

LOGISTIC REGRESSION (binary)

Assume that we have features described by a random vector $X \in \mathbb{R}^M$ and a random variable Y with range in $\{0, 1\}$, $Y \sim \text{Be}(p)$.

- X represents features
- Y represents the label of classification.

We are interested in finding

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

The variable Y is dependent from the variable X .

What we want to reconstruct is thus

$$P(Y=1 \mid X=x) =: p(x)$$

This means that we want to reconstruct the probability of classifying with "1" when we observe some "features" x .

Remark: By the law of total probability

$$\sum_x \mathbb{P}(Y=1 | X=x) \mathbb{P}(X=x) = p \quad \leftarrow \text{the one of } \text{Be}(p)$$

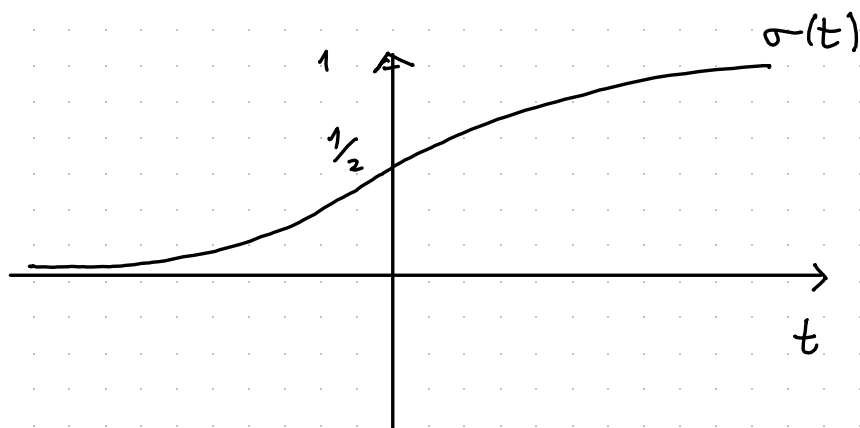
The logistic model makes the following ansatz:

$$\mathbb{P}(Y=1 | X=x) = p(x) \approx \sigma(xw + b) \\ =: q(x; w, b)$$

where

- $x \in \mathbb{R}^{1 \times M}$ is the observed feature
- $w \in \mathbb{R}^{M \times 1}$ are weights
- $b \in \mathbb{R}$ is a bias
- σ is the sigmoid / logistic function

$$\sigma(t) = \frac{1}{1 + e^{-t}} = \frac{e^t}{e^t + 1}$$



Assume that we have observed data $\{(x_i, y_i)\}_{i=0, \dots, N-1}$, where

- $x_0, \dots, x_{N-1} \in \mathbb{R}^{1 \times M}$ are realizations of a random sample drawn from X
- $y_0, \dots, y_{N-1} \in \{0, 1\}$ are the corresponding labels.

We are making an ansatz on the model:

$$\mathbb{P}(Y=1 \mid X=x) = \sigma(xw + b) \\ =: q(x; w, b)$$

where w, b are the parameters.

We want to find the MLE of w, b and compute its realization on the observed data.

The $-\log$ -likelihood is:

$$-\log \left(P \left(\bigcap_{i=0}^{N-1} \{X_i = x_i, Y_i = y_i\} \right) \right) =$$

$$= -\log \left(\prod_{i=0}^{N-1} P(\{X_i = x_i, Y_i = y_i\}) \right) =$$

$$= -\sum_{i=0}^{N-1} \log \left(P(Y_i = y_i | X_i = x_i) P(X_i = x_i) \right)$$

$$= -\sum_{i=0}^{N-1} \log \left(P(Y_i = y_i | X_i = x_i) \right) +$$

$$-\sum_{i=0}^{N-1} \log P(X_i = x_i)$$



This quantity is independent
of the parameters W, b .

It is not interesting for the
optimization problem

We study only the first term.

$$\begin{aligned}
& - \sum_{i=0}^{N-1} \log(P(Y=y_i | X=x_i)) = \\
& = - \sum_{i=0}^{N-1} \log(q(x_i; w, b)^{y_i} (1-q(x_i; w, b))^{1-y_i}) = \\
& = - \sum_{i=0}^{N-1} y_i \log q(x_i; w, b) + (1-y_i) \log(1-q(x_i; w, b))
\end{aligned}$$

The part of the negative log-likelihood depending on the parameters w, b is:

$$\begin{aligned}
L(w, b; \{(x_i, y_i)\}_{i=0, \dots, N-1}) &= \\
&= - \sum_{i=0}^{N-1} y_i \log q(x_i; w, b) + (1-y_i) \log(1-q(x_i; w, b))
\end{aligned}$$

To minimize this, we need numerical methods (we will see them later).

CROSS-ENTROPY

Cross-entropy is a tool in information theory.

Assume that Ω is the sample space (discrete) and P and Q are two probability measures on Ω .

Let us set:

$$p(\omega) := P(\{\omega\}), \quad q(\omega) := Q(\{\omega\})$$

The cross-entropy is defined by

$$H(p, q) = \mathbb{E}_p(-\log q) =$$

in this formula

$$= - \sum_{\omega \in \Omega} p(\omega) \log q(\omega) \quad -\log q: \Omega \rightarrow \mathbb{R} \text{ is treated as a r.v.}$$

Let us understand how to interpret this.
First of all,

$$H(p) = H(p|p) = - \sum_{\omega \in \Omega} p(\omega) \log p(\omega)$$

This is the entropy of the probability distribution $p(\omega)$.

The entropy is the expected amount of information carried by events sampled from the distribution.

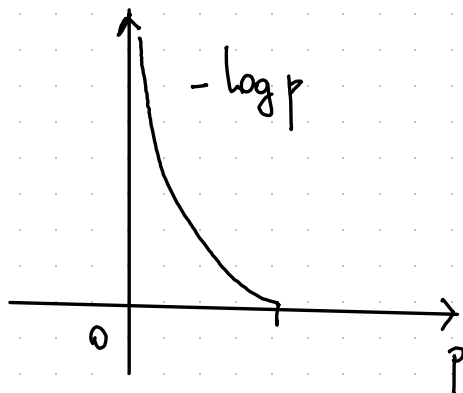
The idea is the following:

- very likely events carry a little amount of information
- very unlikely events carry a lot of information

A function that models this is

$$-\log p(w)$$

for a discrete event x .



In some sense, this measures how much we are "surprised" of observing the discrete event w .

Let us go back to cross entropy:

$$\begin{aligned} H(p, q) &= - \sum_{\omega \in \Omega} p(\omega) \log q(\omega) = \\ &= - \sum_{\omega \in \Omega} p(\omega) \left(\log \frac{q(\omega)}{p(\omega)} + \log p(\omega) \right) = \\ &= - \sum_{\omega \in \Omega} p(\omega) \log \frac{q(\omega)}{p(\omega)} + H(p) \\ &= \underbrace{D_{KL}(p \parallel q)} + H(p) \end{aligned}$$

This is called Kullback-Leibler divergence

We use the interpretation of "surprise" to explain the meaning of the Kullback-Leibler divergence.

First of all, we assume that

P is the probability distribution underlying the described phenomenon.

Remark: $D_{KL}(p \parallel q) \geq 0$. This is known as Gibb's inequality.

A proof is the following:

\log is a concave function.

By Jensen's inequality

$$\sum_{\omega \in \Omega} p(\omega) \log \frac{q(\omega)}{p(\omega)} \leq \log \left(\sum_{\omega \in \Omega} q(\omega) \right) = 0.$$

Equality holds true if $\frac{q(\omega)}{p(\omega)} = 1$ and only if

$\frac{q(\omega)}{p(\omega)}$ is constant.

- If $\frac{q(\omega)}{p(\omega)} \approx 1$, this means

$q(\omega) \approx p(\omega)$ - Using Q to estimate the probability of ω instead of the "true" probability P is a good choice. We pay a low cost.

- If $\frac{q(\omega)}{p(\omega)}$ far from 1, then we are making a mistake. With $-\log \frac{q(\omega)}{p(\omega)}$ we

measure how much we are surprised of seeing the event w drawn from Q relatively to how much we are surprised if the probability distribution is P .
Then

$D_{KL}(P \parallel Q) = \mathbb{E}[-\log q + \log p]$
is the expected excess in surprise.

This can be thought as follows:

- There is an underlying distribution for the data: P
- We use a model Q
- We measure the mistake we are making by using P instead of Q by $D_{KL}(P \parallel Q)$.

Minimizing the cross-entropy with respect to q is equivalent to minimizing the DKL.

Cross-entropy for binary classification

Assume that we have

- X random vector for the "explaining" features
- Y random variable explained by X ,
 $Y \sim \text{Be}(p)$

Here we have the following situation:
The whole vector (X, Y) is distributed according to some law.
If we are working in the discrete setting, there is a joint probability mass function

$$\tilde{p}(x, y) = \mathbb{P}(X=x, Y=y)$$

About the joint probability mass function, we know that:

$$\begin{aligned}\tilde{p}(x, y) &= P(X=x, Y=y) = \\ &= P(Y=y | X=x) P(X=x)\end{aligned}$$

$$\tilde{p}(x, 1) = p(x) P(X=x)$$

$$\tilde{p}(x, 0) = (1-p(x)) P(X=x)$$

Instead of \tilde{p} , the logistic model uses \tilde{q} for which

$p(x)$ is replaced by

$$q(x; w, b) := \sigma(wx + b)$$

Assuming that the distribution of X is the same:

$$\tilde{q}(x, 1; w, b) = q(x) P(X=x)$$

$$\tilde{q}(x, 0; w, b) = (1-q(x)) P(X=x)$$

The cross-entropy is

$$H(\tilde{p}, \tilde{q}) = -\mathbb{E}_{\tilde{p}}(\log \tilde{q})$$

We are thinking at $-\log(\tilde{q})$ as a random variable, in the following sense:

$$\omega \in \Omega \mapsto -\log(\tilde{q}(X(\omega), Y(\omega)))$$

Let us write more explicitly the expression of the cross-entropy:

$$\begin{aligned} H(\tilde{p}, \tilde{q}) &= - \sum_x \sum_y \tilde{p}(x, y) \log \tilde{q}(x, y; w, b) = \\ &= \left\{ Y \sim \text{Be}(p) \right\} = \\ &= - \sum_x \left(\tilde{p}(x, 1) \log \tilde{q}(x, 1; w, b) + \tilde{p}(x, 0) \log \tilde{q}(x, 0; w, b) \right) \end{aligned}$$

Lemma: Let (X, Y) be distributed with \tilde{p} .

$$\begin{aligned} & -\mathbb{E}_{\tilde{p}}[Y \log(q(X; w, b)) \\ & \quad + (1-Y) \log(1 - q(X; w, b))] = \\ & = H(\tilde{p}, \tilde{q}) \end{aligned}$$

Proof:

$$\begin{aligned} & -\sum_x \sum_y \tilde{p}(x, y) (y \log(q(x)) + (1-y) \log(1 - q(x))) \\ & = -\sum_x \tilde{p}(x, 1) \log q(x) + \tilde{p}(x, 0) \log(1 - q(x)) \end{aligned}$$

□

We cannot compute explicitly the cross-entropy, but we have data at our disposal $\{(x_i, y_i)\}_{i=0, \dots, N-1}$. These are realizations of a random sample $\{(X_i, Y_i)\}_{i=0, \dots, N-1}$, i.e., i.i.d. random vectors distributed with \tilde{p} . We consider the random variables:

$$-(Y_i \log q(X_i; w, b) + (1 - Y_i) \log (1 - q(X_i; w, b)))$$

and their empirical average:

$$-\frac{1}{N} \sum_{i=0}^{N-1} (Y_i \log q(X_i; w, b) + (1 - Y_i) \log (1 - q(X_i; w, b)))$$

By the law of large numbers, for $N \rightarrow +\infty$ the empirical average converges to the mean of the population, which, by the Lemma, is precisely $H(\tilde{p}, \tilde{q})$.

(In fact, it is also an unbiased estimator, since its expectation is $H(\tilde{p}, \tilde{q})$).

This allows us to use

$$-\frac{1}{N} \sum_{i=0}^{N-1} (y_i \log(x_i; w, b) + (1 - y_i) \log(1 - q(x_i; w, b)))$$

as an estimate of the cross-entropy.

Minimizing the empirical cross-entropy is therefore equivalent to minimizing the $-\log$ -likelihood of the data.

LOGISTIC REGRESSION (multiclass)

In this situation we have

- a random vector $X: \Omega \rightarrow \mathbb{R}^{1 \times M}$
(the features)
- a random label $Y: \Omega \rightarrow \{0, 1, \dots, K-1\}$

Again, (X, Y) is distributed according to a probability distribution with probability mass function $\tilde{p}: \mathbb{R}^M \times \{0, 1, \dots, K-1\} \rightarrow [0, 1]$, $\tilde{p}(x, y)$.

We are interested in studying

$$P(Y = k \mid X = x)$$

The logistic regression model is based on the ansatz:

$$P(Y = k \mid X = x) \approx q_k(x; \underbrace{W, b}_{\text{parameters}})$$

parameters
that we discuss
later

Let us now understand the structure of $q_k(x; W, b)$.

In the binary logistic regression, we had to define only the probability of falling in class 1, given by

$$\sigma(\underbrace{xW + b}_{\text{logit}})$$

Now we have classes $0, 1, \dots, K-1$.

We make the ansatz that the logit for class $k = 0, \dots, K-1$ is of the form

$$xW_k + b_k$$

where

- $x \in \mathbb{R}^{1 \times M}$ are features
- $w_k \in \mathbb{R}^{M \times 1}$ are weights, column of $W \in \mathbb{R}^{M \times K}$
- $b_k \in \mathbb{R}$ are biases, entry of $b \in \mathbb{R}^K$

Then we transform $xw_k + b_k$ in a number in $[0, 1]$. This must be consistent with the fact that the probabilities sum to one.

$$\sigma_k(xw_k + b_k) = c e^{xw_k + b_k}$$

$$\sum_{h=0}^{K-1} c e^{xw_h + b_h} = 1 \Rightarrow c = \frac{1}{\sum_{h=0}^{K-1} e^{xw_h + b_h}}$$

It follows that

$$\sigma_k(xw + b) = \frac{e^{xw_k + b_k}}{\sum_{h=0}^{K-1} e^{xw_h + b_h}}$$

We have defined a function

$$\sigma: \mathbb{R}^K \rightarrow [0, 1]^K \text{ given by}$$
$$\sigma(a)_k = \sigma_k(a) = \frac{e^{a_k}}{\sum_{h=0}^{K-1} e^{a_h}}$$

This function is called softmax
(in fact, ~~softmax~~ softmax)

As before, we have some observed data $\{(x_i, y_i)\}_{i=0, \dots, N-1}$.

We can compute $-\log$ -likelihood of the data:

$$-\log\left(\mathbb{P}\left(\bigcap_{i=0}^{N-1} \{X_i = x_i, Y_i = y_i\}\right)\right) =$$

$$= -\log\left(\prod_{i=0}^{N-1} \mathbb{P}(X_i = x_i, Y_i = y_i)\right) =$$

$$= -\sum_{i=0}^{N-1} \log(\mathbb{P}(Y_i = y_i | X_i = x_i) \mathbb{P}(X_i = x_i))$$

$$= -\sum_{i=0}^{N-1} \log(\mathbb{P}(Y_i = y_i | X_i = x_i)) +$$

$$-\underbrace{\sum_{i=0}^{N-1} \log(\mathbb{P}(X_i = x_i))}_{\text{this is independent of the parameters}}$$

this is independent
of the parameters

Let us study only the first term:

$$- \sum_{i=0}^{N-1} \log(P(Y_i = y_i | X_i = x_i))$$

To express the probability more explicitly, it is convenient to use a one-hot encoding for the variable Y .

We define the matrix $\tilde{y} \in \mathbb{R}^{N \times K}$ such that the row $\tilde{y}_{i\cdot} \in \mathbb{R}^{1 \times K}$ is given by $\tilde{y}_{i\cdot} = e_k$ if $y_i = e_k$.

$$\begin{aligned} & - \sum_{i=0}^{N-1} \log(P(Y_i = y_i | X_i = x_i)) = \\ & = - \sum_{i=0}^{N-1} \log(P(\tilde{Y}_i = \tilde{y}_{i\cdot} | X_i = x_i)) = \\ & = - \sum_{i=0}^{N-1} \log\left(\prod_{k=0}^{K-1} \sigma_k(xW + b)^{\tilde{y}_{ik}}\right) \\ & = - \sum_{i=0}^{N-1} \sum_{k=0}^{K-1} \tilde{y}_{ik} \log \sigma_k(xW + b) \end{aligned}$$

Recall that $w_0 = 0$, $b_0 = 0$,
hence the 0 class is not affecting
the optimisation.

Hence multiclass logistic regression
has the objective to find

$$W \in \mathbb{R}^{M \times K}, \quad b \in \mathbb{R}^K$$

such that

$$(W, b) \in \operatorname{argmin} L(W, b; \{(x_i, y_i)\}_{i=0, \dots, N-1})$$

where

$$\begin{aligned} L(W, b; \{(x_i, y_i)\}_{i=0, \dots, N-1}) \\ = - \sum_{i=0}^{N-1} \sum_{k=1}^{K-1} \underbrace{\tilde{y}_{ik}}_{\text{one-hot encoding}} \underbrace{\log \sigma_k(xW + b)}_{\text{softmax}} \end{aligned}$$

Let us interpret this using the cross-entropy.

The data $\{(x_i, y_i)\}_{i=0, \dots, N-1}$ is the realization of a random sample $(X_0, Y_0), \dots, (X_{N-1}, Y_{N-1})$ extracted from a probability distribution with joint probability mass function $\tilde{p}(x, y)$.

The random variables Y_i have range $\{0, 1, \dots, K-1\}$.

Instead of $\tilde{p}(x, y)$, we use a model which induces a joint probability distribution $\tilde{q}(x, y; W, b)$

$$\tilde{q}(x, k; W, b) = \sigma_k(xW + b) \mathbb{P}(X = x)$$

The cross-entropy is

$$\begin{aligned} H(\tilde{p}, \tilde{q}) &= -\mathbb{E}_{\tilde{p}}(\log \tilde{q}) = \\ &= -\sum_x \sum_{k=0}^{K-1} \tilde{p}(x, k) \log \tilde{q}(x, k; W, b) \end{aligned}$$

Let us use the notation \tilde{Y} to denote the one-hot encoding of Y .

By $\text{log}\sigma(xW+b)$ we mean the vector $(\log \sigma_k(xW+b))_{k=0, \dots, K-1}$.

Lemma: We have that (entropy of X)

$$\mathbb{E}_{\tilde{p}}(-\tilde{Y} \cdot \text{log}\sigma(xW+b)) = H(\tilde{p}, \tilde{q}) - \overbrace{H(p_x)}^{\text{entropy of } X}$$

Proof: interpreted as a vector

$$\begin{aligned} & \mathbb{E}_{\tilde{p}}(-\tilde{Y} \cdot \text{log}\sigma(Wx+b)) = \\ &= - \sum_x \sum_{k=0}^{K-1} \tilde{p}(x, k) \tilde{k} \cdot \log \sigma(Wx+b) = \\ &= - \sum_x \sum_{k=0}^{K-1} \tilde{p}(x, k) \log \sigma_k(Wx+b) \\ &= - \sum_x \sum_{k=0}^{K-1} \tilde{p}(x, k) \log \frac{\tilde{q}(x, k; W, b)}{P(X=x)} = \end{aligned}$$

...

$$\begin{aligned}
 & H(\tilde{p}, \tilde{q}) \\
 = & \underbrace{- \sum_x \sum_{\kappa=0}^{K-1} \tilde{p}(x, \kappa) \log \tilde{q}(x, \kappa; W, b)}_{\text{cross entropy}} + \\
 & \underbrace{\sum_x \sum_{\kappa=0}^{K-1} \tilde{p}(x, \kappa) \log P(X=x)}_{\text{entropy of } X} \\
 = & \sum_x P(X=x) \log P(X=x) = \\
 = & -H(p_X) \text{ entropy} \\
 & \text{of the distribution} \\
 & X \quad \square
 \end{aligned}$$

By the law of large numbers,

$$-\frac{1}{N} \sum_{i=0}^{N-1} \tilde{Y}_i \cdot \log \sigma(WX_i + b)$$

converges to $H(\tilde{p}, \tilde{q}) - H(p_X)$.

The realization on the dataset is

$$-\frac{1}{N} \sum_{i=0}^{N-1} \sum_{\kappa=0}^{K-1} \tilde{y}_{i\kappa} \log \sigma_{\kappa}(Wx + b).$$

We conclude that minimizing the cross-entropy is equivalent to minimizing the log-likelihood.