

$$\text{Dataset} = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \{+1, -1\} \end{array}$$

Goal:  $h: \mathbb{R}^d \rightarrow \{\pm 1\}$

Performance measure

$$\sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

$$\mathbb{1}(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$\min_{h \in \mathcal{H}_{\text{linear}}} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

$$= \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \mathbb{1}(\text{sign}(w^T x_i) \neq y_i)$$

→ NP-HARD

$$\begin{array}{l} x \in \mathbb{R}^d \\ x' = [x \quad 1] \in \mathbb{R}^{d+1} \\ w' = [\underbrace{w_1 \dots w_d}_w \quad \underbrace{w_{d+1}}_b] \\ w'^T x' = w^T x + b \end{array}$$

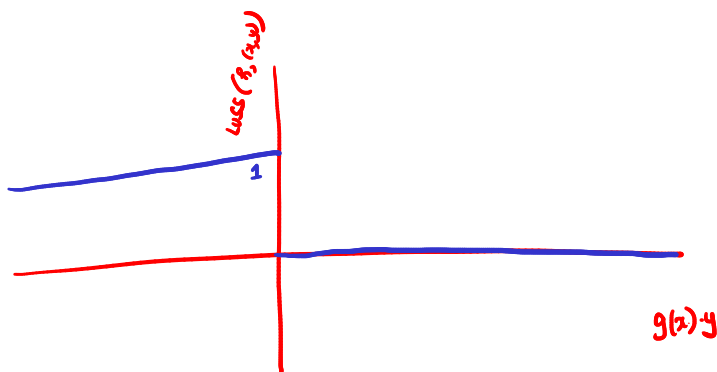
LOSS-FUNCTION VIEW

$$\begin{array}{l} \mathbb{R}^d \rightarrow \{\pm 1\} \\ (x, y) \end{array}$$

$$h: \mathbb{R}^d \rightarrow \{\pm 1\}$$

$$h(x) = \text{sign}(w^T x)$$

$$\mathbb{1}(h(x) \neq y) = \mathbb{1}(\underbrace{\overset{+1}{>0}}_{<0} (w^T x) \cdot \underbrace{\overset{-1}{<0}}_{>0} y < 0)$$



$$h(x) = \text{sign}(g(x))$$

$$\sum_{i=1}^n \mathbb{1}(g(x_i) y_i < 0)$$

ALG 1:

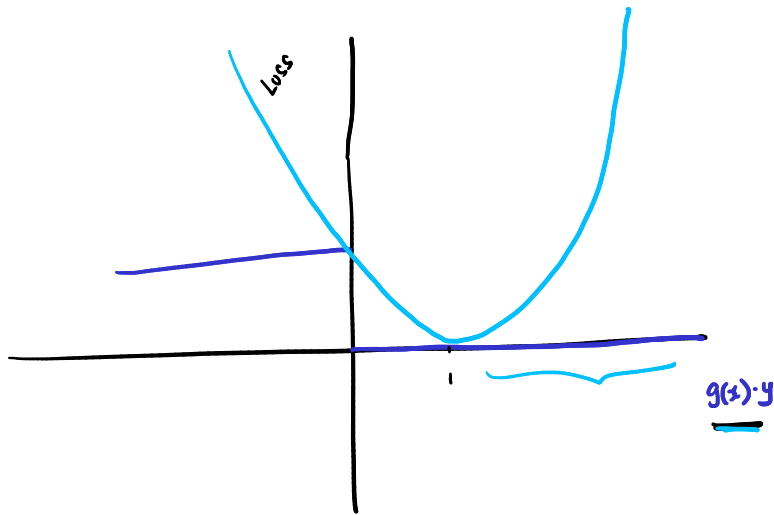
Using Regression for Classification

$$h(x) = \text{sign}(g(x))$$

$$\begin{aligned} \text{Loss}(g, (x, y)) &= \underbrace{(g(x) - \overset{+1}{y})^2}_{\boxed{(g(x)y - 1)^2}} = (g(x))^2 + y^2 - 2g(x)y \\ &= (g(x))^2 + 1 - 2g(x)y \quad \text{--- (1)} \\ &= (g(x) \cdot y)^2 + 1 - 2g(x)y \end{aligned}$$

$$= (g(x))^2 + 1 - 2g(x) \cdot y - 2$$

$$\textcircled{1} = \textcircled{2}$$



- $\rightarrow$  0-1 loss
- $\rightarrow (g(x) \cdot y - 1)^2 \rightarrow$  SQUARED LOSS

### SUPPORT VECTOR MACHINES

$$\min_{w, \xi} \quad \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad \begin{cases} (w^T x_i) y_i + \xi_i \geq 1 \\ \xi_i \geq 0 \end{cases}$$

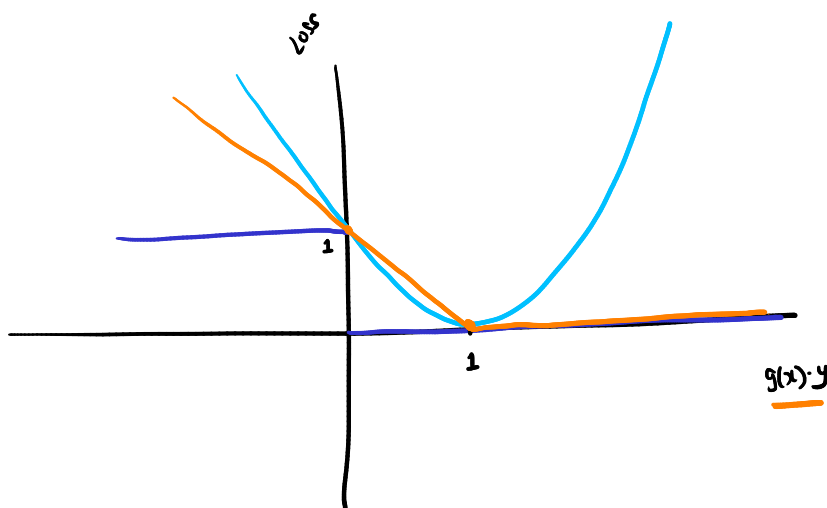
$$\min_{w, \xi} \quad \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad \begin{cases} \xi_i \geq 1 - (w^T x_i) y_i \\ \xi_i \geq 0 \end{cases}$$

$$\xi_i \geq \max(0, 1 - (w^T x_i) y_i)$$

$$\begin{aligned} & \text{Regularizer} \\ & \downarrow \\ & \min_{w \in \mathbb{R}^d} \quad \underbrace{\frac{1}{2} \|w\|^2}_{\text{Model dependent}} + \underbrace{c \sum_{i=1}^n \max(0, 1 - (w^T x_i) y_i)}_{\substack{\text{Data dependent} \\ \text{Loss}}} \end{aligned}$$

REGULARIZATION + HINGE LOSS



- 0-1 loss
- Squared loss  $(g(x) - 1)^2$

$$\text{HINGE LOSS} = \max(0, 1 - g(x) \cdot y)$$

# LOGISTIC REGRESSION

$$\sigma(a) = \frac{1}{1+e^{-a}}$$

$$z_i \in \{0,1\}$$

$$\max_w \prod_{i=1}^n (\sigma(\omega^T x_i))^{z_i} (1 - \sigma(\omega^T x_i))^{(1-z_i)}$$

$$z_i = 1 \quad \text{if } y_i = +1$$

$$z_i = 0 \quad \text{if } y_i = -1$$

$$\max_w \sum_{i=1}^n z_i \log(\sigma(\omega^T x_i)) + (1-z_i) \log(1 - \sigma(\omega^T x_i))$$

$$= \min_w \sum_{i=1}^n \left[ -z_i \log(\sigma(\omega^T x_i)) + (z_i - 1) \log(1 - \sigma(\omega^T x_i)) \right]$$

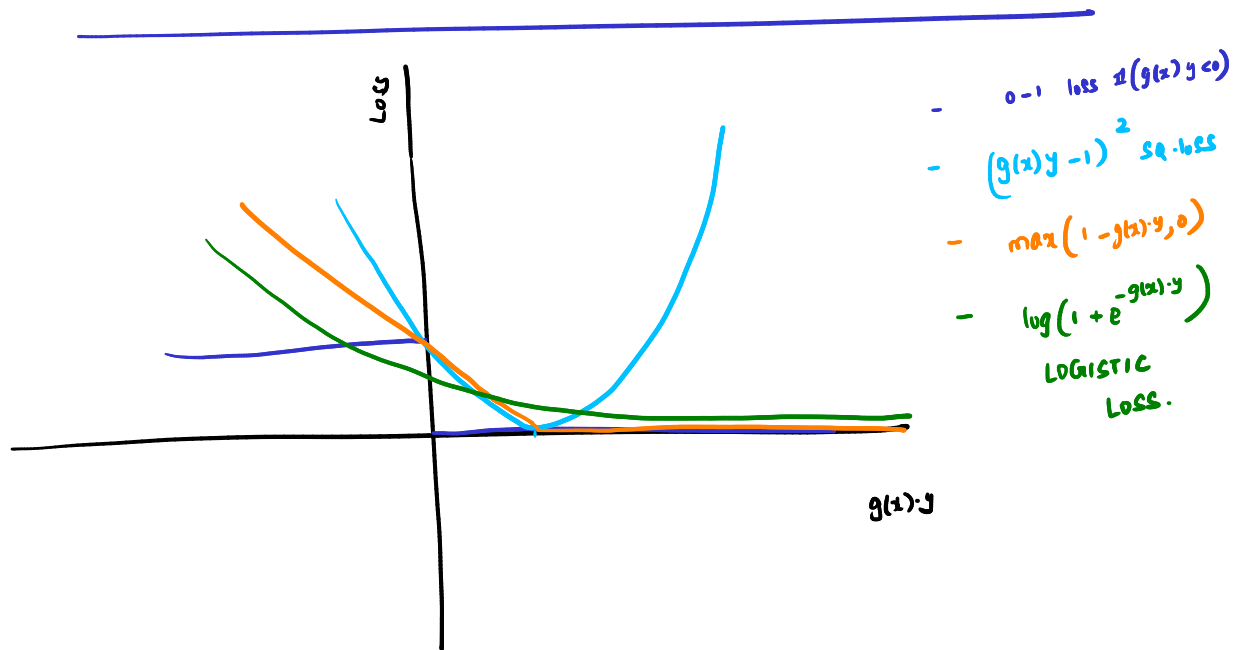
Loss for a single point  $z_i = 1$  ( $y_i = +1$ )

$$\begin{aligned} &= -\log(\sigma(\omega^T x_i)) = -\log\left(\frac{1}{1+e^{-\omega^T x_i}}\right) \\ &= \log(1+e^{-\omega^T x_i}) = \boxed{\log(1+e^{-(\omega^T x_i)y_i})} \end{aligned}$$

Loss for a single point  $z_i = 0$  ( $y_i = -1$ )

$$\begin{aligned} &= -\log(1 - \sigma(\omega^T x_i)) \\ &= -\log\left(1 - \frac{1}{1+e^{-\omega^T x_i}}\right) \\ &= -\log\left(\frac{e^{-\omega^T x_i}}{(1+e^{-\omega^T x_i})/e^{-\omega^T x_i}}\right) \\ &= -\log\left(\frac{1}{e^{\omega^T x_i} + 1}\right) \\ &= \log(1+e^{\omega^T x_i}) \\ &= \boxed{\log(1+e^{-(\omega^T x_i)y_i})} \end{aligned}$$

$$= \min_w \sum_{i=1}^n \log \left( 1 + e^{\underbrace{-\tilde{w}^T x_i y_i}_{g(x_i) y_i}} \right)$$



### CONCLUSIONS

- 0-1 loss is NP-hard to minimize
- Different algorithms use different "surrogate" loss
- Surrogates are convex and hence easy to minimize.

### PERCEPTRON

$$w_{t+1} = w_t + \underline{x_t y_t}$$

### HINGE LOSS

$$\ell(w, (x, y)) = \max(0, \underline{-\tilde{w}^T x y})$$

max-hinge

$$\nabla_w \ell_{\text{hinge}} = \begin{cases} -xy & \text{if } (w^T x)y < 0 \\ 0 & \text{if } (w^T x)y > 0 \\ [-1, 0]xy & \text{if } (w^T x)y = 0 \end{cases}$$

$\swarrow$  mistake  
 $(w^T x)y < 0$   
 $(w^T x)y > 0$   
 $(w^T x)y = 0$

$\hookrightarrow$  chooses  $-xy$  when mistake,  $0$  otherwise.

$$w_{t+1} = w_t - \underbrace{\left(\overset{\frac{1}{2}}{\eta_t}\right)}_{\text{step size}} \underbrace{\nabla_w \ell_{\text{hinge}}(w_t)}_{\text{gradient}}$$

$$= w_t - (-x_t y_t)$$

- Perceptron can be interpreted as S.G.D with modified hinge loss with step size = 1

---

BOOSTING

$$\bullet \text{ Loss}(h, (x, y)) = \frac{e^{-y h(x)}}{e} \hookrightarrow \text{Exponential loss.}$$