

# Logistic Regression

## 1. Linear Classification

Data Vector  $x \in \mathbb{R}^{d+1}$

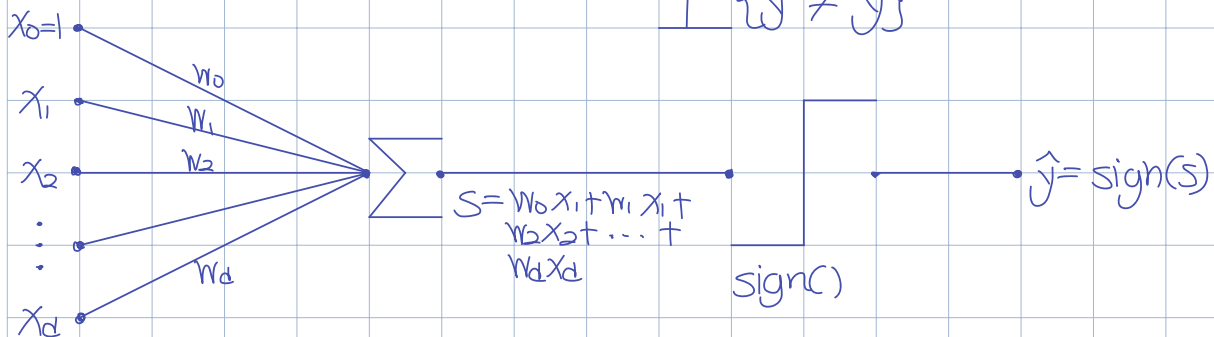
Label  $y \in \{-1, +1\}$

Prediction  $\hat{y} = \text{sign}(w^T x)$

$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

Loss Function

$$1 \{ \hat{y} \neq y \}$$



## 2. Linear Regression

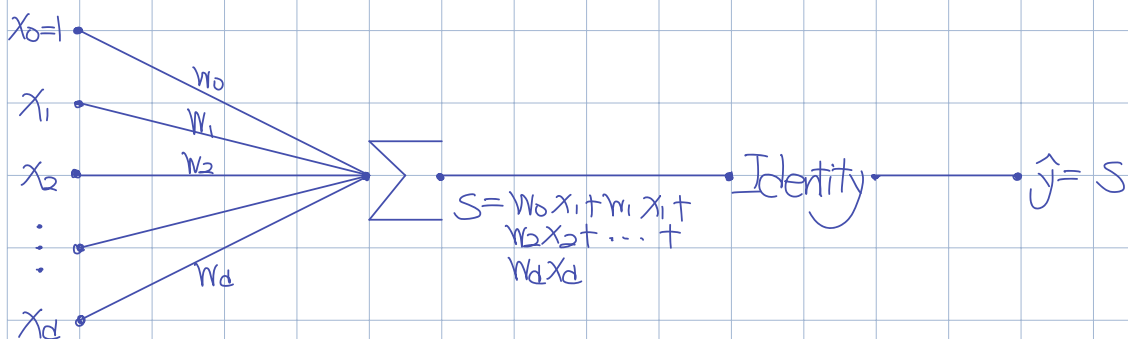
Data Vector  $x \in \mathbb{R}^{d+1}$

Label  $y \in \mathbb{R}$

$$\hat{y} = w^T x$$

Loss Function

$$(y - \hat{y})^2$$



### 3. Logistic Regression

Data Vector  $x \in \mathbb{R}^{d+1}$

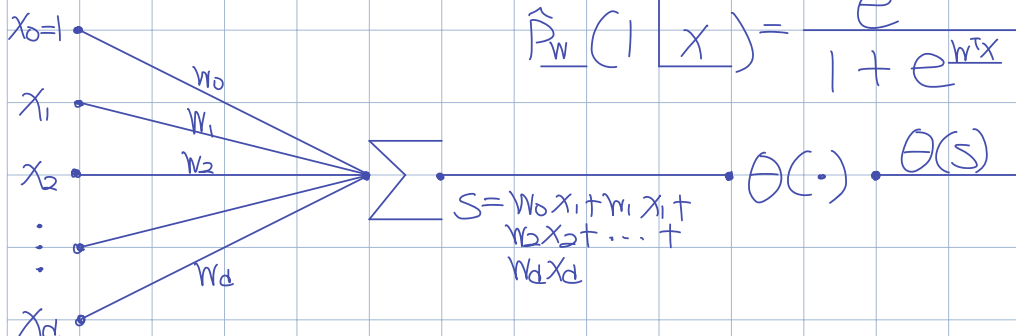
Label  $y \in \{+1, -1\}$

$$Pr(y=1 | x) = \theta(s) = \theta(w^T x)$$

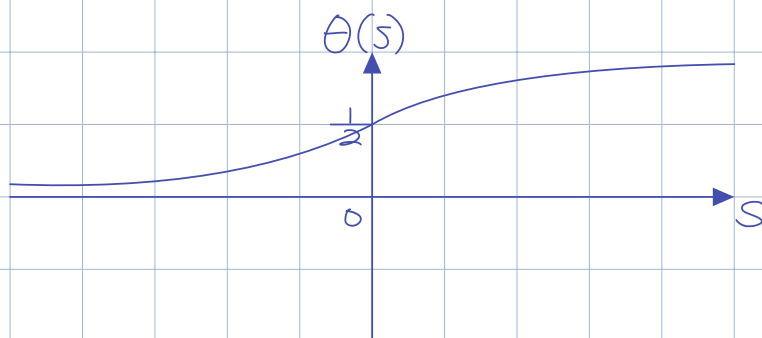
Same

$$= \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$\hat{P}_w(1 | x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$$



$$\theta(s) = \frac{e^s}{1 + e^s} \quad (\text{Non Linear})$$



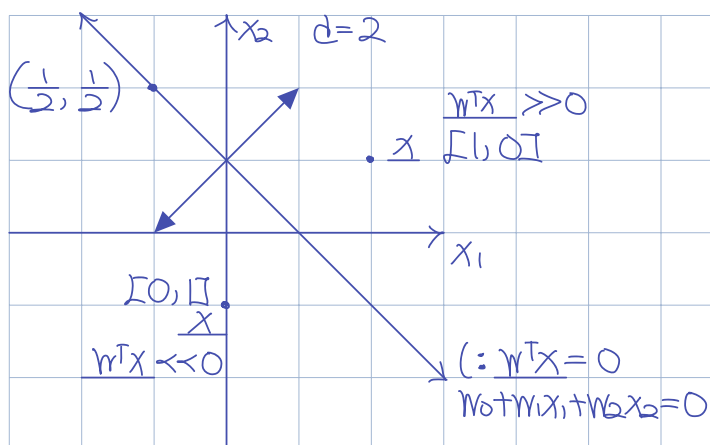
Classification Problem

Sigmoid Function

$$\hat{P}_w(-1 | x) = 1 - \hat{P}_w(1 | x) = 1 - \frac{e^{w^T x}}{1 + e^{w^T x}}$$

$$= \frac{1}{1 + e^{w^T x}} = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

$$\hat{P}_w(y | x) = \frac{e^{y w^T x}}{1 + e^{y w^T x}}, \quad y \in \{-1, +1\}$$



Training Set

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$x_i \in \mathbb{R}^{d_H}, y_i \in \{-1, +1\}$$

Model Parameter:  $w \in \mathbb{R}^{d_H}$

Given  $(x, y)$

Output:  $[\hat{p}_w(1|x), \hat{p}_w(-1|x)]$

$$\hat{p}_w(y|x) = \frac{e^{y w^T x}}{1 + e^{y w^T x}}, y \in \{-1, +1\}$$

Loss Function

Log Loss Function

$$-\log \hat{p}_w(y|x)$$

Example, given  $x$  and  $w$

output probability vector of

logistic regression is  $[0.8, 0.2]$

$$= [\hat{p}_w(1|x), \hat{p}_w(-1|x)]$$

If  $y=1$ , loss  $= -\log 0.8 \approx 0.22$

If  $y=-1$ , loss  $= -\log 0.2 \approx 1.61$

Output = [0.999, 0.001]

If  $y = +1$ ,  $\text{loss} = -\log 0.999 = 10^{-4}$

If  $y = -1$ ,  $\text{loss} = -\log 0.001 = 10$

$$E_n(w) = -\log \hat{P}_w(y_n | x_n)$$

loss on datapoint  
( $x_n, y_n$ )

$$= -\log \frac{e^{y_n w^T x_n}}{1 + e^{y_n w^T x_n}}$$

$$= -\log \frac{1}{1 + e^{-y_n w^T x_n}}$$

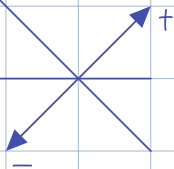
$$= \log(1 + e^{-y_n w^T x_n})$$

Logarithm base e

$$E_n(w) = \log(1 + e^{-y_n w^T x_n})$$

$d=2$   
 $w^T x_n \gg 0$

$w^T x_n \ll 0$



$y_n$	$y_n w^T x_n$	$E_n(w)$
+1	$\gg 0$	$\approx 0$
+1	$\ll 0$	Large
-1	$\gg 0$	$\approx 0$
-1	$\ll 0$	Large

$$E(w) = \frac{1}{N} \sum_{n=1}^N E_n(w) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n w^T x_n})$$

## Training Phase

$$\underline{w} = \underset{\underline{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} E_n(\underline{w})$$

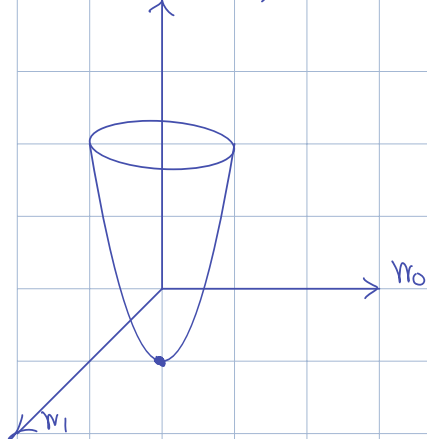
Recall

Linear Regression

$$E_n(\underline{w}) = \frac{1}{N} \sum_{n=1}^N (y_n - \underline{w}^T \underline{x}_n)^2$$

$E_n(\underline{w})$  is convex in  $\underline{w}$

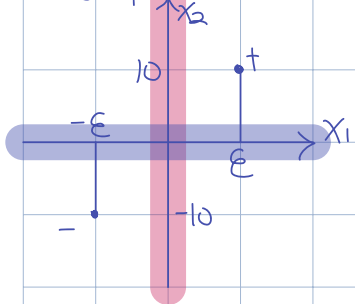
$d=1$



$$\nabla_{\underline{w}} E_n(\underline{w}) = 0$$

Equation In  $\underline{w}$

Example,  $N=2$



$$\underline{x}_1 = (1, \epsilon, 10), y_1 = +1$$

$$\underline{x}_2 = (1, -\epsilon, -10), y_2 = -1$$

$$\epsilon = 10^{-4}$$

Classifier:  $x_1 = 0$

$$w = (0, 1, 0)$$

Classifier:  $x_2 = 0$  (Preferred)

$$w = (0, 0, 1)$$

Compute  $E_n(w)$  for classifier 1 and 2

Classifier 1

$$\begin{aligned} E_n(w) &= \frac{1}{2} (\log(1 + e^{-y_1 w^T x_1}) + \log(1 + e^{-y_2 w^T x_2})) \\ &= \frac{1}{2} (\log(1 + e^{-\epsilon}) + \log(1 + e^{-\epsilon})) \\ &= \log(1 + e^{-\epsilon}) \\ &\approx 0.693 \end{aligned}$$

Classifier 2

$$\begin{aligned} E_n(w) &= \frac{1}{2} (\log(1 + e^{-y_1 w^T x_1}) + \log(1 + e^{-y_2 w^T x_2})) \\ &= \frac{1}{2} (\log(1 + e^{-10}) + \log(1 + e^{-10})) \\ &\approx 5 \times 10^{-4} \end{aligned}$$

$$w_1 = (0, 0, 10)$$

$$w_1^T x = 0$$

$$10x_2 = 0, x_2 = 0$$

$$E_{in}(w_1) < 5 \times 10^{-4}$$

$$\underset{\|w\|}{\operatorname{argmin}} (E_{in}(w) + \lambda \|w\|^2)$$

Regularized Loss Function

### Maximum Likelihood Viewpoint

Training set

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$\Pr(\text{Label Sequence} \mid \text{Data Vector Sequence})$$

$$= \Pr(y_1, y_2, \dots, y_N \mid x_1, x_2, \dots, x_N)$$

Select a model that maximizes probability assigned to label sequence given data vector sequence

Assume that all data samples are generated independently of one another

$$\Pr(y_1, y_2, \dots, y_N \mid x_1, x_2, \dots, x_N)$$

$$= \prod_{n=1}^N \hat{P}_{w_n}(y_n \mid x_n)$$

Probability assigned to data vector  $x_n$

## Max-Likelihood Objective

Select  $w$  that maximizes

$$\prod_{n=1}^N \hat{P}_w(y_n | x_n)$$

$$w^* = \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmax}} \prod_{n=1}^N \hat{P}_w(y_n | x_n)$$

$$= \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmax}} \log \left( \prod_{n=1}^N \hat{P}_w(y_n | x_n) \right)$$

$$= \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmax}} \sum_{n=1}^N \log \hat{P}_w(y_n | x_n)$$

$$= \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N -\log \hat{P}_w(y_n | x_n)$$
$$\log(1 + e^{-y_n w^T x_n})$$

$$= \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} J_n(w)$$



## Cross Entropy Viewpoint

$S = \{s_1, s_2, s_3, \dots, s_M\}$  be a discrete alphabet

Let  $P = (p(s_1), p(s_2), \dots, p(s_M))$

$Q = (q(s_1), q(s_2), \dots, q(s_M))$

be two probability vectors over  $S$

$$CE(P, Q) = - \sum_{i=1}^M P(s_i) \log q(s_i)$$

Log Loss Function Can be viewed as a cross entropy

$$E_n(\underline{w}) = - \log \hat{P}_w(y_n | x_n)$$

$$E_n(\underline{w}) = - \left\{ 1_{\{y_n = +1\}} \log \hat{P}_w(+1 | x_n) + \right. \\ \left. 1_{\{y_n = -1\}} \log \hat{P}_w(-1 | x_n) \right\}$$

$$P_n = (1_{\{y_n = +1\}}, 1_{\{y_n = -1\}})$$

$$Q_n = (\hat{P}_w(+1 | x_n), \hat{P}_w(-1 | x_n))$$

$$\left. \begin{array}{l} y_n = +1 \\ P_n = (1, 0) \\ y_n = -1 \\ P_n = (0, 1) \end{array} \right\} \begin{array}{l} \text{Ideal Output for } (x_n, y_n) \\ \text{Atomic } P_n \end{array}$$

$Q \equiv$  Output generated by our model

$$E_n(w) = \mathbb{E}(P_n, Q_n) \quad P_n = (P_n^{(1)}, P_n^{(2)})$$

Knowledge Distillation