

Lec 7 Logistic Regression

1. Recap:

↳ supervised learning set up:

$$\text{Given } \mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$$

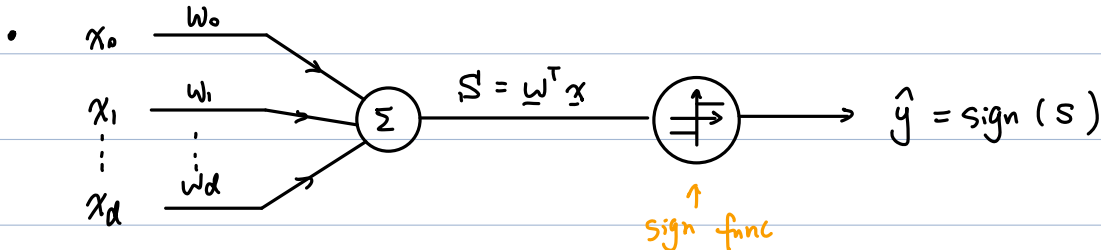
unknown target func. $y = f(\underline{x})$

Hypothesis: $\hat{y} = h(\underline{x})$ where $h \in \mathcal{H}$ // \mathcal{H} is the set of hypothesis

↳ Linear classification:

$$\bullet \underline{x} \in \mathbb{R}^{d+1}, \quad y = \{+1, -1\}$$

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



• The error for each of the sample is:

$$e(\underline{w}) = 1 \{ y \neq \text{sign}(s) \}$$

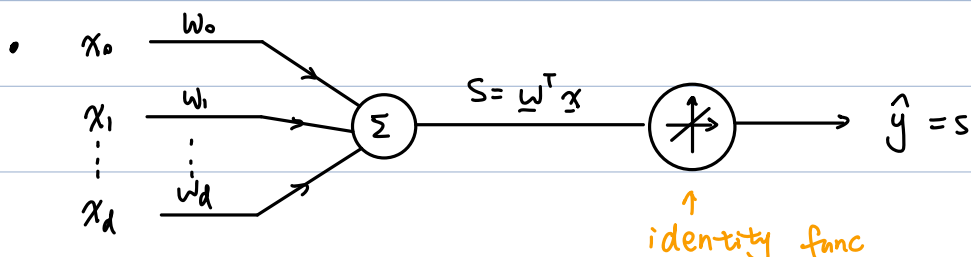
↗ // whether y is equal to estimation or not

↳ Linear Regression:

$$\bullet \underline{x} \in \mathbb{R}^{d+1}, \quad y \in \mathbb{R}$$

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

// compare w/ classification, this don't take the $\text{sign}(\cdot)$



$$\bullet e(\underline{w}) = (y - s)^2$$

All above is deterministic hypothesis

2. What if we want randomness in h ?

↳ example:

- \underline{x} = average (fat, sugar) in diet.

y = heart attack

- This is not deterministic. Nobody say this amount of fat & sugar will definitely have heart attack

- We want to say

large $\underline{w}^T \underline{x} \Rightarrow$ more likely that $y = +1$

- \therefore Want new target function

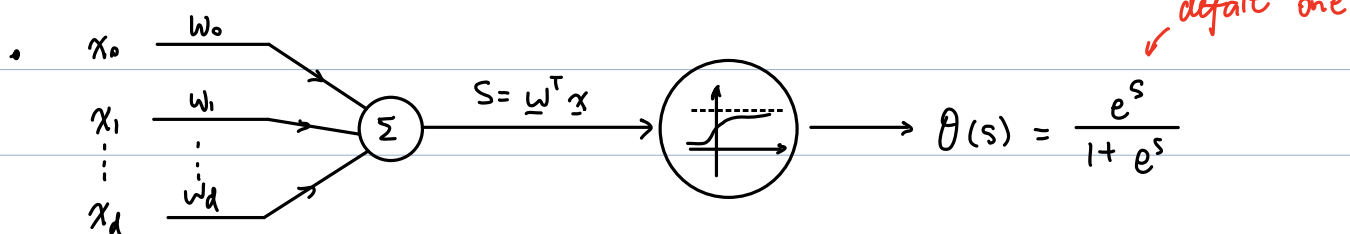
$$f_p(\underline{x}) = \Pr\{y = +1 \mid \underline{x}\}$$

"soft decision"

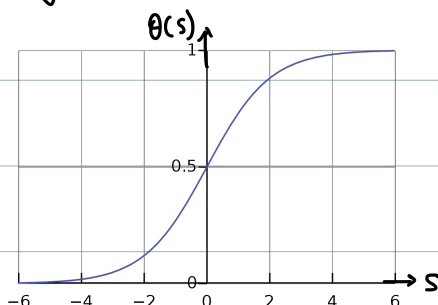
3. Logistic Regression:

↳ setup:

- $\underline{x} \in \mathbb{R}^{d+1}$, $y \in \{+1, -1\}$



- $\theta(s)$: sigmoid function // "sigmoid" means S-shaped



// another $\theta(s)$: tanh

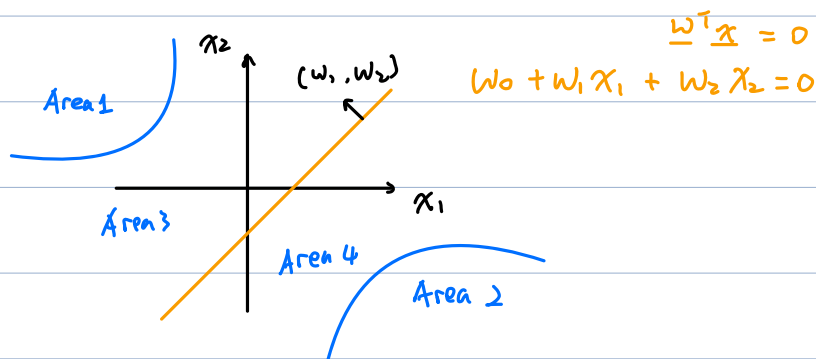
$$\tanh = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

↳ Hypothesis: (guess)

$$Pr\{y = +1 | \underline{x}\} = \theta(s) = \theta(\underline{w}^T \underline{x}) = \frac{e^{\underline{w}^T \underline{x}}}{1 + e^{\underline{w}^T \underline{x}}}$$

↳ Note:

① Geometry: ($d=2$)



• Area 1: large $\|\underline{w}^T \underline{x}\|$, positive

$$\Rightarrow \theta(s) = \frac{e^{\underline{w}^T \underline{x}}}{1 + e^{\underline{w}^T \underline{x}}} = \frac{\infty}{1 + \infty} = 1$$

• Area 2: large $\|\underline{w}^T \underline{x}\|$, negative

$$\Rightarrow \theta(s) = \frac{e^{\underline{w}^T \underline{x}}}{1 + e^{\underline{w}^T \underline{x}}} = \frac{0}{1 + 0} = 0$$

• Area 3: small $\|\underline{w}^T \underline{x}\|$, positive

• Area 4: small $\|\underline{w}^T \underline{x}\|$, negative

$$\textcircled{2} Pr\{y = -1 | \underline{x}\} = 1 - Pr\{y = +1 | \underline{x}\}$$

$$= 1 - \frac{e^{\underline{w}^T \underline{x}}}{1 + e^{\underline{w}^T \underline{x}}} = \frac{1}{1 + e^{\underline{w}^T \underline{x}}} = \frac{e^{-\underline{w}^T \underline{x}}}{1 + e^{-\underline{w}^T \underline{x}}} = \theta(-\underline{w}^T \underline{x}) = \theta(-s)$$

4. loss function of logistic regression:

↳ Notation:

$$\hat{P}_{\underline{w}}(1 | \underline{x}) = \theta(\underline{w}^T \underline{x})$$

the estimate for $\Pr\{y=+1 | \underline{x}\}$

// $\hat{P}_{\underline{w}}$ means is a func. of \underline{w}

$$\hat{P}_{\underline{w}}(-1 | \underline{x}) = \theta(-\underline{w}^T \underline{x})$$

• Combine the two above:

$$\hat{P}_{\underline{w}}(y | \underline{x}) = \theta(y \underline{w}^T \underline{x}), \quad \text{for } y \in \{+1, -1\}$$

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

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

↳ Error criterion:

• For n th sample in training data, we define the error

$$\begin{aligned} e_n(\underline{w}) &= -\log[\hat{P}_{\underline{w}}(y_n | \underline{x}_n)] \\ &= \log(1 + e^{-y \underline{w}^T \underline{x}_n}) \end{aligned}$$

// log loss function

↳ Note:

• If $\underline{w}^T \underline{x} \gg 0$ (i.e. very far away from the line) and $y = +1$

$$\Rightarrow -y \underline{w}^T \underline{x} = -\infty \Rightarrow e^{-y \underline{w}^T \underline{x}} = 0$$

$$\Rightarrow \text{loss} = \log(1 + 0) = 0$$

• If $\underline{w}^T \underline{x} \ll 0$, and $y = -1 \Rightarrow \text{loss} = 0$

true label
↓

J. Example:

Q: suppose output of logistic regression is

$$\hat{P}_{\underline{w}}(-1 | \underline{x}_n) = 0.8 \quad \hat{P}_{\underline{w}}(+1 | \underline{x}_n) = 0.2$$

$$\text{if } y_n = +1, \quad e_n(\underline{w}) = -\log(0.2) \approx 1.61$$

$$\text{if } y_n = -1, \quad e_n(\underline{w}) = -\log(0.8) \approx 0.22$$

↳ Note: Loss is smaller if $y_n = -1$

because our $\hat{P}_{\underline{w}}$ assigns higher prob. for $y_n = -1$

$$\text{suppose } \hat{P}_{\underline{w}}(-1 | \underline{x}_n) = 0.999, \quad \hat{P}_{\underline{w}}(+1 | \underline{x}_n) = 0.001,$$

$$\text{if } y_n = +1, \quad e_n(\underline{w}) = -\log(0.001) = 10$$

$$\text{if } y_n = -1, \quad e_n(\underline{w}) = -\log(0.999) = 10^{-4}$$

Summary: Logic Regression

Given $\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$, find $\underline{w} \in \mathbb{R}^{d+1}$,

to minimize

$$\begin{aligned} E_{\text{in}}(\underline{w}) &= \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) \\ &= \frac{1}{N} \sum_{n=1}^N \log(1 - e^{-y_n \underline{w}^T \underline{x}_n}) \end{aligned}$$