Machine Learning I

Lucas Fabian Naumann, David Stein, Bjoern Andres

Machine Learning for Computer Vision TU Dresden



https://mlcv.cs.tu-dresden.de/courses/25-winter/ml1/

Winter Term 2025/2026

Contents. This part of the course is about the supervised learning of linear functions, more specifically, about logistic regression.

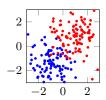
- We introduce the problem by defining labeled data, a family of functions and a probability measure whose maximization motivates a regularizer and a loss function.
- ► We show: This supervised learning problem is convex. It can be solved, e.g., by the steepest descent algorithm.

We consider labeled data with **real features**. More specifically, we consider some finite, non-empty set V, called the set of features, and labeled data T=(S,X,x,y) such that $X=\mathbb{R}^V$. Hence:

$$x \colon S \to \mathbb{R}^V \tag{1}$$

$$y \colon S \to \{0, 1\} \tag{2}$$

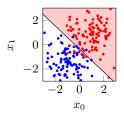
Example.



We consider **linear functions**. More specifically, we consider $\Theta=\mathbb{R}^V$ and $f:\Theta\to\mathbb{R}^X$ such that

$$\forall \theta \in \Theta \ \forall \hat{x} \in X \colon \quad f_{\theta}(\hat{x}) = \langle \theta, \hat{x} \rangle = \sum_{v \in V} \theta_v \, \hat{x}_v \tag{3}$$

Example.

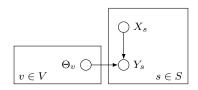


We introduce a probabilistic model:

- ▶ For any sample $s \in S$, let X_s be a random variable whose value is a vector $x_s \in \mathbb{R}^V$, the **feature vector** of s
- ▶ For any sample $s \in S$, let Y_s be a random variable whose value is a binary number $y_s \in \{0,1\}$, the **label** of s
- For any $v \in V$, let Θ_v be a random variable whose value is a real number $\theta_v \in \mathbb{R}$, a parameter of the linear function we seek to learn

We assume that the joint probability factorizes according to:

$$P(X, Y, \Theta) = \prod_{s \in S} (P(Y_s \mid X_s, \Theta)P(X_s)) \prod_{v \in V} P(\Theta_v)$$
 (4)



We attempt to learn parameters by maximizing the conditional probability

$$P(\Theta \mid X, Y) = \frac{P(X, Y, \Theta)}{P(X, Y)}$$

$$= \frac{P(Y \mid X, \Theta) P(X) P(\Theta)}{P(X, Y)}$$

$$\propto P(Y \mid X, \Theta) P(\Theta)$$

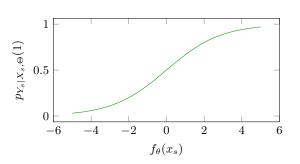
$$= \prod_{s \in S} P(Y_s \mid X_s, \Theta) \prod_{v \in V} P(\Theta_v) .$$

We attempt to infer labels by maximizing the conditional probability

$$P(Y \mid X, \Theta) = \prod_{s \in S} P(Y_s \mid X_s, \Theta) .$$

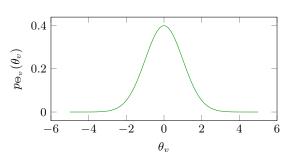
► Sigmoid distribution

$$\forall s \in S: \qquad p_{Y_s|X_s,\Theta}(1) = \frac{1}{1 + 2^{-f_{\theta}(x_s)}}$$
 (5)



▶ Normal distribution with $\sigma \in \mathbb{R}^+$:

$$\forall v \in V: \qquad p_{\Theta_v}(\theta_v) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\theta_v^2/2\sigma^2} \tag{5}$$



Lemma. Estimating maximally probable parameters θ , given attributes x and labels y, i.e.,

$$\underset{\theta \in \mathbb{R}^m}{\operatorname{argmax}} \quad p_{\Theta|X,Y}(\theta, x, y)$$

is equivalent of the supervised learning problem

$$\min_{\theta \in \Theta} \quad \lambda R(\theta) + \sum_{s \in S} L(f_{\theta}(x_s), y_s) \tag{6}$$

with L, R and λ such that

$$\forall r \in \mathbb{R} \ \forall \hat{y} \in \{0, 1\} \colon \quad L(r, \hat{y}) = -\hat{y}r + \log(1 + 2^r) \tag{7}$$

$$\forall \theta \in \Theta \colon \qquad R(\theta) = \|\theta\|_2^2$$
 (8)

$$\lambda = \frac{\log e}{2\sigma^2} \ . \tag{9}$$

It is called the l_2 -regularized **logistic regression problem** with respect to $x,\ y$ and $\sigma.$

Proof. Firstly,

$$\underset{\theta \in \mathbb{R}^{m}}{\operatorname{argmax}} \quad p_{\Theta|X,Y}(\theta, x, y) \\
= \underset{\theta \in \mathbb{R}^{m}}{\operatorname{argmax}} \quad \prod_{s \in S} p_{Y_{s}|X_{s},\Theta}(y_{s}, x_{s}, \theta) \prod_{v \in V} p_{\Theta_{v}}(\theta_{v}) \\
= \underset{\theta \in \mathbb{R}^{m}}{\operatorname{argmax}} \quad \sum_{s \in S} \log p_{Y_{s}|X_{s},\Theta}(y_{s}, x_{s}, \theta) + \sum_{v \in V} \log p_{\Theta_{v}}(\theta_{v}) \tag{10}$$

Secondly,

$$\log p_{Y_s|X_s,\Theta}(y_s,x_s,\theta)$$

$$= y_s \log p_{Y_s|X_s,\Theta}(1,x_s,\theta) + (1-y_s) \log p_{Y_s|X_s,\Theta}(0,x_s,\theta)$$

$$= y_s \log \frac{p_{Y_s|X_s,\Theta}(1,x_s,\theta)}{p_{Y_s|X_s,\Theta}(0,x_s,\theta)} + \log p_{Y_s|X_s,\Theta}(0,x_s,\theta)$$
(11)

Thus, with (5) and (4):

$$\underset{\theta \in \mathbb{R}^m}{\operatorname{argmin}} \quad \sum_{s \in S} \left(-y_s \langle \theta, x_s \rangle + \log \left(1 + 2^{\langle \theta, x_s \rangle} \right) \right) + \frac{\log e}{2\sigma^2} \|\theta\|_2^2$$
 (12)

Lemma. The objective function

$$\varphi(\theta) = \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_{\theta}(x_s), y_s)$$
 (13)

of the l_2 -regularized logistic regression problem is convex.

Proof. Exercise!

The l_2 -regularized logistic regression problem can be solved, e.g., by the steepest descent algorithm with a tolerance parameter $\epsilon \in \mathbb{R}_0^+$:

Algorithm. Steepest descent with line search

```
\begin{array}{l} \theta := 0 \\ \text{repeat} \\ d := \nabla \varphi(\theta) \\ \eta := \mathop{\mathrm{argmin}}_{\eta' \in \mathbb{R}} \varphi(\theta - \eta' d) \\ \theta := \theta - \eta d \\ \text{if } \|d\| < \epsilon \\ \text{return } \theta \end{array} \tag{line search}
```

Lemma: Estimating maximally probable labels y, given attributes x^\prime and parameters θ , i.e.,

$$\underset{y \in \{0,1\}^S}{\operatorname{argmax}} \quad p_{Y|X,\Theta}(y, x', \theta) \tag{14}$$

is equivalent to the inference problem

$$\min_{y' \in \{0,1\}^S} \sum_{s \in S} L(f_{\theta}(x_s), y'_s) . \tag{15}$$

It has the solution

$$\forall s \in S' : \quad y_s = \begin{cases} 1 & \text{if } f_{\theta}(x'_s) > 0 \\ 0 & \text{otherwise} \end{cases}$$
 (16)

Proof. Firstly,

Secondly,

$$\min_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left(-y_s f_{\theta}(x'_s) + \log\left(1 + 2^{f_{\theta}(x'_s)}\right) \right) = \sum_{s \in S'} \max_{y_s \in \{0,1\}} y_s f_{\theta}(x'_s) .$$

Summary.

- ► The l₂-regularized logistic regression problem is a supervised learning problem wrt. the family of linear functions.
- ▶ It can be derived from a statistical model with the sigmoid distribution as the likelihood as the normal distribution as the prior.
- ▶ It is a convex optimization problem that can be solved, e.g., by the steepest descent algorithm.