# Logistic Regression

Still dealing w/ linear model    $\underline{w}^T \underline{x}$



loss: $\mathbb{1}\{y \neq \hat{y}\}$

sign → $\hat{y} = \text{sign}(s)$

same linear model "frontend"

lin. reg → $\hat{y} = s$    loss: $(y-s)^2$

Logistic Regression

sigmoid → $\hat{y} = \Theta(s)$    this is just probability of estimate

$$\Theta = \frac{e^s}{1+e^s}$$



confidence/probability

Interpretation:

$$\Theta(s) = \mathbb{P}(y=1 \mid \underline{x})$$

or

probability estimate    $\hat{P}_{\underline{w}}(y \mid \underline{x}) = \text{Bernoulli}\left(y \mid \Theta(\underline{w}^T\underline{x})\right)$

$$\hat{P}_{\underline{w}}(y=1 \mid \underline{x}) = \Theta(\underline{w}^T\underline{x}) = \frac{e^{\underline{w}^T\underline{x}}}{1+e^{\underline{w}^T\underline{x}}}$$
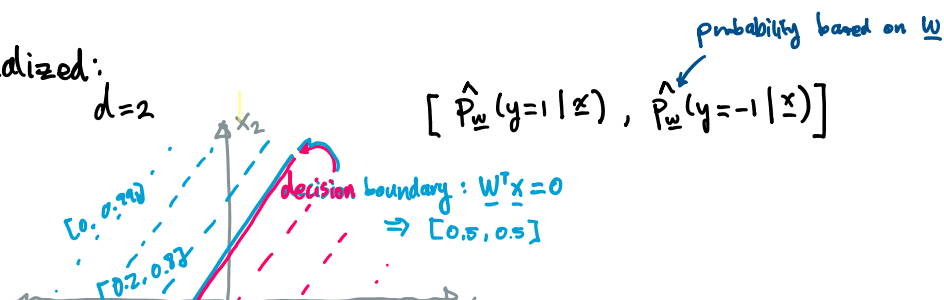
$$\hat{P}_{\underline{w}}(y=-1 \mid \underline{x}) = 1-(\ldots) = \frac{e^{-\underline{w}^T\underline{x}}}{1+e^{-\underline{w}^T\underline{x}}}$$
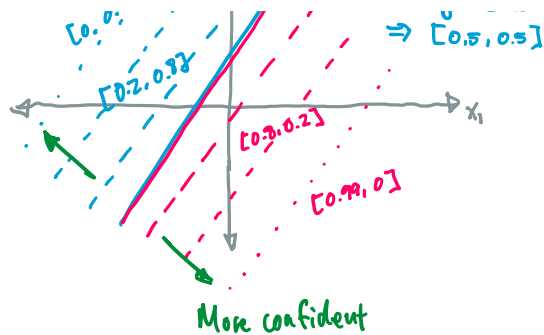
compact ⟹ $\hat{P}_{\underline{w}}(y \mid \underline{x}) = \dfrac{e^{y\underline{w}^T\underline{x}}}{1+e^{y\underline{w}^T\underline{x}}}$

Log Odds (Logit) function!    $\log\left(\dfrac{\Theta(s)}{1-\Theta(s)}\right) = \underline{W}^T\underline{s}$

Visualized:

$d=2$    probability based on $\underline{w}$

$\left[\ \hat{P}_{\underline{w}}(y=1 \mid \underline{x})\ ,\ \hat{P}_{\underline{w}}(y=-1 \mid \underline{x})\ \right]$



decision boundary: $W^Tx = 0$

⟹ $[0.5, 0.5]$

$[0, 0.990]$

$[0.2, 0.88]$

[0.0.] ⇒ [0.5, 0.5]
[0.2, 0.88]
[0.9, 0.2]
[0.98, 0] x₁

**More confident**

**Loss Function:** Measures uncertainty

$$\text{Loss} = -\log \hat{p}_{\underline{w}}(y|\underline{x})$$

this is entropy : shows "uncertainty" or "surprise"

ex.  $\hat{p} = 0.8$

| $\hat{p}$ | $y=1$ | $y=-1$ |
|---|---|---|
| 0.8 | 0.22 | 1.61 |
| $10^{-3}$ | 10 | $10^{-4}$ |

(true value)

⟵ we check the probability of the true 'y' value

(?)

**Training Set** $\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots\}$

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

Loss on $(\underline{x}_n, y_n) = \log\left(1 + e^{-y_n \underline{w}^T \underline{x}_n}\right)$
for each set of data

$$E_{in}(\underline{w}) = \frac{1}{N}\sum_{n=1}^{N} e_n(\underline{w})$$

To find $\underline{w}^*$ :  $\underline{w}^* = \underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{argmin}} \ E_{in}(\underline{w})$

⇒ $E_{in}(\underline{w})$ is convex function of $\underline{w}$

⇒ With Regularization :

$E_{in}(\underline{w}) + \lambda \|\underline{w}\|^2$ is also convex for $\lambda \geq 0$

↳ So this means there's an unique global minimum

⚠ but unlike linear regression, we don't have a closed form solution.

**Maximum Likelihood Viewpoint**

$$\hat{p}_{\underline{w}}(y_n|\underline{x}_n) = \Theta(y_n \underline{w}^T \underline{x}_n) = \frac{1}{1 + e^{y_n \underline{w}^T \underline{x}_n}}$$

Likelihood:

$$\hat{P}_{\underline{w}}(y_1, \ldots, y_N \mid \underline{x}_1 \ldots x_N) = \prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid x_n)$$

$0.5 \times 0.6 \times 0.1 \times \cdots$

we want to find $\underline{w}$ such that we maximizes <u>likelyhood</u>

AKA:

$$\Rightarrow \quad \underline{w}^* = \underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{argmax}} \prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n)$$

<span style="color:blue">log of products $\equiv$ sum of logs (easier to compute)</span>

$$= \underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{argmax}} \ \log\left( \prod_{n=1}^{N} \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n) \right)$$

$$= \underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{argmax}} \ \sum_{n=1}^{N} \log\left( \hat{P}_{\underline{w}}(y_n \mid \underline{x}_n) \right)$$

$$= \underset{\underline{w} \in \mathbb{R}^{d+1}}{\text{argmin}} \left(\frac{1}{N}\right) \sum_{n=1}^{N} \underbrace{-\log\left( \hat{P}_{\underline{w}}(y_n \mid x_n) \right)}_{e_n(\underline{w})}$$

$$\underbrace{\phantom{-\log\left( \hat{P}_{\underline{w}}(y_n \mid x_n) \right)}}_{E_{in}(\underline{w})}$$

# Cross Entropy Viewpoint

discrete alphabets
$$\mathcal{S} = \{s_1, s_2 \ldots s_m\}$$

probability dist.
$$P = \{p_1, p_2 \ldots p_m\}$$

<span style="color:blue">this looks like entropy</span>

Cross Entropy $\ CE(P, Q) = -\sum_{i=1}^{M} p_i \cdot \log(q_i)$

<span style="color:blue">$H(P) = -\sum_{i=1}^{m} p_i (\log(p_i))$</span>

$$= H(P) + \underline{D_{KL}(P \| Q)}$$

distance b/w $P$ & $Q$

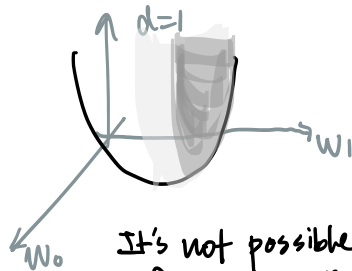▷ The log-loss function can be viewed as cross entropy

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

$$P(\{y_n = +1\}, \{y_n = -1\}) = -1\{y_n = +1\} \cdot \log \hat{P}_{\underline{w}}(+1 \mid x_n) - 1\{y_n = -1\} \log \hat{P}_{\underline{w}}(-1 \mid x_n)$$

$$P(\{y_n=+1\}, \{y_n=-1\}) = -1\{y_n=+1\} \cdot \log \hat{p}_{\underline{w}}(+1|x_n) - 1\{y_n=-1\} \log \hat{p}_{\underline{w}}(-1|x_n)$$

$$\underset{p_i}{} \quad \underset{\log(q_i)}{} \quad \underset{p_i}{} \quad \underset{\log(q_i)}{}$$

How do we minimize $E_{in}(\underline{w})$?

Similar to linear reg. $E_{in}(\underline{w})$ is concave:



$d=1$

$w_1$

$w_0$

It's not possible to solve for $\underline{w}^*$
for $E_{in}(w^*) = \emptyset$

$\Rightarrow$ iterative solution using GRADIENT DESCENT