# Recap of Soft-SVM

$$\min_{w,b} \quad \frac{\|w\|^2}{2} + \check{C}\|\xi\|_1$$

Subject to (A) $y_i(\langle w, x_i\rangle + b) \geq 1 - \xi_i$

(B) $\xi_i \geq 0 \qquad \forall i \in [m]$

From KKT conditions

$$w = \sum_{i=1}^{m} \alpha_i y_i x_i$$

$x_i$ : support vectors for $\alpha_i \neq 0$
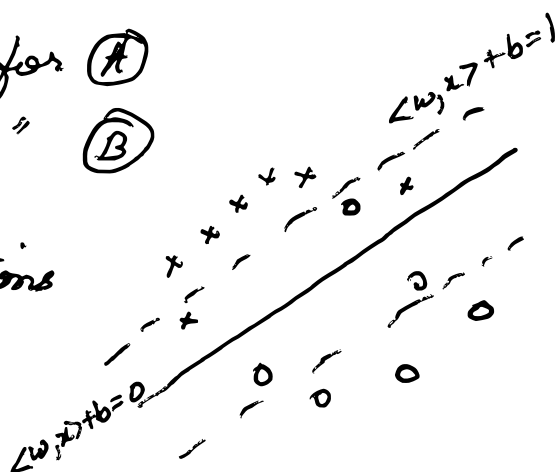
$\alpha_i$ : Dual variables for (A)

$\beta_i$ : "    "    "    (B)

Other KKT conditions

$\alpha_i + \beta_i = C$

$\alpha_i = 0$    or    $y_i(\langle w, x_i\rangle + b) = 1 - \xi_i$

$\beta_i = 0$    or    $\xi_i = 0 \leftarrow$

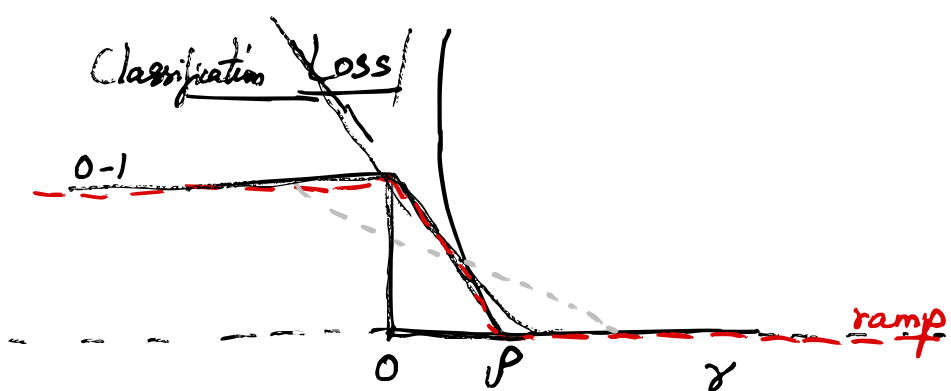$\alpha_i \neq 0 \Rightarrow y_i(\langle w, x_i\rangle + b) = 1 - \xi_i$
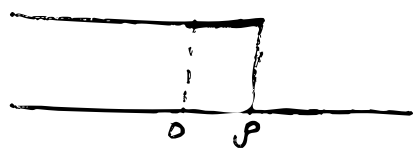
$\xi_i = 0$

$y_i(\langle w, x_i\rangle + b) = 1$

$x_i$ are on marginal hyperplane

$\xi_i > 0$

$\beta_i = 0$

$\alpha_i = C$

$y_i(\langle w, x_i\rangle + b) = 1 - \xi_i$

$\langle w, x\rangle + b = 1$

$\langle w, x\rangle + b = 0$

Classification Loss



$$\gamma(x,y,\omega,b) = y(\langle \omega, x \rangle + b)$$



$$\ell_{hinge}((x,y),(\omega,b)) = \max\{0, 1 - \frac{\gamma}{\rho}\}$$

$$\rho = 1$$

$$\ell_{quad\text{-}hinge} = (\ell_{hinge})^2$$

$$\ell_{ramp}((x,y),(\omega,b))$$

$$= \min\{1, \max\{0, 1 - \frac{\gamma}{\rho}\}\}$$

$$\mathbb{1}\{\gamma((x,y),(\omega,b)) \leq 0\}$$
$$\leq \underline{\ell_{ramp}((x,y),(\omega,b))} \leq \mathbb{1}\{\gamma((x,y),(\omega,b)) \leq \rho\}$$

$$(\rho > 0)$$

$$\hat{R}_{S,\rho}(h) = \frac{1}{m} \sum_{i=1}^{m} \ell_{ramp}(z_i, h)$$

$$z_i = (x_i, y_i)$$

$$h \in \mathcal{H} = \{ h(\cdot, \omega, b) : h(x, \omega, b) = \langle \omega, x \rangle + b \}$$

$$R_\rho(h) = \mathbb{E}_{S \sim \mathcal{D}^m} \hat{R}_{S,\rho}(h)$$

(Form of) Generalization bound
with $Pr \geq 1 - \delta$ over $S \sim \mathcal{D}^m$

$$R_\rho(h) \leq \hat{R}_{S,\rho}(h) + \sqrt{\frac{f(\mathcal{H}, \mathcal{D})}{m}} \checkmark$$
$$+ \sqrt{\frac{\log 1/\delta}{2m}}$$

$$R(h) \leq \hat{R}_{S,\rho}(h) + \text{---} \cdot$$

(generalization for 0-1 loss)

$$\leq \hat{R}_{S,\rho,hinge}(h) + \cdots$$

# Generalization bound based on Rademacher Complexity

Koltchinskii & Panchenko 2002
Bartlett and Mendelson 2002

$$R(h) \leq \hat{R}_{S,\rho}(h) + \frac{2}{\rho} \text{Rad}_S(\mathcal{H})$$
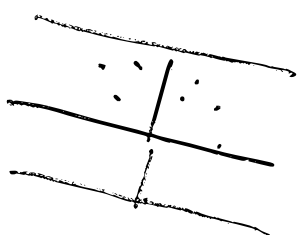$$+ 3\sqrt{\frac{\log(2/\delta)}{2m}}$$

for
$$\mathcal{H} = \{x \to \langle \omega, x \rangle : \|\omega\| < \Lambda\}$$

(0-1)  $\forall h \in \mathcal{H}$
$$R(h) \leq \hat{R}_{S,\rho}(h) + \frac{2}{\rho} \frac{\Lambda \, r}{\sqrt{m}}$$
$$+ 3\sqrt{\frac{\log 2/\delta}{2m}}$$

$\frac{\Lambda r}{\rho}$ is small and at the same time $\hat{R}_{S,\rho}(h)$ (hinge loss) is small $\Longrightarrow$ good generalization

Want:
$$Rad_S(\mathcal{H}) \leq \sqrt{\frac{r^2 \Lambda^2}{m}}$$

where $r = \sup_{x \sim \mathcal{D}} \|x\|$ ✓

$\underline{\Lambda} = \sup_{\omega} \|\omega\|$ ✓

$$Rad_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_\sigma \sup_{h \in \mathcal{H}} \sum_{i=1}^{m} \sigma_i h(x_i)$$

$\sigma = \{\sigma_1, \ldots, \sigma_m\}$ iid

$\mathcal{H} = \{h(\cdot, \omega) : h(x, \omega) = \langle \omega, x \rangle, \|\omega\| \leq \Lambda\}$

**Proof**

$$Rad_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_\sigma \sup_{\|\omega\| \leq \Lambda} \sum_{i=1}^{m} \sigma_i \langle \omega, x_i \rangle$$

$$= \frac{1}{m} \mathbb{E}_\sigma \sup_{\|\omega\| \leq \Lambda} \langle \omega, \sum_{i=1}^{m} \sigma_i x_i \rangle$$

$(|\langle a, b \rangle| \leq \|a\| \|b\|)$

$$\leq \frac{\Lambda}{m} \mathbb{E}_\sigma \| \sum_{i=1}^{m} \sigma_i x_i \| \quad \rightarrow \circledtimes$$

$\left[ \text{(Jensen's inequality: } \mathbb{E}f(X) \geq f(\mathbb{E}X) \right.$
   when $f$ is convex)

$\mathbb{E} X^2 \geq (\mathbb{E}X)^2$

$X = \| \sum_{i=1}^{m} \sigma_i x_i \|$

$\left. \left( \mathbb{E}_\sigma \| \sum_{i=1}^{m} \sigma_i x_i \|^2 \right)^{1/2} \geq \left( \mathbb{E}_\sigma \| \sum_{i=1}^{m} \sigma_i x_i \| \right) \rightarrow \circledast \right]$

$$Rad_S(\mathcal{H}) \leq \frac{\Lambda}{m} \left( \mathbb{E}_\sigma \| \sum_{i=1}^{m} \sigma_i x_i \|^2 \right)^{1/2}$$

$$= \frac{\Lambda}{m} \left( \mathbb{E}_\sigma \sum_{d, i=1}^{m} \sigma_i \sigma_j \langle x_i, x_j \rangle \right)^{1/2}$$

$\left( \sigma_i \text{ iid} \quad \mathbb{E}[\sigma_i \sigma_j] = \mathbb{E}[\sigma_i] \mathbb{E}[\sigma_j] \right.$
   $\left. = 0 \right)$

$$\leq \frac{\Lambda}{m} \left( \sum_{i=1}^{m} \|x_i\|^2 \right)^{1/2}$$

$$\leq \frac{\Lambda}{m} (r^2 m)^{1/2} = \frac{r \Lambda}{\sqrt{m}}$$

$$\boxed{Rad_S(\mathcal{H}) \leq \frac{r \Lambda}{\sqrt{m}}}$$

Rademacher complexity

$$Rad_m(\mathcal{H}) = \mathbb{E}_{S \sim \mathcal{D}^m} Rad_S(\mathcal{H})$$

Generalization of SVM:

$$R(h) \leq \hat{R}_{S,\rho}(h) + \frac{2}{\rho}\sqrt{\frac{r^2 \Lambda^2}{m}}$$

$$+ \ 3\sqrt{\frac{\log 2/\delta}{2m}}$$

for all $h \in \mathcal{H} = \{x \rightarrow \langle w, x \rangle :$
$$\|w\| \leq \Lambda \}$$

and $\|x\| < r$

with probability at least $1 - \delta$
over $S \sim D^m$

$$\frac{r\Lambda}{\rho} \qquad \text{vs} \qquad \hat{R}_{S,\rho}(h)$$

# Activity tasks

→ trains a SVM

Components
→ Training data

→ M .

→ Loss ←

→ SGD ;

→ hyperparameters
   (problem definition:
      algorithm / optimizer.)

Code

   c : loss
   lr : algorithm