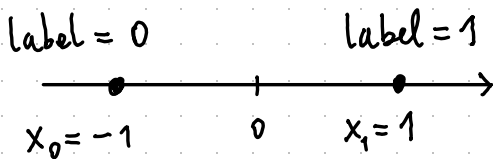


QUESTION: Does a minimum always exist for the cross-entropy loss in classification?

ANSWER: In general, no.

Example: Consider a dataset of $N=2$ points.

$$\text{features: } \begin{cases} x_0 = -1 \\ x_1 = 1 \end{cases} \quad \text{labels: } \begin{cases} y_0 = 0 \\ y_1 = 1 \end{cases}$$



$w \in \mathbb{R}$ weight

$b \in \mathbb{R}$ bias

The cost for this binary classification is:

$$\begin{aligned} L(w, b) &= - \sum_{i=0}^{N-1} \left(y_i \log(\sigma(wx_i + b)) + (1 - y_i) \log(1 - \sigma(wx_i + b)) \right) \\ &= \underbrace{-\log(1 - \sigma(-w + b))}_{i=0} - \underbrace{\log(\sigma(w + b))}_{i=1} \\ &= -\log((1 - \sigma(-w + b)) \sigma(w + b)) \end{aligned}$$

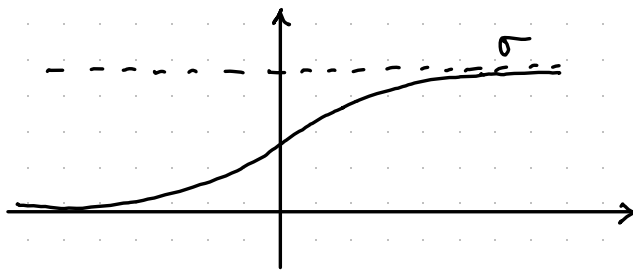
Let us show that $L(w, b)$ does not have a minimum.

Since $0 < \sigma < 1$, we have
 $-\log((1-\sigma)\sigma) > 0$.

Hence, the value 0 is not reached for any w and b .

However, choosing $b=0$ and letting $w \rightarrow +\infty$, we have

$$\lim_{w \rightarrow +\infty} \sigma(w) = 1, \quad \lim_{w \rightarrow +\infty} \sigma(-w) = 0$$



This implies

$$\lim_{w \rightarrow +\infty} (1 - \sigma(-w)) \sigma(w) = 1$$

and thus

$$\lim_{w \rightarrow +\infty} L(w, 0) = 0.$$

It means that $\inf_{w, b} L(w, b) = 0$,
but it is not reached.

Intuition: The dataset is very discriminatory.
From the dataset, it seems that
the probability distribution behind
data is

$$P(-1, 0) = 1$$

$$P(-1, 1) = 0$$

$$P(+1, 1) = 1$$

$$P(+1, 0) = 0$$

This can be approximated using the
logistic regression.

$$Q(-1, 0; w, b) = \sigma(-w + b) P(X = -1)$$

$$Q(-1, 1; w, b) = (1 - \sigma(-w + b)) P(X = -1)$$

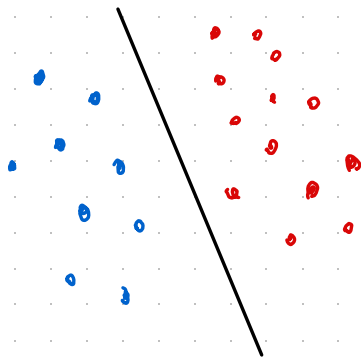
$$Q(+1, 0; w, b) = (1 - \sigma(w + b)) P(X = 1)$$

$$Q(+1, 1; w, b) = \sigma(w + b) P(X = 1)$$



We can approximate
the probabilities P
with $w \rightarrow +\infty$.

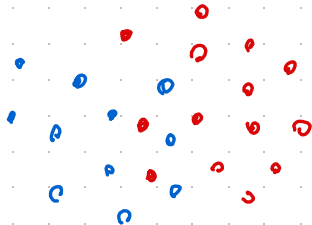
The same reasoning works in a
more general case when data can be
linearly separated.



In some sense, the logistic regression model can "potentially" model perfectly the actual probability distribution, but in a limit $w \rightarrow +\infty$.

Remark: Even if the minimum does not exist, a numerical optimization algorithm may allow to approach the infimum value.

In the next example we show that, when the dataset is not linearly separated, the minimum exists.



Example:

$$x_0 = -1$$

$$y_0 = 0$$

$$x_1 = 1$$

$$y_1 = 1$$

$$N = 3$$

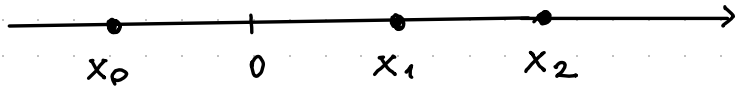
$$x_2 = 2$$

$$y_2 = 0$$

label = 0

label = 1

label = 0



$$L(w, b) = - \sum_{i=0}^{N-1} \left(y_i \log(\sigma(wx_i + b)) + (1 - y_i) \log(1 - \sigma(wx_i + b)) \right)$$

$$= \underbrace{-\log(1 - \sigma(-w + b))}_{i=0} - \underbrace{\log(\sigma(w + b))}_{i=1} +$$

$$\underbrace{-\log(1 - \sigma(2w + b))}_{i=2} =$$

$$= -\log\left(1 - \frac{1}{e^{w-b} + 1}\right) - \log\left(\frac{1}{e^{-w-b} + 1}\right) +$$

$$-\log\left(1 - \frac{1}{e^{-2w-b} + 1}\right) =$$

$$\begin{aligned}
&= -\log\left(\frac{e^{w-b}}{e^{w-b}+1}\right) - \log\left(\frac{1}{e^{-w-b}+1}\right) + \\
&\quad - \log\left(\frac{\bar{e}^{2w-b}}{\bar{e}^{2w-b}+1}\right) = \\
&= -\log\left(\frac{1}{1+e^{-w+b}}\right) - \log\left(\frac{1}{1+e^{-w-b}}\right) + \\
&\quad - \log\left(\frac{1}{1+e^{2w+b}}\right) = \\
&= \log(1+e^{-w+b}) + \log(1+e^{-w-b}) + \log(1+e^{2w+b})
\end{aligned}$$

We exploit the "misclassified" point.

Consider a sequence (w_k, b_k) such that $|w_k| + |b_k| \rightarrow +\infty$.

In the following cases, we assume $|w_k| + |b_k|$ is large.

We want to show that $L(w_k, b_k)$ is large. This allows us to search for a minimum in a bounded region, which exists by Weierstrass' theorem.

Assume that

$$\log(1 + e^{-w_k + b_k}) + \log(1 + e^{-w_k - b_k}) + \log(1 + e^{2w_k + b_k})$$

is bounded from above. (We know $L > 0$)

Since $t \mapsto \log(1 + e^t)$ is increasing, it means that

$$\begin{cases} -w_k + b_k \leq M & (1) \\ -w_k - b_k \leq M & (2) \\ 2w_k + b_k \leq M & (3) \end{cases}$$

for some $M > 0$ independent of k .

Summing (1)+(2):

$$-2w_k \leq 2M \Rightarrow -w_k \leq M \Rightarrow -M \leq w_k$$

Summing (2)+(3)

$$w_k \leq 2M$$

Hence $|w_k| \leq 2M$.

From (1), $b_k \leq M + w_k \leq 3M$

From (2), $-b_k \leq M + w_k \leq 3M$

It follows that $|w_k| + |b_k| \leq 4M$ is bounded, contradicting the fact that it is large.

In fact, we can say more precisely that

$$\lim_{|w_k, b_k| \rightarrow +\infty} L(w_k, b_k) = +\infty.$$

To show this, assume that

$$|w_k| + |b_k| > 4M$$

Then one of the three inequalities

$$-w_k + b_k \leq M \quad (1)$$

$$-w_k - b_k \leq M \quad (2)$$

$$2w_k + b_k \leq M \quad (3)$$

must be violated (for, otherwise, we have shown that $|w_k| + |b_k| \leq 4M$)

If (1) is violated, $-w_k + b_k > M \Rightarrow$
 $\Rightarrow \log(1 + e^{-w_k + b_k}) > \log(1 + e^M)$

If (2) is violated,

$$\log(1 + e^{-w_k - b_k}) > \log(1 + e^M)$$

If (3) is violated,

$$\log(1 + e^{2w_k + b_k}) > \log(1 + e^M)$$

In either case,

$$L(w_k, b_k) > \underbrace{\log(1 + e^M)}$$

Since this term $\rightarrow +\infty$
as $M \rightarrow +\infty$, we
have shown concavity
of L .