Primal problem :

$$\min_{w, b, \xi} \quad \frac{1}{2} \|w\|^2 + \underline{C} \sum_{i=1}^{m} \xi_i$$

Subject to

$$y_i (w_i \cdot x_i + b) \geq 1 - \boxed{\xi_i}$$

$$\gamma_i \quad \text{and} \quad \xi_i \geq 0 \qquad i \in [m].$$

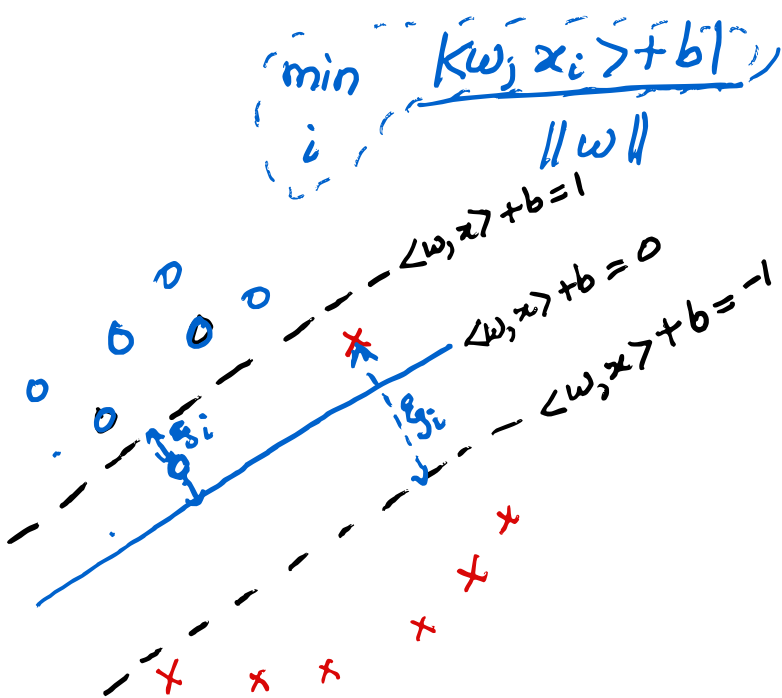$\xi = [\xi_1, ..., \xi_m]^T$

$\xi_i = 0$  HARD SVM          Soft SVM

Today

1. SVM convex optimization

2. what are support vectors ?

3. Analysis : generalization bounds
   for SVM

Material :   Ch 5 of Mohri et al
             SVM  in  Hastie, Tibshirani

Margin : $\dfrac{1}{\|w^*\|}$

$$\min_{i} \frac{|\langle w, x_i \rangle + b|}{\|w\|}$$



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

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

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

$\xi_i$

$\xi_i$

Necessary and sufficient conditions
for existence of unique
to convex optimization problems:

$w^*$ is a minimizer of
Primal problem, iff
$\rightarrow \exists \, w \in \omega \; g_i(w) \leq 0$ ( Slater's condition )
$\rightarrow \nabla_w \mathcal{L}(w^*, \alpha^*) = 0$

Complementarity constraints
$\rightarrow \quad \sum_{i=1}^{m} \alpha_i^* \, g_i(w^*) = 0$

$\Rightarrow \quad \alpha_i^* \, g_i(w^*) = 0$

Primal problem
$$\min_w \; f(w)$$
$$g_i(w) \leq 0 \qquad i \in [m]$$

$$\mathcal{L}(w, \alpha) = f(w) + \sum_{i=1}^{m} \alpha_i \, g_i(w)$$

$$\mathcal{L}(\omega, b, \xi, \alpha, \beta)$$

$$= \frac{1}{2}\|\omega\|^2 + C\sum_{i=1}^{m}\xi_i +$$

$$\sum_{i=1}^{m}\alpha_i\left(1 - \xi_i - y_i(\langle \omega_i, x_i\rangle + b)\right)$$

$$- \sum_{i=1}^{m}\beta_i\xi_i$$

$(\alpha_i, \beta_i), \ i \in M$  Dual variables

Constraints:

$\checkmark \quad y_i(\langle \omega, x_i\rangle + b) \geqslant 1 - \xi_i$

$\checkmark \quad -\xi_i \leq 0$

$$\nabla_\omega \mathcal{L} = 0 \Rightarrow \omega = \sum_{i=1}^{m}\alpha_i y_i x_i \ \checkmark$$

$$\nabla_b \mathcal{L} = 0 \qquad \sum_{i=1}^{m}\alpha_i y_i = 0 \ \checkmark$$

$$\nabla_{\xi_i}\mathcal{L} = 0 \qquad \alpha_i + \beta_i = C$$

Complementarity:

$$\left(y_i(\langle \omega, x_i\rangle + b) = 1 - \xi_i\right) \text{ or}$$

$$\alpha_i = 0$$

$$\xi_i = 0 \text{ or } \beta_i = 0$$

$$\forall \ i \in [m]$$

$x_i$ are called support vectors
for any $i$ when $\alpha_i \neq 0$

$$\mathcal{L} = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

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

## Dual problem

$$\max_{\alpha} \quad \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

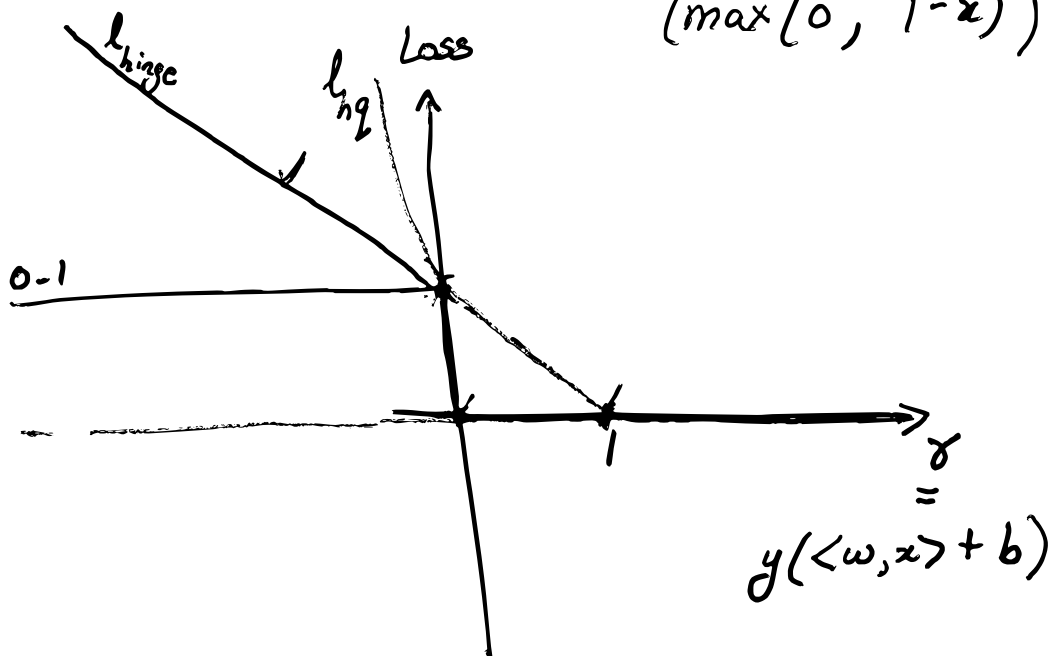$$C \geq \alpha_i \geq 0 \qquad \forall \; i \in [m]$$

## QP

# Margin loss function

$$\min\left(1,\ \max\left(0,\ 1 - \frac{y(\langle w, x\rangle + b)}{\rho}\right)\right) \checkmark$$
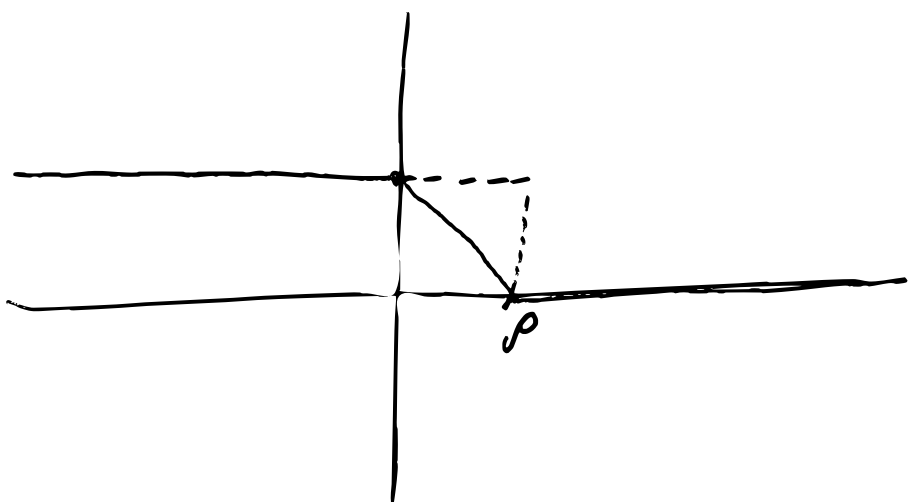
Hinge loss:

$$\ell_{hinge}(x) = \max(0,\ 1-x)$$

$$(\max(0,\ 1-x))^2$$



$l_{hinge}$    $l_{hq}$   Loss

0-1

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

0-1 loss minimization is
NP-hard



$\rho$

$$\to \hat{R}_{S,\rho}(h) = \sum_{i=1}^{m} \min\left(1,\ \max\left(0,\ \frac{(\langle w, x_i\rangle + b)y_i}{\rho}\right)\right)$$

$(w, b)$

$$\leq \sum_{i=1}^{m} 1_{\{(x_i, y_i)\in S:\ (\langle w, x_i\rangle + b)y_i \geq \rho\}}$$

$$\underbrace{\qquad\qquad}_{\hat{R}_S(h)}$$

$$\underset{S\sim D^m}{\mathbb{E}}\ \hat{R}_S(h)$$

$$\| \quad R_S(h) \leq \hat{R}_{S,\rho}(h) + \underline{\qquad}$$

$$\leq \hat{R}_S(h) + \underline{\qquad}$$

If there is an "appropriate" margin
$\rho$ for the data distribution $D$
then, SVM problem generalizes
with a bound that does not
"explicitly" depend on $d$.

$x,\ w \in \mathbb{R}^d$

# Rademacher complexity

$$\text{Rad}(\mathcal{H}) = \frac{1}{m} \underset{\sigma}{E}\left[\sup_{h \in \mathcal{H}} \sum_{i=1}^{m} \sigma_i h(x_i)\right]$$

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

$$\sigma_i = \begin{cases} 1 & \text{Probability } \frac{1}{2} \\ -1 & \text{Probability } \frac{1}{2} \end{cases}$$

$$\sum_{i=1}^{m} \sigma_i h(x_i) : \quad \text{"correlation" between function output \& noise}$$

$$\mathcal{H} = \left\{ h(\cdot, w, b) : \begin{array}{c} \langle w, x \rangle + b \\ = h(x) \end{array} \right\}$$

$$\mathcal{H} = \left\{ h(\cdot, w, b) : h(x) = \langle w, x \rangle + b, \quad \|h(x)\| < \Lambda \right\}$$

Class of Classifiers is

$$\tilde{\mathcal{H}} : \phi \circ \mathcal{H}$$

e.g.

$\phi : \text{sgn}$

$\phi : \text{hinge loss}$

$$\underline{\tilde{\mathcal{H}}} = \left\{ h(\cdot, w, b) : \underset{h(x) =}{\min\{1, \max\{\frac{\langle w, x \rangle + b}{\rho}, 0\}\}} \right\}$$

$\underline{Ex:}$
$\frac{1}{\rho}$

$\rho$

Thm: If $\Phi$ is $\ell$-lipschitz,

$$\text{Rad}_S(\overline{\Phi \circ \mathcal{H}}) \leq \ell \, \text{Rad}_S(\mathcal{H})$$

$\Phi$ is $\ell$-lipschitz if

$$|\Phi(x) - \Phi(y)| \leq \ell \, \|x - y\|$$

$$\forall x, y$$

$$\ell := \sup_x \|\nabla \Phi(x)\|$$

Next time: 1) $\text{Rad}_S(\mathcal{H}) \leq C$

$\downarrow$

linear

2) Use thm