## Recap : Logistic Regression
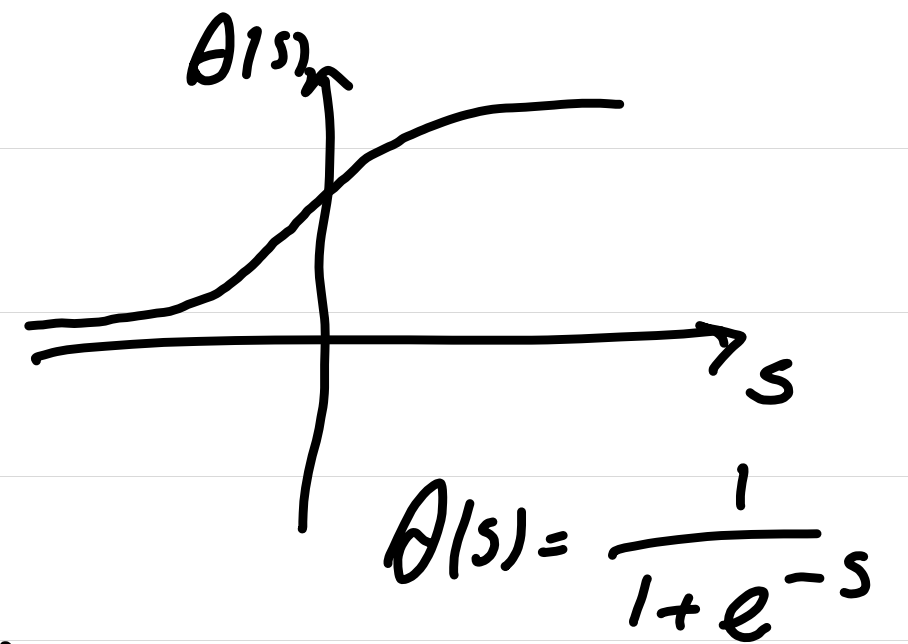
$$\hat{P}_{\underline{w}}(y_n | \underline{x}_n) = \theta(y_n \underline{w}^T \underline{x}_n)$$

$$\theta(s) = \frac{1}{1 + e^{-s}}$$

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = \log\left(1 + e^{-y_n \underline{w}^T \underline{x}_n}\right)$$

Today: We will see why we use this loss function?
(Mathematically speaking)

We want to predict with randomness

$$\to \left(: \hat{P}_{\underline{w}}(y \mid \underline{x}) = \frac{1}{1 + e^{-y \underline{w}^T \underline{x}}}\right.$$

We saw why it makes sense

We need a $E_{in}(\underline{w})$

We use log-loss
$$e_n(w) = -\log \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$$

Today: Why does log-loss make sense?

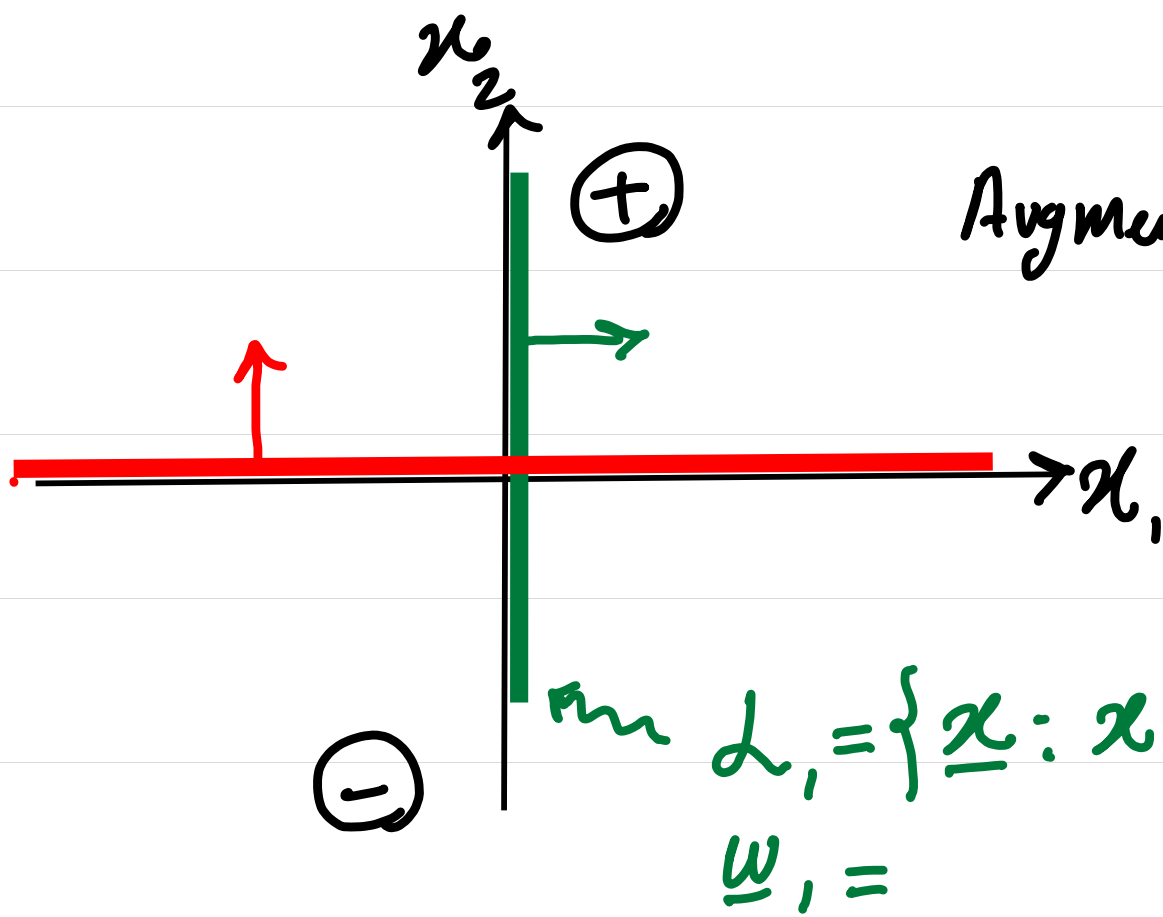It can distinguish better different hyperplanes: Numerical examples

# Benefits Over Linear Classification

E.g.: $d=2$

$N=2$

$d_2 = \{\underline{x} : x_2 = 0\}$

$\underline{w}_2 =$

Augmented $\nearrow$ $\underline{x}_1 = (1, 0.001, 10), y_1 = +1$

$\underline{x}_2 = (1, -0.001, -10), y_2 = -1$

$d_1 = \{\underline{x} : x_1 = 0\}$

$\underline{w}_1 =$

■ For linear classification, which line is better?

$\left( e_n(\underline{w}) = \mathbb{1}(y_n \neq \text{sign}(\underline{w}^T \underline{x}_n)) \right)$ : They are the same

$$E_{in}(\underline{w}_1) = E_{in}(\underline{w}_2) = 0$$

■ Intuitively, which line is better?

For logistic regression, which line is better?

$$E_{in}(\underline{w}_1) = \frac{1}{2}\left[\log\left(1+e^{-y_1 \underline{w}_1^T x_1}\right) + \log\left(1+e^{-y_2 \underline{w}_1^T x_2}\right)\right.$$

$$= \frac{1}{2}\left[\log\left(\phantom{xxxx}\right) + \log\left(\phantom{xxxx}\right)\right] \approx$$

$$E_{in}(\underline{w}_2) = \frac{1}{2}\left[\log\left(1+e^{-y_1 \underline{w}_2^T x_1}\right) + \log\left(1+e^{-y_2 \underline{w}_2^T x_2}\right)\right)$$

$$= \frac{1}{2}\left[\log\left(\phantom{xxxx}\right) + \log\left(\phantom{xxxx}\right)\right] \approx$$

So,

☐ $\mathcal{L}_1$ is preferred          ☐ $\mathcal{L}_2$ is preferred

■ **What about** $\underline{W}_3 = (0, 0, 100)$ ?

⌦ This is the same line as $d_2$.

⌦ What about $E_{in}(\underline{W}_3)$ ? $\frac{1}{2}(\log(1 + e^{-100}) + \log(1 + e^{-100}))$

  ● $E_{in}(\underline{W}_3)$ is much lower!

  ● But $\underline{W}_3$ and $\underline{W}_2$ are the same lines.

■ We can fix the norm of $\|W\|_2^2$ to be 1.

⌦ But this constraint makes the optimization more difficult <span style="color:red">challenging</span> ⟿ <span style="color:red">$\min\limits_{\underline{W}} E_{in}(\underline{W})$</span>

<span style="color:red">s.t. $\|\underline{W}\| = 1$</span>

■ Typically, we regularize logistic regression

$$\min_{\underline{w}} E_{in}(\underline{W}) + \lambda \|W\|_2^2$$

- So far, we have seen numerical demostration why log-loss makes sense.
- Time for more rigorous math.
  - There are two mathematical explanation for log-loss
    1. When we minimize log-loss, we maximize the likelihood
    2. when we minimize log-loss, we minimize cross-entropy

We want to predict with randomness

$\rightarrow$

$\mathcal{H}: \hat{P}_{\underline{w}}(y|\underline{x}) = \dfrac{1}{1 + e^{-y \underline{w}^T \underline{x}}}$

We saw why it makes sense

$\rightarrow$

We need a $E_{in}(\underline{w})$

We use log-loss

$e_n(w) = -\log \hat{P}_{\underline{w}}(y_n|\underline{x}_n)$

Today: Why does log-loss make sense?

It can distinguish better different hyperplanes: Numerical examples

Maximum likelihood interpretation:
when we minimize log-loss, we maximize likelihood

■ Let $D = \{(\underline{x}_1, y_1), \ldots, (\underline{x}_N, y_N)\}$ be the observed datapoint.

■ Consider $P(y_1, y_2, \ldots, y_N \mid \underline{x}_1, \underline{x}_2, \ldots, \underline{x}_N)$

$$= \mathbb{P}\left[1^{st} \text{ label is } y_1, \ 2^{nd} \text{ label is } y_2, \ldots \mid 1^{st} \text{ example} = \underline{x}_1, \ 2^{nd} \text{ example} = \underline{x}_2, \ldots\right]$$

$= \text{"likelihood of observing this particular labels } y_1, \ldots, y_N$
$\text{given datapoint } \underline{x}_1, \ldots, \underline{x}_N \text{"}$

■ This is the joint distribution.

■ Assuming having I.I.D. examples, <span style="color:blue">independent and identically distributed</span>

$$P(y_1, y_2, \ldots, y_N \mid \underline{x}_1, \underline{x}_2, \ldots, \underline{x}_N) = \prod_{n=1}^{N} P(y_n \mid \underline{x}_N)$$

<span style="color:red">← unknown target probability distribution</span>

$\left( \begin{array}{l} \text{We want a} \\ \hat{P}_{\underline{w}} \in \mathcal{H} \text{ s.t.} \end{array} \right) \quad \approx \prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$

<span style="color:blue">likelihood</span>

<span style="color:red">We want to find $\underline{w}$ that maximizes $\prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$</span>

$$\Longleftrightarrow \underset{\underline{w}}{\text{maximize}} \ \frac{1}{N} \log \prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$$

<span style="color:blue">$e_n(\underline{w})$</span>

$$\Longleftrightarrow \underset{\underline{w}}{\max} \ \frac{1}{N} \sum_{n=1}^{N} \log \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n) \Longleftrightarrow \underset{\underline{w}}{\min} \ \frac{1}{N} \sum_{n=1}^{N} -\log(y_n \mid \underline{x}_n)$$

$$\Longleftrightarrow \underset{\underline{w}}{\min} \ E_{in}(\underline{w})$$

We want to predict with randomness

$\rightarrow$

rC: $\hat{P}_{\underline{w}}(y \mid \underline{x}) = \dfrac{1}{1 + e^{-y \underline{w}^T \underline{x}}}$

We saw why it makes sense

$\rightarrow$

We need a $E_{in}(\underline{w})$

We use log-loss
$e_n(w) = -\log \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$

Today: Why does log-loss make sense?

It can distinguish better different hyperplanes: Numerical examples

Maximum likelihood interpretation:
when we minimize log-loss, we maximize likelihood

Cross-entropy interpretation:
When we minimize log-loss, We minimize Cross-entropy

# Cross-Entropy Interpretation

Defn: Suppose $P$ and $Q$ are to distribution over the same Sample space $S$

Sample space $S$: The set of all values a r.v. can take

E.g.: poisson R.V. with mean:
$$S = \{0, 1, \dots\}$$
$$\text{distribution}: P(K) = \frac{\lambda^K}{K!} e^{-\lambda}$$

The cross-entropy b/w $P$ and $Q$ is
$$CE(P, Q) = -\sum_{K \in S} P(K) \log Q(K)$$

(measures the difference b/w $P$ and $Q$)

- For the $n^{th}$ example, Consider the "true" distribution

$P_n = \left( \mathbb{P}[y_n = +1], \mathbb{P}[y_n = -1] \right)$ <span style="color:blue">(distribution of the true labels)</span>

$$= \begin{cases} (1, 0), & \text{if } y_n = 1 \\ \\ (0, 1), & \text{if } y_n = -1 \end{cases}$$

$$= \left( \mathbb{1}(y_n = +1), \mathbb{1}(y_n = -1) \right)$$

- For $n^{th}$ example, Consider the estimated distribution
$$Q_n = \left( \hat{P}_{\underline{w}} (1 | \underline{x}_n) , \hat{P}_{\underline{w}} (-1 | \underline{x}_n) \right)$$ <span style="color:blue">"estimated distribution of $y_n$ given example $\underline{x}_n$"</span>

- The "closer" $Q_n$'s are to $P_n$'s, the better it is.
- what does "closeness" mean ? Cross-Entropy

$$CE(P_n, Q_n) = - \left( \mathbb{1}_{(y_n = +1)} \hat{P}_{\underline{w}} (1 | \underline{x}_n) + \mathbb{1}_{(y_n = -1)} \hat{P}_{\underline{w}} (-1 | \underline{x}_n) \right)$$

$$= - \log \hat{P}_{\underline{w}} (y_n | \underline{x}_n) = e_n (\underline{w})$$

<span style="color:blue">Min $E_{in}(\underline{w}) \Longleftrightarrow$ Min average "distance" b/w $P_n$'s and $Q_n$'s.</span>

$$\frac{1}{N} \sum_{n=1}^{N} CE(P_n, Q_n)$$

- Hopefully we are all convinced that log-loss is a reasonable loss function

- But, what is our learning algorithm to find the $\underline{w}$ that minimizes $E_{in}(\underline{w})$

  ↳ Next time. Gradient descent.