

Homework 8
Mouselinos Spyridon
February 2020

Exercise 1

In logistic regression in the form of two-class case we assume that the log-ratio of the posteriors is a linear function of x .

The assumed model is thus:

$$\ln \frac{P(w_1|x)}{P(w_0|x)} = \theta^T x$$

We also know that because we have two classes $w_1, w_0 \rightarrow P(w_1|x) + P(w_0|x) = 1$

so if $\sigma(t) = \frac{1}{1+e^{-t}}$, then:

$$P(w_1|x) = \sigma(\theta^T x)$$

$$P(w_0|x) = 1 - \sigma(\theta^T x)$$

Exercise 1 (cont)

Now, given a set of N training examples where $Y = (y_n, x_n)$ with $y_n \in \{0, 1\}$

and $x_n \in \mathbb{R}^d$. If we assume statistical independence and try to solve the problem by ML we have:

$$\text{Likelihood}(\theta) = \prod_{n=1}^N P(y_n | x_n; \theta) \Rightarrow$$

$$\Rightarrow \text{Log Likelihood}(\theta) = \ln \left(\prod_{n=1}^N P(y_n | x_n; \theta) \right)$$

$$\Rightarrow L(\theta) = \ln \prod_{n=1}^N P(y_n | x_n; \theta) =$$

$$= \ln \prod_{n=1}^N P(1 | x_n; \theta)^{y_n} P(0 | x_n; \theta)^{1-y_n}$$

Just like a Bernoulli Trial

$$\Rightarrow \ln \prod_{n=1}^N \sigma(\theta^T x_n)^{y_n} (1 - \sigma(\theta^T x_n))^{1-y_n}$$

$$= \ln \prod_{n=1}^N s_n^{y_n} (1 - s_n)^{1-y_n}$$

where $s_n = \sigma(\theta^T x_n)$

Exercise 1 (cont)

Now the negative log likelihood is

$$NL(\theta) = - \sum_{n=1}^N y_n \ln s_n + (1-y_n) \ln(1-s_n)$$

The gradient wrt θ is:

~~$$NL(\theta) = - \sum_{n=1}^N y_n \ln s_n + (1-y_n) \ln(1-s_n)$$~~

$$\nabla_{\theta} NL(\theta) = - \sum_{n=1}^N \frac{d}{d\theta} (y_n \ln s_n) + \frac{d}{d\theta} [(1-y_n) \ln(1-s_n)]$$

$$\nabla_{\theta} NL(\theta) = - \sum_{n=1}^N y_n \frac{d}{d\theta} \ln s_n + (1-y_n) \frac{d}{d\theta} \ln(1-s_n)$$

But:

$$\begin{aligned} \frac{d}{d\theta} \ln s_n &= \frac{d}{d\theta} (s_n) \cdot \frac{1}{s_n} = \frac{1}{\sigma(\theta^T x)} \frac{d}{d\theta} \sigma(\theta^T x) \\ &= \frac{1}{\sigma(\theta^T x)} \cdot (\sigma(\theta^T x) (1 - \sigma(\theta^T x)) x) = (1 - \sigma(\theta^T x)) x \end{aligned}$$

In the same fashion:

$$\frac{d}{d\theta} \ln(1-s_n) = \frac{1}{1-s_n} \frac{d}{d\theta} (1-s_n)$$

$$= \frac{1}{1-\sigma(\theta^T x)} (-\sigma(\theta^T x) (1-\sigma(\theta^T x)) x) = (-\sigma(\theta^T x) x)$$

So:

$$\nabla_{\theta} L(\theta) = - \sum_{n=1}^N (1-\sigma(\theta^T x)) y_n x_n + (1-y_n) (-\sigma(\theta^T x) x)$$

$$\nabla_{\theta} L(\theta) = - \sum_{n=1}^N y_n x_n - y_n \sigma(\theta^T x) x_n + (-\sigma(\theta^T x) x + \sigma(\theta^T x) x_n y_n)$$

$$\nabla_{\theta} L(\theta) = - \sum_{n=1}^N (y_n - \sigma(\theta^T x)) x_n = \sum_{n=1}^N (\sigma(\theta^T x) - y_n) x_n$$

So finally:

$$\nabla_{\theta} NL(\theta) = X^T (S - y), \text{ where}$$

$$X^T = [x_1, x_2, \dots, x_N], \quad S^T = [s_1, s_2, \dots, s_N]$$

$$\text{and } y^T = [y_1, y_2, \dots, y_N]$$

A gradient update is given by:

$$\theta^{(\text{time step } i)} = \theta^{(\text{time step } i-1)} - \left(\mu_i \nabla NL(\theta) \right) \Big|_{\theta = \theta^{(i-1)}}$$

Meaning that:

$$\theta^t = \theta^{t-1} - \mu_i X^T (s^{t-1} - y)$$

μ_i is the learning rate parameter
and t the time step of each
update.

Exercise 2

a) Is in the HW8.ipynb.

b) As shown in Ex1. the gradient descent step is:

$$\theta_i = \theta_{i-1} - \mu \nabla_{\theta} J(\theta)$$

where $J(\theta)$ is simply the cost function.

~~$$J(\theta) = \sum_{n=1}^N (y_n - \theta^T x_n)^2$$~~

However here under MSE:

$$J(\theta) = \sum_{n=1}^N (y_n - f(\theta^T x))^2$$

So:

$$\nabla_{\theta} J(\theta) = -2 \sum_{n=1}^N (y_n - f(\theta^T x)) \cdot (f(\theta^T x_n) (1 - f(\theta^T x_n)))$$

• $x_n \cdot n$

Exercise 2 (cont)

So:

$$\theta_i = \theta_{i-1} + 2\mu \sum_{n=1}^N (y_n - f(\theta^T x)) (f(\theta^T x) (1 - f(\theta^T x)))$$

(c) In order for the model to respond with a clear 0 we should have $f(\theta^T x) \approx 0$

thus $\frac{1}{1+e^{-\theta^T x \cdot a}} = 0$ that can happen

[illegible]

that is impossible for finite values of x and a .

the same: applies for the "1" response

Exercise 2 (cont)

(a) Having in mind the sigmoid function we can say that for a given set of (x, a) the model will respond with a value $f(\theta^T x)$. That value is the probability of the input to belong to class 0 or 1 on the threshold of 0.5, if $f(x^T \theta) \geq 0.5$ our model says class 1. if $f(x^T \theta) < 0.5$ class 0 respectively.

Exercise 2 (cont)

(e) A way of forcing our values close to 0 or 1 would be to use a large value of a as we showed on question (c).

Exercise 3

In this problem we are somewhat lucky, since the order in our dataset enables us to perform a "sort of" binary search in $O(\log(n))$ in order to find the threshold.

We keep 3 pointers:

low \leftarrow The most negative point

mid \leftarrow The middle

high \leftarrow The most positive point

Iteratively we check if the mid and the high point have different classes $(-1, +1)$. If they don't

~~the high point~~ the low is now the

mid point. We keep until the

low is at -1 the high at $+1$

and have convergence