

## Homework 4. Supervised learning

In Lab 5 you will use supervised learning to solve two classification problems. To prepare for the lab, in the homework you will compute some of the necessary expressions that needed to implement one of the learning algorithms—logistic regression.

(Binary) logistic regression provides a discriminative model that estimates the probability of selecting each of two actions,  $\mathcal{A} = \{0, 1\}$ , given a set of examples,  $\mathcal{D} = \{(\mathbf{x}_1, a_1), \dots, (\mathbf{x}_N, a_N)\}$ , with  $a_n \in \mathcal{A}, n = 1, \dots, N$  and  $\mathbf{x}_n \in \mathbb{R}^P, n = 1, \dots, N$ . In particular, each state  $\mathbf{x}_n$  is described as a vector with  $P$  attributes, where attribute  $p$  is a real value. Logistic regression assumes that the probability of action  $a = 1$  takes the form

$$\pi(1 \mid \mathbf{x}; \mathbf{w}, b) \stackrel{\text{def}}{=} \mathbb{P}[a = 1 \mid \mathbf{x} = \mathbf{x}; \mathbf{w}, b] = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}}, \quad (1)$$

where  $\mathbf{w}$  and  $b$  are the parameters to be learned.

It is a common practice to write

$$\mathbf{w}^\top \mathbf{x} + b = \underbrace{\begin{bmatrix} \mathbf{w}^\top & b \end{bmatrix}}_{\text{new } \mathbf{w}} \underbrace{\begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}}_{\text{new } \mathbf{x}}.$$

In other words, we augment  $\mathbf{x}$  with a component  $x_{P+1} \equiv 1$  so that the parameter  $b$  can be treated as one of the weights. Using this simplification, (1) becomes just

$$\pi(1 \mid \mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}},$$

where now  $\mathbf{w}, \mathbf{x} \in \mathbb{R}^{P+1}$ . We adopt this practice throughout.

Training logistic regression consists in finding the parameters  $\mathbf{w}$  that minimize the *negative log likelihood* of the data. In other words, and assuming the examples  $(x_n, a_n)$  to be independent, we want to compute

$$\mathbf{w}^* = \underset{\mathbf{w} \in \mathbb{R}^{P+1}}{\operatorname{argmin}} \{-\log \mathbb{P}[\mathcal{D}; \mathbf{w}]\} = \underset{\mathbf{w} \in \mathbb{R}^{P+1}}{\operatorname{argmin}} \left\{ -\log \prod_{n=1}^N \mathbb{P}[a_n \mid \mathbf{x}_n; \mathbf{w}] \right\}. \quad (2)$$

You will start by writing (2) as a function of the data in  $\mathcal{D}$  and the probability  $\pi(1 \mid \mathbf{x}, \mathbf{w})$ ; you will then compute the operators needed to implement a local search algorithm (namely, the Newton-Raphson algorithm).

### Exercise 1.

(a) Show that

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^{P+1}} \ell(\mathcal{D}; \pi),$$

where

$$\ell(\mathcal{D}; \pi) \stackrel{\text{def}}{=} \sum_{n=1}^N [a_n \log \pi(1 \mid \mathbf{x}_n; \mathbf{w}) + (1 - a_n) \log(1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w}))].$$

(b) Show that the gradient of  $\ell(\mathcal{D}; \pi)$  with respect to  $\mathbf{w}$  (i.e., the vector of first-order derivatives) is given by

$$\mathbf{g} = \sum_{n=1}^N \mathbf{x}_n (a_n - \pi(1 \mid \mathbf{x}_n; \mathbf{w}))$$

(c) Show that the Hessian of  $\ell(\mathcal{D}; \pi)$  with respect to  $\mathbf{w}$  (i.e., the matrix of second order derivatives) is given by

$$\mathbf{H} = - \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \pi(1 \mid \mathbf{x}_n; \mathbf{w}) (1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w})).$$

#### Solution 1:

(a) We have that

$$\pi(1 \mid \mathbf{x}_n; \mathbf{w})^{a_n} = \begin{cases} \pi(a_n \mid \mathbf{x}_n; \mathbf{w}) & \text{if } a_n = 1; \\ 1 & \text{otherwise.} \end{cases}$$

and

$$(1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w}))^{1-a_n} = \begin{cases} \pi(a_n \mid \mathbf{x}_n; \mathbf{w}) & \text{if } a_n = 0; \\ 1 & \text{otherwise.} \end{cases}$$

Therefore,

$$\pi(a_n \mid \mathbf{x}_n; \mathbf{w}) = \pi(1 \mid \mathbf{x}_n; \mathbf{w})^{a_n} (1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w}))^{1-a_n}.$$

Replacing in (2) and using the properties of the logarithm yields the desired result.

(b) We have that

$$\begin{aligned}
\mathbf{g} &= \nabla \ell(\mathcal{D}; \mathbf{w}) \\
&= \sum_{n=1}^N -a_n \frac{-\mathbf{x}_n e^{-\mathbf{w}^\top \mathbf{x}_n}}{1 + e^{-\mathbf{w}^\top \mathbf{x}_n}} + (1 - a_n) \left( -\mathbf{x}_n - \mathbf{x}_n \frac{-\mathbf{x}_n e^{-\mathbf{w}^\top \mathbf{x}_n}}{1 + e^{-\mathbf{w}^\top \mathbf{x}_n}} \right) \\
&= \sum_{n=1}^N a_n \mathbf{x}_n (1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w})) + (1 - a_n) \mathbf{x}_n \pi(1 \mid \mathbf{x}_n; \mathbf{w}) \\
&= \sum_{n=1}^N \mathbf{x}_n (a_n - \pi(1 \mid \mathbf{x}_n; \mathbf{w})).
\end{aligned}$$

(c) We have that

$$\begin{aligned}
\mathbf{H} &= \nabla^2 \ell(\mathcal{D}; \mathbf{w}) \\
&= \sum_{n=1}^N \mathbf{x}_n \frac{-\mathbf{x}_n^\top e^{-\mathbf{w}^\top \mathbf{x}_n}}{(1 + e^{-\mathbf{w}^\top \mathbf{x}_n})^2} \\
&= - \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \pi(1 \mid \mathbf{x}_n; \mathbf{w})(1 - \pi(1 \mid \mathbf{x}_n; \mathbf{w})).
\end{aligned}$$