# Lecture Notes: Binary logistic regression

## CS375: NLP / Williams College / Spring 2023

Let's derive our loss function (also sometimes call the called *objective function*) for binary logistic regression: negative log likelihood (also called *cross entropy*). The loss function is how we fit the weights, $\theta$ for our logistic regression classifier.

**Notation.** We begin with a bit of notation. Our *training data* consists of a matrix, $X$, the input text data, and $\vec{y}$, the labels (typically annotations on the documents by humans). For simplicity, we drop the $\vec{\phantom{y}}$ and call $\vec{y}$ just $y$ going forward. We index into a single element of these (a single document and a single label) with $i$, resulting in $(\vec{x}_i, y_i)$ where $y_i \in \{0, 1\}$. In a bag-of-words feature representation, $|\vec{x}| = |V|$ where $V$ is the set of vocabulary words.

**Set-up.** Recall, the weights in our logistic regression model are $\theta$ and our model gives a prediction for the probability of the positive class given the input

$$\hat{p}_i := P_\theta(y = 1 | x_i) = \frac{1}{1 + e^{-x_i \cdot \theta}} \tag{1}$$

Now let's use the fact that $y_i$ is binary to write a generic probability that incorporates the case in which $y_i = 1$ and $y_i = 0$. We can rewrite this as

$$p_\theta(y_i | x_i) = \hat{p}_i^{y_i} (1 - \hat{p}_i)^{1 - y_i} \tag{2}$$

Why does this work? Well suppose the true label $y_i = 1$ then we have

$$p_\theta(y_i | x_i) = \hat{p}_i^1 (1 - \hat{p}_i)^{1-1} \tag{3}$$
$$= \hat{p}_i \tag{4}$$

If instead $y_i = 0$ then

$$p_\theta(y_i | x_i) = \hat{p}_i^0 (1 - \hat{p}_i)^{1-0} \tag{5}$$
$$= 1 - \hat{p}_i \tag{6}$$

Both of these are true for how we defined $\hat{p}_i$ in Equation 1.

**Loss function.** We want to create a function that says

$$L(\hat{p}_i, y_i) = \text{``Distance'' of our predicted probability to the true class} \tag{7}$$

We want want to minimize loss and minimize distance.

We'll use principles of maximum likelihood estimation to define this loss function. We want to set parameters $\theta$ that maximize the likelihood of the data, $P_\theta(y_i | x_i)$. For reasons we'll see later, this is equivalent to saying we want to minimize the negative likelihood (NL):

$$NL(\hat{p}_i, y_i) = -P_\theta(y_i | x_i) \tag{8}$$
$$= \hat{p}_i^{y_i} (1 - \hat{p}_i)^{1 - y_i} \tag{9}$$

by substituting Equation 2.

Like we've seen before, it'll be easier to use logs in implementation so we take the log of both sides and this becomes the negative log likelihood (NLL)

$$NLL(\hat{p}_i, y_i) = -\log\left(\hat{p}_i^{y_i}(1-\hat{p}_i)^{1-y_i}\right) \tag{10}$$

$$= -y_i\log\hat{p}_i - (1-y_i)\log(1-\hat{p}_i) \tag{11}$$

To maximize this probability across the entire dataset, we can take the average across all examples $i$,

$$NLL(\vec{\hat{p}}, \vec{y}) = \frac{1}{n}\sum_{i=1,2,\cdots,n}\left(-y_i\log\hat{p}_i - (1-y_i)\log(1-\hat{p}_i)\right) \tag{12}$$

**Intution.** Why does the negative log likelihood work?

Suppose the true label is $y_i = 1$ and our model predicts $\hat{p}_i = 0.99$.

Then from Equation 11 we have

$$NLL(0.99, 1) = -1\log(0.99) - (1-1)\log(1-0.99) \tag{13}$$

$$= -1\log(0.99) \tag{14}$$

$$= -1*(-0.01) \tag{15}$$

$$= 0.01 \tag{16}$$

which is very close to zero and intuitively we consider this "successful."

If instead we predicted for this same example $i$, $\hat{p}_i = 0.6$ we have

$$NLL(0.6, 1) = -1\log(0.6) - (1-1)\log(1-0.6) \tag{17}$$

$$= -1\log(0.6) \tag{18}$$

$$= -1*(-0.51) \tag{19}$$

$$= 0.51 \tag{20}$$

this gives us a higher number (which is worse) indicating that we were not as successful with our model.

Try out similar examples for yourself for $y_i = 0$.

---

We use **gradient descent** to find the weights $\theta$ that minimize the negative log likelihood objective function. Let's keep things simple and just look at this gradient for a single document $i$

$$\hat{\theta} = \operatorname*{argmin}_{\theta} NLL(\hat{p}_i, y_i) \tag{21}$$

$$= \operatorname*{argmin}_{\theta}\left(-y_i\log\hat{p}_i - (1-y_i)\log(1-\hat{p}_i)\right) \tag{22}$$

$$= \operatorname*{argmin}_{\theta}\left(-y_i\log\left(\frac{1}{1+e^{-x_i\cdot\theta}}\right) - (1-y_i)\log\left(1 - \frac{1}{1+e^{-x_i\cdot\theta}}\right)\right) \tag{23}$$

Now we need the derivative with respect to $\theta$

$$\frac{\partial}{\partial\theta}\left(-y_i\log\left(\frac{1}{1+e^{-x_i\cdot\theta}}\right) - (1-y_i)\log\left(1 - \frac{1}{1+e^{-x_i\cdot\theta}}\right)\right) \tag{24}$$

To be continued in class examples...