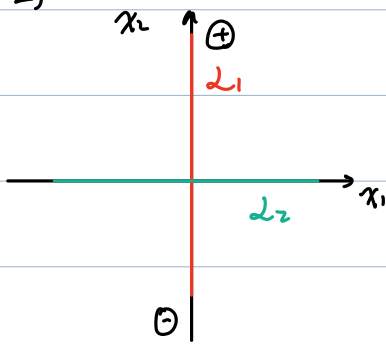↳ **Main goal:**

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

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

---

## 1. Benefit over linear classification:

↳ e.g. $(d=2)$



$N = 2$

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

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

$\mathcal{L}_1 : x_1 = 0 \qquad \underline{w}_1 = (0, 1, 0)$

$\mathcal{L}_2 : x_2 = 0 \qquad \underline{w}_2 = (0, 0, 1)$

$\qquad\qquad\qquad\qquad\qquad\quad \uparrow \quad \uparrow \quad \nwarrow$
$\qquad\qquad\qquad\qquad\qquad w_{x_0} \;\; w_{x_1} \;\; w_{x_2}$

- **Linear classification:**

$$e_n(\underline{w}) = \mathbb{1}(y_n \neq \text{sign}(\underline{w}^T \underline{x}_n))$$

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

It does not tell us which line is better.

- **Logistic Regression:**

$$= \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0.001 & 10 \end{bmatrix}$$

$$E_{in}(\underline{w}_1) = \frac{1}{2}\left[ \log(1 + e^{-(1)\cdot \underline{w}_1^T \underline{x}_1}) + \log(1 + e^{-(-1)\cdot \underline{w}_1^T \underline{x}_2}) \right]$$

$$= \frac{1}{2}\left[ \log(1 + e^{-0.001}) + \log(1 + e^{-0.001}) \right]$$

$$\approx 0.693$$

$$E_{in}(\underline{w}_2) = \frac{1}{2}\left[ \log(1 + e^{-\underline{w}_2^T \underline{x}_1}) + \log(1 + e^{+\underline{w}_2^T \underline{x}_2}) \right]$$

$$= \frac{1}{2}\left[ \log(1 + e^{-10}) + \log(1 + e^{-10}) \right]$$

$$\approx 5 \times 10^{-5}$$

$\Rightarrow \mathcal{L}_2$ is clearly preferred.

- **Note:** what if $\underline{w}_2 = (0, 0, 100)$?

  i.e. $0(x_0) + 0(x_1) + 100(x_2) = 0.$

  It's the same line as $\mathcal{L}_2$

  But the loss $E_{in}(\underline{w})$ is much lower

  $\Rightarrow$ Logistic regression typically requires some regularization.

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

## 2. Maximum likelihood viewpoint

↳ **set up:**

Let $\{(x_1, y_1), (x_2, y_2) \cdots (x_N, y_N)\}$ be the trainig dataset

consider $P(y_1, y_2 \cdots y_N | \underline{x}_1, \underline{x}_2, \cdots \underline{x}_N) \triangleq Pr[1^{st} \text{ label is } y_1 \cdots | 1^{st} \text{ example is } x_1, \cdots]$

$$= \prod_{n=1}^{N} P(y_n | \underline{x}_n)$$

by Hypothesis $\longrightarrow$

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

↳ **Maximum Likelihood approach:**

want to find $\underline{w} \in \mathbb{R}^{d+1}$ that maximize $P(y_1, y_2 \cdots y_N | x_1, x_2 \cdots x_N)$

$\Leftrightarrow$ Maximize $\frac{1}{N} \log \prod_{n=1}^{N} \hat{P}_{\underline{w}} (y_n | \underline{x}_n)$

$$= \frac{1}{N} \sum_{n=1}^{N} \log \hat{P}_{\underline{w}} (y_n | \underline{x}_n)$$

$\Leftrightarrow$ Minimize $\frac{1}{N} \sum_{n=1}^{N} -\log \hat{P}_{\underline{w}} (y_n | \underline{x}_n)$

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

$$= E_{in} (\underline{w})$$

## 3. Cross-Entropy Viewpoint:

↳. If we write out the training error:

$$E_{in}(\underline{w}) = -\frac{1}{N}\sum_{n=1}^{N}\left[\; \mathbb{1}(y_n = +1)\;\log \hat{P}_{\underline{w}}(1\mid \underline{x}_n)\right. \qquad \text{// +ve case}$$

$$\left. + \mathbb{1}(y_n = -1)\;\log \hat{P}_{\underline{w}}(-1\mid \underline{x}_n)\right] \qquad \text{// -ve case}$$

This is exactly the form of cross-entropy

↳ Def:

Suppose $P$ and $Q$ are two prob. distributions over $X = \{x_1, \cdots, x_m\}$

// e.g. Poission with mean $\alpha$.

The set of possible value $X = \{0, 1, 2, \cdots\}$

$$P(x) = \frac{\alpha^x}{x!}\, e^{-\alpha}$$

The cross-entropy between $P$ and $Q$ (measurement of distance) is:

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

↳. Relation between CE and $E_{in}(\underline{w})$:

For the nth example, consider the distribution

$$P_n = (\Pr\{y_n = +1\}, \; \Pr\{y_n = -1\}) \qquad \text{// true prob. distribution for } y_n$$

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

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

Let $Q_n = (\hat{P}_{\underline{w}}(+1\mid \underline{x}_n), \; \hat{P}_{\underline{w}}(-1\mid \underline{x}_n))$    // Our estimate for prob. distrb.

of $y_n$ given $\underline{x}_n$

$$\Rightarrow \boxed{E_{in}(\underline{w}) = \frac{1}{N}\sum_{n=1}^{N} CE(P_n, Q_n)}$$

∴ Minimizing $E_{in}(\underline{w})$ = Minimizing distance between $P_n$ and $Q_n$

↑ true distrb of $y_n$      ↑ estimate distrb. of $y_n$