

# Machine Learning I

Bjoern Andres, Shengxian Zhao

Machine Learning for Computer Vision  
TU Dresden



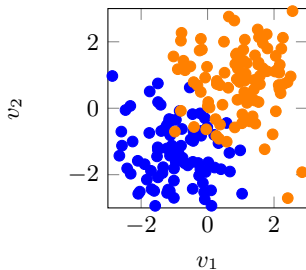
Winter Term 2022/2023

## Deciding with Linear Functions

**Contents.** This part of the course is about a special case of supervised learning: the supervised learning of linear functions by **logistic regression**.

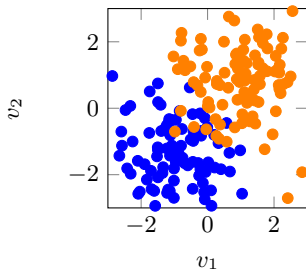
- ▶ We state the problem by defining labeled data, the family of functions and a **probability distribution** whose maximization motivates a regularizer and a loss function
- ▶ We show: This supervised learning problem is convex and can thus be solved by means of the **steepest descent algorithm**.

## Deciding with Linear Functions



We consider **real attributes**. More specifically, we consider some finite set  $V \neq \emptyset$  and labeled data  $T = (S, X, x, y)$  with  $X = \mathbb{R}^V$ .

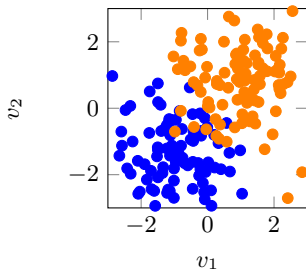
## Deciding with Linear Functions



We consider **real attributes**. More specifically, we consider some finite set  $V \neq \emptyset$  and labeled data  $T = (S, X, x, y)$  with  $X = \mathbb{R}^V$ .

Hence,  $x: S \rightarrow \mathbb{R}^V$  and  $y: S \rightarrow \{0, 1\}$ .

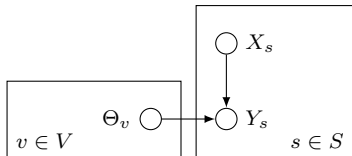
## Deciding with Linear Functions



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

$$\forall \theta \in \Theta \quad \forall \hat{x} \in X: \quad f_{\theta}(\hat{x}) = \langle \theta, \hat{x} \rangle = \sum_{v \in V} \theta_v \hat{x}_v \quad (1)$$

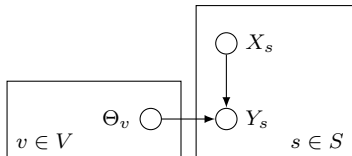
## Deciding with Linear Functions



### *Random Variables*

- 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 **attribute vector** of  $s$

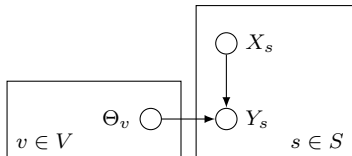
## Deciding with Linear Functions



### *Random Variables*

- ▶ 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 **attribute 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  $s$

## Deciding with Linear Functions

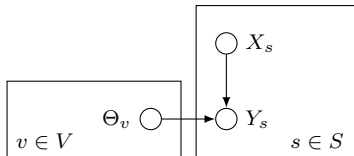


### *Random Variables*

- ▶ 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 **attribute 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  $\Theta_v$  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



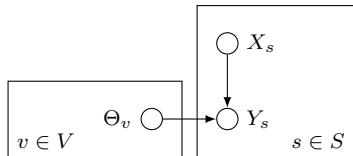
## Deciding with Linear Functions



### *Factorization*

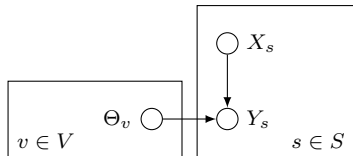
$$P(X, Y, \Theta) = \prod_{s \in S} (P(Y_s | X_s, \Theta) P(X_s)) \prod_{v \in V} P(\Theta_v) \quad (2)$$

## Deciding with Linear Functions



*Factorization*

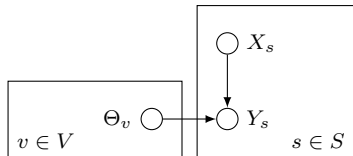
## Deciding with Linear Functions



*Factorization*

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

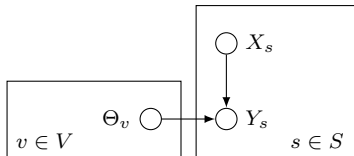
## Deciding with Linear Functions



*Factorization*

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

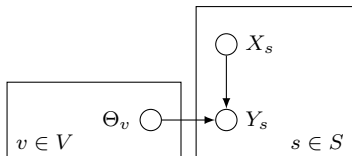
## Deciding with Linear Functions



### *Factorization*

$$\begin{aligned} P(\Theta \mid X, Y) &= \frac{P(X, Y, \Theta)}{P(X, Y)} \\ &= \frac{P(Y \mid X, \Theta) P(X) P(\Theta)}{P(X, Y)} \end{aligned}$$

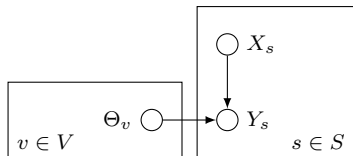
## Deciding with Linear Functions



### *Factorization*

$$\begin{aligned} P(\Theta \mid X, Y) &= \frac{P(X, Y, \Theta)}{P(X, Y)} \\ &= \frac{P(Y \mid X, \Theta) P(X) P(\Theta)}{P(X, Y)} \\ &\propto P(Y \mid X, \Theta) P(\Theta) \end{aligned}$$

## Deciding with Linear Functions



### *Factorization*

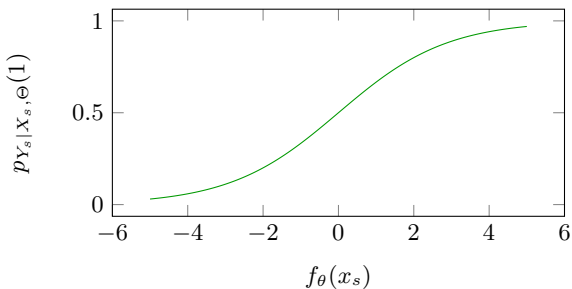
$$\begin{aligned} P(\Theta \mid X, Y) &= \frac{P(X, Y, \Theta)}{P(X, Y)} \\ &= \frac{P(Y \mid X, \Theta) P(X) P(\Theta)}{P(X, Y)} \\ &\propto P(Y \mid X, \Theta) P(\Theta) \\ &= \prod_{s \in S} P(Y_s \mid X_s, \Theta) \prod_{v \in V} P(\Theta_v) \end{aligned}$$

# Deciding with Linear Functions

## *Distributions*

### ► Logistic distribution

$$\forall s \in S: \quad p_{Y_s|X_s, \Theta}(1) = \frac{1}{1 + 2^{-f_{\theta}(x_s)}} \quad (3)$$



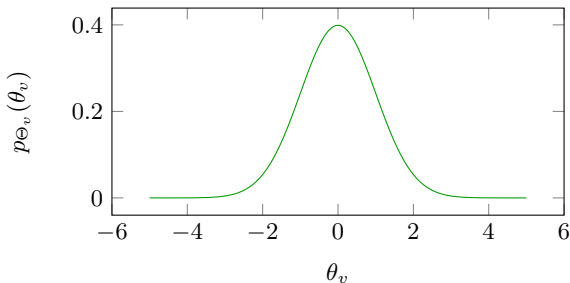


## Deciding with Linear Functions

### *Distributions*

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

$$\forall v \in V : \quad p_{\Theta_v}(\theta_v) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\theta_v^2/2\sigma^2} \quad (3)$$



## Deciding with Linear Functions

**Lemma.** Estimating maximally probable parameters  $\theta$ , given attributes  $x$  and labels  $y$ , i.e.,

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

is equivalent to the supervised learning problem

$$\min_{\theta \in \Theta} \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_{\theta}(x_s), y_s) \quad (4)$$

with  $L$ ,  $R$  and  $\lambda$  such that

$$\forall r \in \mathbb{R} \quad \forall \hat{y} \in \{0, 1\}: \quad L(r, \hat{y}) = -\hat{y}r + \log(1 + 2^r) \quad (5)$$

$$\forall \theta \in \Theta: \quad R(\theta) = \|\theta\|_2^2 \quad (6)$$

$$\lambda = \frac{\log e}{2\sigma^2} . \quad (7)$$

## Deciding with Linear Functions

**Lemma.** Estimating maximally probable parameters  $\theta$ , given attributes  $x$  and labels  $y$ , i.e.,

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

is equivalent to the supervised learning problem

$$\min_{\theta \in \Theta} \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_{\theta}(x_s), y_s) \quad (4)$$

with  $L$ ,  $R$  and  $\lambda$  such that

$$\forall r \in \mathbb{R} \quad \forall \hat{y} \in \{0, 1\}: \quad L(r, \hat{y}) = -\hat{y}r + \log(1 + 2^r) \quad (5)$$

$$\forall \theta \in \Theta: \quad R(\theta) = \|\theta\|_2^2 \quad (6)$$

$$\lambda = \frac{\log e}{2\sigma^2} . \quad (7)$$

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

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{\theta \in \mathbb{R}^m} p_{\Theta|X,Y}(\theta, x, y) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \prod_{s \in S} p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) \prod_{v \in V} p_{\Theta_v}(\theta_v) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \sum_{s \in S} \log p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) + \sum_{v \in V} \log p_{\Theta_v}(\theta_v) \end{aligned} \quad (8)$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{\theta \in \mathbb{R}^m} p_{\Theta|X,Y}(\theta, x, y) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \prod_{s \in S} p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) \prod_{v \in V} p_{\Theta_v}(\theta_v) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \sum_{s \in S} \log p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) + \sum_{v \in V} \log p_{\Theta_v}(\theta_v) \end{aligned} \quad (8)$$

Secondly,

$$\begin{aligned} & \log p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) \\ &= y_s \log p_{Y_s|X_s, \Theta}(1, x_s, \theta) + (1 - y_s) \log p_{Y_s|X_s, \Theta}(0, x_s, \theta) \\ &= y_s \log \frac{p_{Y_s|X_s, \Theta}(1, x_s, \theta)}{p_{Y_s|X_s, \Theta}(0, x_s, \theta)} + \log p_{Y_s|X_s, \Theta}(0, x_s, \theta) \end{aligned} \quad (9)$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{\theta \in \mathbb{R}^m} p_{\Theta|X,Y}(\theta, x, y) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \prod_{s \in S} p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) \prod_{v \in V} p_{\Theta_v}(\theta_v) \\ &= \operatorname{argmax}_{\theta \in \mathbb{R}^m} \sum_{s \in S} \log p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) + \sum_{v \in V} \log p_{\Theta_v}(\theta_v) \end{aligned} \quad (8)$$

Secondly,

$$\begin{aligned} & \log p_{Y_s|X_s, \Theta}(y_s, x_s, \theta) \\ &= y_s \log p_{Y_s|X_s, \Theta}(1, x_s, \theta) + (1 - y_s) \log p_{Y_s|X_s, \Theta}(0, x_s, \theta) \\ &= y_s \log \frac{p_{Y_s|X_s, \Theta}(1, x_s, \theta)}{p_{Y_s|X_s, \Theta}(0, x_s, \theta)} + \log p_{Y_s|X_s, \Theta}(0, x_s, \theta) \end{aligned} \quad (9)$$

Thus, with (3) and (4):

$$\operatorname{argmin}_{\theta \in \mathbb{R}^m} \sum_{s \in S} \left( -y_s \langle \theta, x_s \rangle + \log \left( 1 + 2^{\langle \theta, x_s \rangle} \right) \right) + \frac{\log e}{2\sigma^2} \|\theta\|_2^2 \quad (10)$$

## Deciding with Linear Functions

**Lemma.** The objective function

$$\varphi(\theta) = \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_{\theta}(x_s), y_s) \quad (11)$$

of the  $l_2$ -regularized logistic regression problem is convex.

## Deciding with Linear Functions

**Lemma.** The objective function

$$\varphi(\theta) = \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_\theta(x_s), y_s) \quad (11)$$

of the  $l_2$ -regularized logistic regression problem is convex.

*Proof.* Exercise!



## Deciding with Linear Functions

**Lemma.** The objective function

$$\varphi(\theta) = \lambda R(\theta) + \frac{1}{|S|} \sum_{s \in S} L(f_{\theta}(x_s), y_s) \quad (11)$$

of the  $l_2$ -regularized logistic regression problem is convex.

*Proof.* Exercise!

The problem can be solved by the steepest descent algorithm with a tolerance parameter  $\epsilon \in \mathbb{R}_0^+$ :

---

```
 $\theta := 0$ 
repeat
   $d := \nabla \varphi(\theta)$ 
   $\eta := \operatorname{argmin}_{\eta' \in \mathbb{R}} \varphi(\theta - \eta' d)$  (line search)
   $\theta := \theta - \eta d$ 
  if  $\|d\| < \epsilon$ 
    return  $\theta$ 
```

---

## Deciding with Linear Functions

**Lemma:** Estimating maximally probable labels  $y$ , given attributes  $x'$  and parameters  $\theta$ , i.e.,

$$\operatorname{argmax}_{y \in \{0,1\}^S} p_{Y|X,\Theta}(y, x', \theta) \quad (12)$$

is equivalent to the inference problem

$$\min_{y' \in \{0,1\}^S} \sum_{s \in S} L(f_\theta(x_s), y'_s) . \quad (13)$$

It has the solution

$$\forall s \in S' : y_s = \begin{cases} 1 & \text{if } f_\theta(x'_s) > 0 \\ 0 & \text{otherwise} \end{cases} . \quad (14)$$

## Deciding with Linear Functions

*Proof.* Firstly,

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

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \end{aligned}$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \end{aligned}$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x'_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x'_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x'_s, \theta) \right) \end{aligned}$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x'_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x'_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x'_s, \theta) \right) \\ &= \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( -y_s f_\theta(x'_s) + \log \left( 1 + 2^{f_\theta(x'_s)} \right) \right) \end{aligned}$$

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ = & \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x'_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x'_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x'_s, \theta) \right) \\ = & \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( -y_s f_\theta(x'_s) + \log \left( 1 + 2^{f_\theta(x'_s)} \right) \right) \\ = & \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} L(f_\theta(x'_s), y_s) . \end{aligned}$$



## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x'_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x'_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x'_s, \theta) \right) \\ &= \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( -y_s f_\theta(x'_s) + \log \left( 1 + 2^{f_\theta(x'_s)} \right) \right) \\ &= \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} L(f_\theta(x'_s), y_s) . \end{aligned}$$

Secondly,

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

## Deciding with Linear Functions

*Proof.* Firstly,

$$\begin{aligned} & \operatorname{argmax}_{y \in \{0,1\}^{S'}} p_{Y|X,\Theta}(y, x', \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \prod_{s \in S'} p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \log p_{Y_s|X_s,\Theta}(y_s, x'_s, \theta) \\ &= \operatorname{argmax}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( y_s \log \frac{p_{Y_s|X_s,\Theta}(1, x'_s, \theta)}{p_{Y_s|X_s,\Theta}(0, x'_s, \theta)} + \log p_{Y_s|X_s,\Theta}(0, x'_s, \theta) \right) \\ &= \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} \left( -y_s f_\theta(x'_s) + \log \left( 1 + 2^{f_\theta(x'_s)} \right) \right) \\ &= \operatorname{argmin}_{y \in \{0,1\}^{S'}} \sum_{s \in S'} L(f_\theta(x'_s), y_s) . \end{aligned}$$

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  $l_2$ -regularized logistic regression problem is a supervised learning problem w.r.t. the family of linear functions.
- ▶ It is motivated by a Bayesian statistical model with the logistic 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.