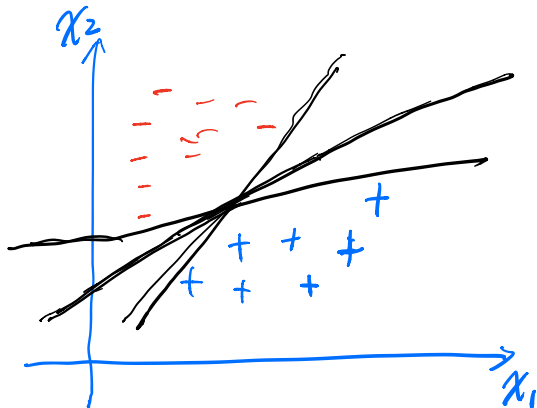


Logistic Loss ✓

Hinge Loss.

Used by Support Vector Machine (SVM)

Main idea: Margin Maximization.



Linearly separable

Infinite linear models  
w/ perfect accuracy.

Want to maximize the distance between  
decision boundary & examples.

$w \in \mathbb{R}^d$  : weight vector

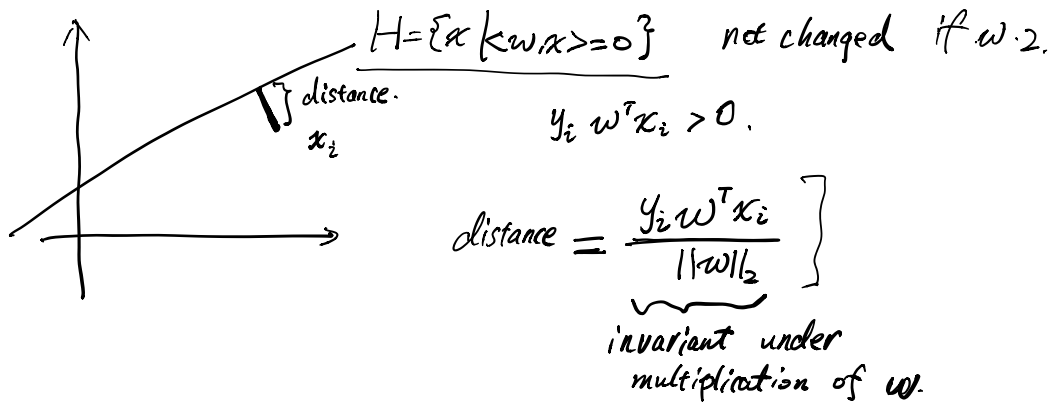
Decision boundary :  $H = \{x \in \mathbb{R}^d \mid \langle w, x \rangle = 0\}$   
↑  
Hyperplane

Suppose  $w$  is perfectly accurate

for all examples  $(x_i, y_i) \in \mathbb{R}^d \times \{\pm 1\}$

$$y_i w^T x_i > 0$$

Distance between  $x$  &  $H$ .



Idea of margin / distance maximization.

$$\max_w \min_i \frac{y_i w^T x_i}{\|w\|_2}$$

closest distance

optimization does not care about  $\|w\|_2$

might as well consider  $w$ 's

such that  $\min_i y_i w^T x_i = 1$

$\Leftrightarrow \max_w \frac{1}{\|w\|_2}$  such that  $\min_i y_i w^T x_i = 1$ . (2)

$\Leftrightarrow \min_w \frac{1}{2} \|w\|_2^2$  such that  $y_i (w^T x_i) \geq 1 \quad \forall i$

(3)

such that  $\min_i y_i (w^T x_i) = 1 \rightarrow$  take  $\frac{w}{2}$ .

> margin maximization formulation  
 "in a perfect world"

What if not linearly separable?

$$\min_w \frac{1}{2} \|w\|_2^2 \quad \text{such that} \quad \boxed{y_i (w^T x_i) \geq 1} \quad \forall_i \text{ all}$$

"Hard" Margin SVM (may not be feasible)

"Soft" Margin SVM.

$$\min_w \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \xi_i$$

such that

constri-  $\forall_i, y_i (w^T x_i) \geq 1 - \xi_i$  "X<sub>i</sub>"

get rid of?  $\forall_i, \xi_i \geq 0$  slack variable

(Equivalent form)

$$\min_w \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max\{0, 1 - y_i w^T x_i\}$$

"Hinge loss"

for some  $\lambda$

$$\min_w \underbrace{\lambda \|w\|_2^2}_{\text{Regularization}} + \sum_{i=1}^n \max\{0, 1 - y_i w^T x_i\}$$

Hinge loss ERM

# Constrained OPT

# Unconstrained OPT

$$\min_w \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \xi_i$$
 such that  

$$\forall i, y_i (w^T x_i) \geq 1 - \underbrace{\xi_i}_{\text{slack variable}}$$

$$\forall i, \xi_i \geq 0.$$

$$\min_w \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max\{0, 1 - y_i w^T x_i\}$$
 "Hinge loss"

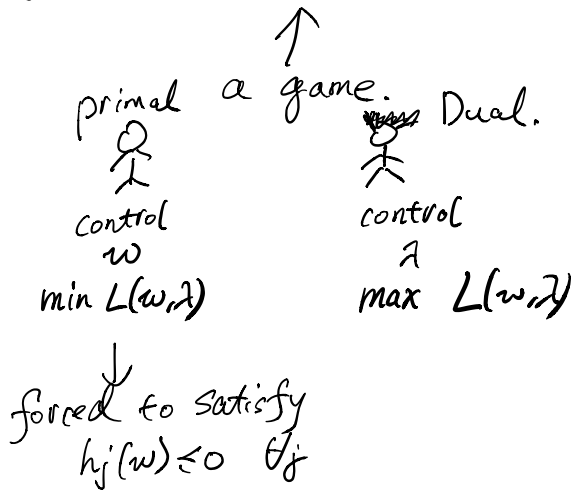
Detour: Lagrange Duality.

$$\min_w F(w) \quad \text{such that} \quad \underline{h_j(w) \leq 0}, \quad \forall j \in [m]$$
  
 Constrained OPT ↓  
{1, 2, ..., m}

"Put constraints into obj"  $w \leftarrow$  "primal —"

For each  $j$ , introduce  $\lambda_j \geq 0 \leftarrow$  "dual variable"

Lagrangian: 
$$L(w, \lambda) = F(w) + \sum_{j=1}^m \lambda_j h_j(w)$$



Suppose  $w$  is not a feasible solution  
 $h_j(w) > 0$   
 $\lambda_j \leftarrow$  "infinity"  
 $L(w, \lambda) \rightarrow \infty$

Zero-sum game view

$\min_w$	$\max_\lambda$	$L(w, \lambda)$	"minmax"	
<div style="border: 1px solid black; border-radius: 50%; padding: 5px; display: inline-block;"> <math>\max_\lambda</math> </div>		$\min_w$	$L(w, \lambda)$	"maxmin"
left		$\longrightarrow$	right	

minmax  $\hat{=}$  minimization goes first  
maximization "best response"

maxmin  $\hat{=}$  opposite

Play second is "usually" better.

Weak Duality:  $\max_w \min_\lambda L(w, \lambda) \leq \min_w \max_\lambda L(w, \lambda)$   
favors min player. favors max player.

Strong Duality  $\hat{=}$   $\max_w \min_\lambda L(w, \lambda) = \min_w \max_\lambda L(w, \lambda)$   
under "mild" condition (e.g., soft SVM, Slater's Condition)

$$\lambda^* = \arg \max_\lambda \left( \min_w L(w, \lambda) \right)$$

$$w^* = \arg \min_w \left( \max_\lambda L(w, \lambda) \right)$$

Strong Duality  $\Rightarrow$

$$\begin{aligned}
 \boxed{F(w^*)} &= \min_w \max_\lambda L(w, \lambda) && \text{(Strong Duality)} \\
 &= \max_\lambda \min_w L(w, \lambda) && \\
 &= \min_w L(w, \lambda^*) && \text{(Defn } \lambda^*) \\
 &\leq L(w^*, \lambda^*) && \text{(Defn min)} \\
 &= \left( F(w^*) + \sum_i \lambda_i^* h_i(w^*) \right) && \text{(Defn of } L)
 \end{aligned}$$

$$\sum_j \lambda_j^* h_j(w^*) \geq 0$$

$\min_w F(w)$  such that  $h_j(w) \leq 0 \quad \forall j \in [m]$   
 $w^*$  is feasible  $\lambda_j^* \geq 0$

$$\sum_j \lambda_j^* h_j(w^*) \leq 0$$

$$\sum_j \lambda_j^* h_j(w^*) = 0$$

Each  $\lambda_j^* h_j(w^*) \leq 0$   
 $\Rightarrow \forall j, \lambda_j^* h_j(w^*) = 0$

### KKT conditions

(Characterization of  $w^*, \lambda^*$ )

- $\forall j, \lambda_j^* h_j(w^*) = 0$  (complementary slackness)
- $w^*$  is the minimizer of  $L(w, \lambda^*)$   
 $\nabla_w L(w^*, \lambda^*) \rightarrow 0$  (Stationarity)
- $\lambda_j \geq 0, h_j(w^*) \leq 0$  for all  $j$  (Feasibility)

Dated back Kuhn-Tucker '1951

Karush '1939 in unpublished master thesis