition of the ramp-loss into convex and concave parts.

B.4. Weighted ramp-loss classifier (CCCP)

Classification with the ramp-loss is difficult, due to the non-convexity of the loss function. One popular method to perform optimization is to split the non-convex function into a convex and concave part. The concave part is then upper-bounded by a linear function, and optimization is iteratively performed: minimization of the upper-bound, and tightening of the upper-bound around the new minima. We minimize the ramp-loss problem here using this approach. This is a straightforward application of the convex-concave procedure (CCCP) in Yuille & Rangarajan (2002) and is essentially the same as Collobert et al. (2006).

We wish to minimize the following non-convex objective function:

$$J(\boldsymbol{\alpha}, b) = \frac{1}{n} \sum_{i=1}^n w_i \ell_R \left(y_i \sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha}, \quad (16)$$

where the ramp loss $\ell_R(z)$ is defined as

$$\ell_R(z) = \max \left(0, \min \left(1, \frac{1}{2} - \frac{1}{2}z \right) \right) = \frac{1}{2} \max(0, \min(2, 1 - z)).$$

By defining the following (slightly more general) hinge loss

$$H_\epsilon(z) = \frac{1}{2} \max(0, \epsilon - z),$$

the ramp loss $\ell_R(z)$ can be decomposed as:

$$\ell_R(z) = H_1(z) - H_{-1}(z).$$

This is illustrated in Fig. 6. The objective Eq. (16) can therefore be decomposed as

$$\begin{aligned} J(\boldsymbol{\alpha}, b) &= J_{\text{vex}}(\boldsymbol{\alpha}, b) + J_{\text{cave}}(\boldsymbol{\alpha}, b), \\ J_{\text{vex}}(\boldsymbol{\alpha}, b) &= \frac{1}{n} \sum_{i=1}^n w_i H_1 \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha}, \\ J_{\text{cave}}(\boldsymbol{\alpha}, b) &= -\frac{1}{n} \sum_{i=1}^n w_i H_{-1} \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) \end{aligned}$$

The following self-evident relation can be used to upper-bound the concave part

$$\begin{aligned} tz - f(z) &\leq \sup_{y \in \mathbb{R}} yt - f(y) \\ \Rightarrow f(z) &\geq tz - f^*(t), \end{aligned} \quad (17)$$

where

$$f^*(t) = \sup_{y \in \mathbb{R}} yt - f(y).$$

The inequality in Eq.(17) is known as the *Fenchel inequality* and the function $f^*(z)$ is known as the *Fenchel dual* or *convex conjugate*. Applying the above inequality to $H_\epsilon(z)$, we can obtain a bound as

$$\begin{aligned} H_\epsilon(z) &\geq zt - H_\epsilon^*(t), \\ -H_\epsilon(z) &\leq H_\epsilon^*(t) - zt, \end{aligned}$$

where $H_\epsilon^*(t)$ is the Fenchel dual of $H_\epsilon(z)$. The Fenchel dual of $H_{-1}(t)$ is (the full calculation is given in Appendix B.4.3)

$$H_{-1}^*(t) = \begin{cases} -t & -\frac{1}{2} \leq t \leq 0, \\ \infty & \text{otherwise.} \end{cases}$$

We can minimize the upper-bound as

$$\arg \min_t H_{-1}^*(t) - tz = \begin{cases} t = 0 & z > -1. \\ t = -\frac{1}{2} & z \leq -1. \end{cases}$$

The concave part is then bounded, with the parameter \mathbf{a} as

$$\bar{J}_{\text{cave}}(\boldsymbol{\alpha}, b, \mathbf{a}) = \frac{1}{n} \sum_{i=1}^n w_i \left(H_{-1}^*(a_i) - a_i y_i \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) \right),$$

where $J_{\text{cave}}(\boldsymbol{\alpha}, u) \leq \bar{J}_{\text{cave}}(\boldsymbol{\alpha}, b, \mathbf{a})$, for any \mathbf{a} .

B.4.1. TIGHTENING OF THE UPPER-BOUND

The upperbound is minimized (tightened) when

$$a_i = \begin{cases} -\frac{1}{2} & y_i \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) \leq -1, \\ 0 & \text{otherwise.} \end{cases}$$

B.4.2. MINIMIZING THE OBJECTIVE

We wish to minimize the convex part and the upper bound $\bar{J}(\boldsymbol{\alpha}, u, \mathbf{a}) = J_{\text{vex}}(\boldsymbol{\alpha}, u) + \bar{J}_{\text{cave}}(\boldsymbol{\alpha}, u, \mathbf{a})$ with respect to \mathbf{a} . This gives an objective of

$$\bar{J}(\boldsymbol{\alpha}, b, \mathbf{a}) = \frac{1}{n} \sum_{i=1}^n w_i H_{-1} \left(y_i \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) \right) + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha} - \frac{1}{n} \sum_{i=1}^n w_i a_i y_i \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right).$$

We define the following matrices:

$$\begin{aligned} \Phi_{i,\ell} &= y_i k(\mathbf{x}_i, \mathbf{c}_\ell), \\ \bar{\Phi}_{i,\ell} &= w_i a_i y_i k(\mathbf{x}_i, \mathbf{c}_\ell), \end{aligned}$$

The QP for this is then

$$\begin{aligned} \min_{\boldsymbol{\alpha}, b, \boldsymbol{\xi}} & \frac{1}{2n} \mathbf{w}^\top \boldsymbol{\xi} + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \boldsymbol{\alpha} - \frac{1}{n} \mathbf{1}^\top \bar{\Phi} \boldsymbol{\alpha} - b \frac{1}{n} \sum_{i=1}^n w_i a_i y_i. \\ \text{s.t.} & \xi_i \geq 0 \quad \forall i = 1, \dots, n \\ & \xi_i \geq 1 - y_i \left(\sum_{\ell=1}^m \alpha_\ell \varphi_\ell(\mathbf{x}_i) + b \right) \quad \forall i = 1, \dots, n. \end{aligned}$$

We define again

$$\boldsymbol{\gamma} = \begin{bmatrix} \boldsymbol{\alpha} \\ b \\ \boldsymbol{\xi} \end{bmatrix}.$$

The quadratic term is

$$H = \begin{bmatrix} \lambda I_{m \times m} & O_{m \times 1} & O_{n \times n} \\ O_{1 \times n} & 0 & O_{1 \times n} \\ O_{n \times n} & O_{n \times 1} & O_{n \times n} \end{bmatrix}.$$

The linear term is

$$\mathbf{f} = \begin{bmatrix} -\frac{1}{n} \bar{\Phi}^\top \mathbf{1} \\ -\frac{1}{n} \sum_{i=1}^n w_i a_i y_i \\ \frac{1}{2n} \mathbf{w} \end{bmatrix}$$

The lower-bound is

$$lb = \begin{bmatrix} -\infty_{m \times 1} \\ -\infty \\ \mathbf{0}_{n \times 1} \end{bmatrix}.$$

The linear term is

$$-\Phi \alpha - b \mathbf{y} - \xi \preceq -\mathbf{1}_{n \times 1}.$$

This gives a matrix of

$$L = \begin{bmatrix} -\Phi & -\mathbf{y} & -I_{n \times n} \end{bmatrix},$$

and \mathbf{k} is

$$\mathbf{k} = [-\mathbf{1}_{n \times 1}].$$

B.4.3. CALCULATION OF THE FENCHEL DUAL OF $H_\epsilon(z)$

In this section, we briefly give the derivation of the Fenchel dual of $H_\epsilon(z)$

$$\begin{aligned} H_\epsilon^*(t) &= \sup_v tv - H_\epsilon(v) \\ &= \sup_v tv - \frac{1}{2} \max(0, \epsilon - v). \end{aligned}$$

To make the above easier, we split the domain of the v :

$$\begin{aligned} H_\epsilon^*(t) &= \max \left(\sup_{v \leq \epsilon} tv - \frac{1}{2} \max(0, \epsilon - v), \sup_{v > \epsilon} tv - \frac{1}{2} \max(0, \epsilon - v) \right), \\ &= \max \left(\sup_{v \leq \epsilon} tv - \frac{1}{2} (\epsilon - v), \sup_{v > \epsilon} tv \right). \end{aligned}$$

For the first part:

$$\begin{aligned} \sup_{v \leq \epsilon} tv - \frac{1}{2} (\epsilon - v) &= \sup_{v \leq \epsilon} v \left(t + \frac{1}{2} \right) - \frac{1}{2} \epsilon, \\ &= \begin{cases} \epsilon t & t \geq -\frac{1}{2}, \\ \infty & t < -\frac{1}{2} \end{cases} \end{aligned}$$

The second part is

$$\sup_{t > \epsilon} tv = \begin{cases} \epsilon v & t \leq 0, \\ \infty & t > 0. \end{cases}$$

Putting these two together gives:

$$H_\epsilon^*(t) = \begin{cases} \epsilon t & -\frac{1}{2} \leq t \leq 0, \\ \infty & \text{otherwise.} \end{cases}$$