

# Chapter ML:III

## III. Linear Models

- ❑ Logistic Regression
- ❑ Loss Computation in Detail
- ❑ Overfitting
- ❑ Regularization
- ❑ Gradient Descent in Detail

# Logistic Regression

## Binary Classification Problems

Setting:

- $X$  is a multiset of feature vectors from an inner product space  $\mathbf{X}$ ,  $\mathbf{X} \subseteq \mathbf{R}^p$ .
- $C = \{0, 1\}$  is a set of two classes. Similarly:  $\{-1, 1\}$ ,  $\{\ominus, \oplus\}$ ,  $\{\text{no}, \text{yes}\}$ , etc.
- $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times C$  is a multiset of examples.

Learning task:

- Fit  $D$  using a logistic function  $y()$ .

# Logistic Regression

## Binary Classification Problems

Setting:

- $X$  is a multiset of feature vectors from an inner product space  $\mathbf{X}$ ,  $\mathbf{X} \subseteq \mathbf{R}^p$ .
- $C = \{0, 1\}$  is a set of two classes. Similarly:  $\{-1, 1\}$ ,  $\{\ominus, \oplus\}$ ,  $\{\text{no}, \text{yes}\}$ , etc.
- $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times C$  is a multiset of examples.

Learning task:

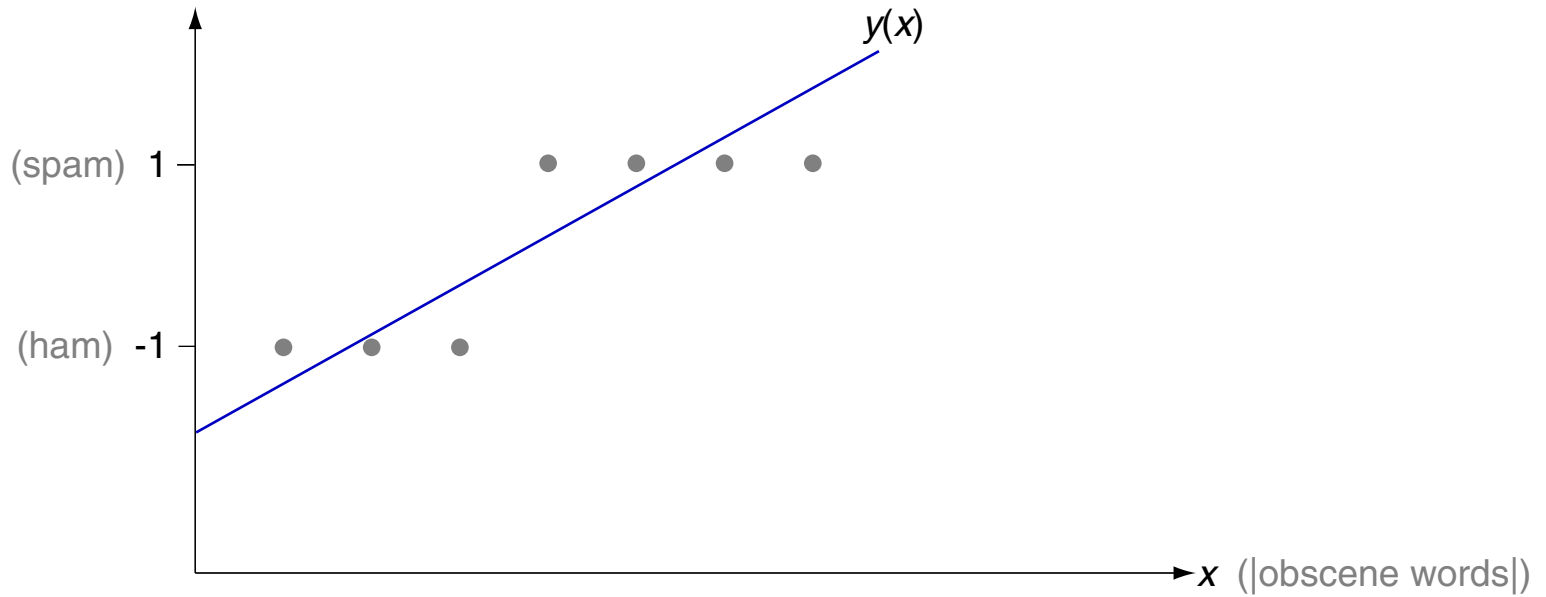
- Fit  $D$  using a logistic function  $y()$ .

Examples for binary classification problems:

- Is an email spam or ham?
- Is a patient infected or healthy?
- Is a bank customer creditworthy or not?

# Logistic Regression

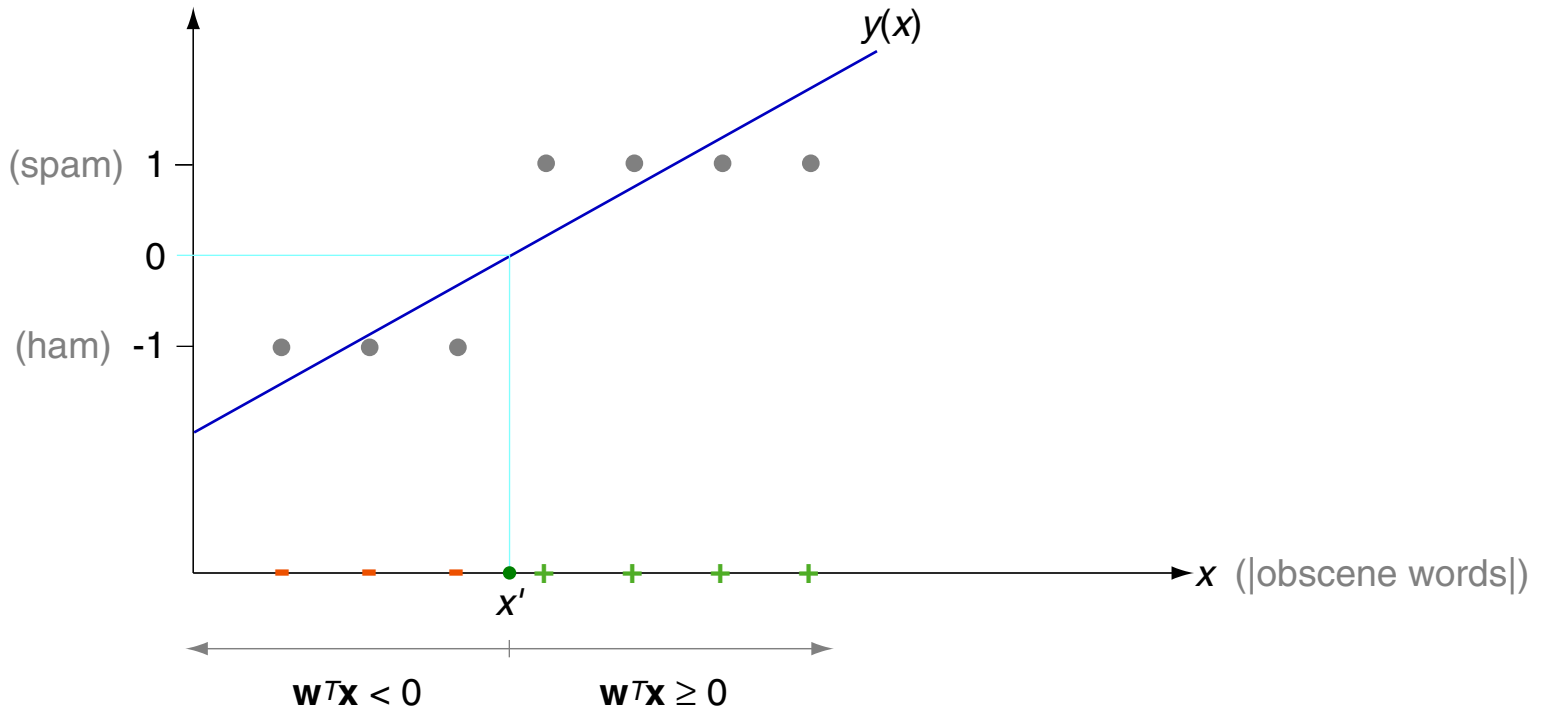
## Linear Regression



□ Linear regression:  $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$

# Logistic Regression

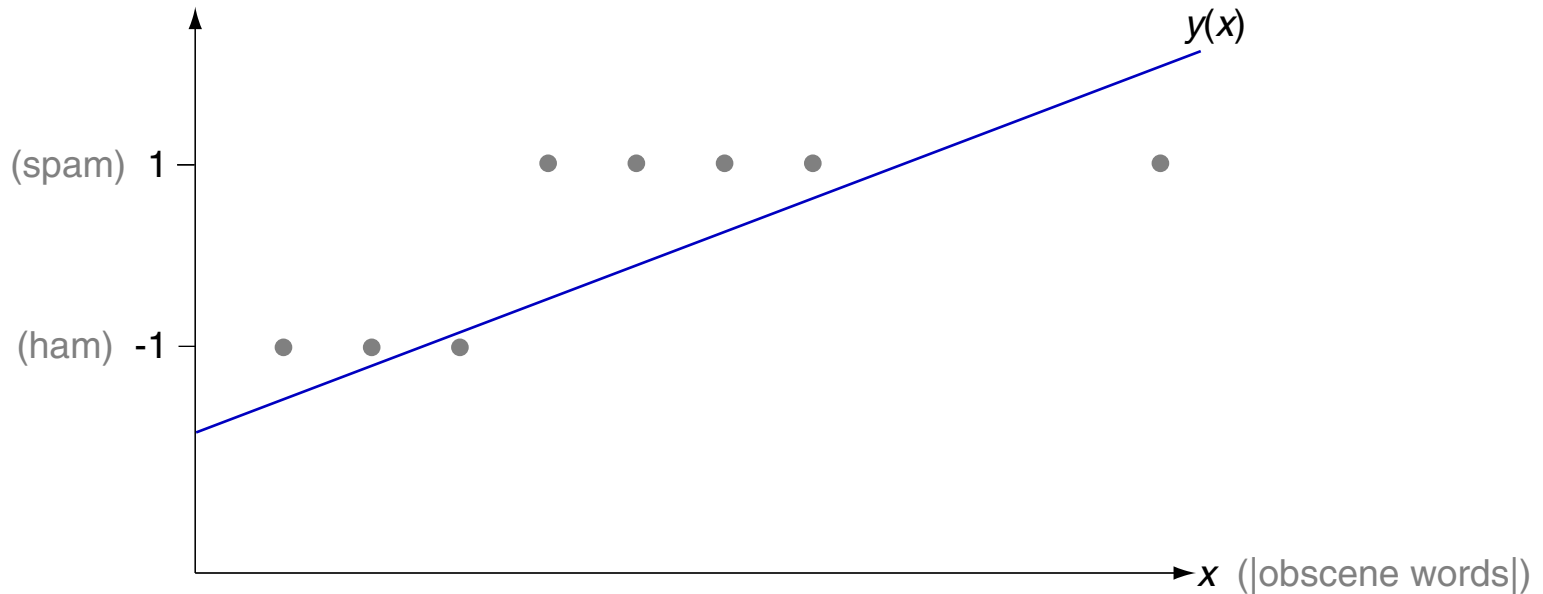
## Linear Regression (continued)



- Linear regression:  $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$
- Classification: Predict  $\begin{cases} \text{"spam"}, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ \text{"ham"}, & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$

# Logistic Regression

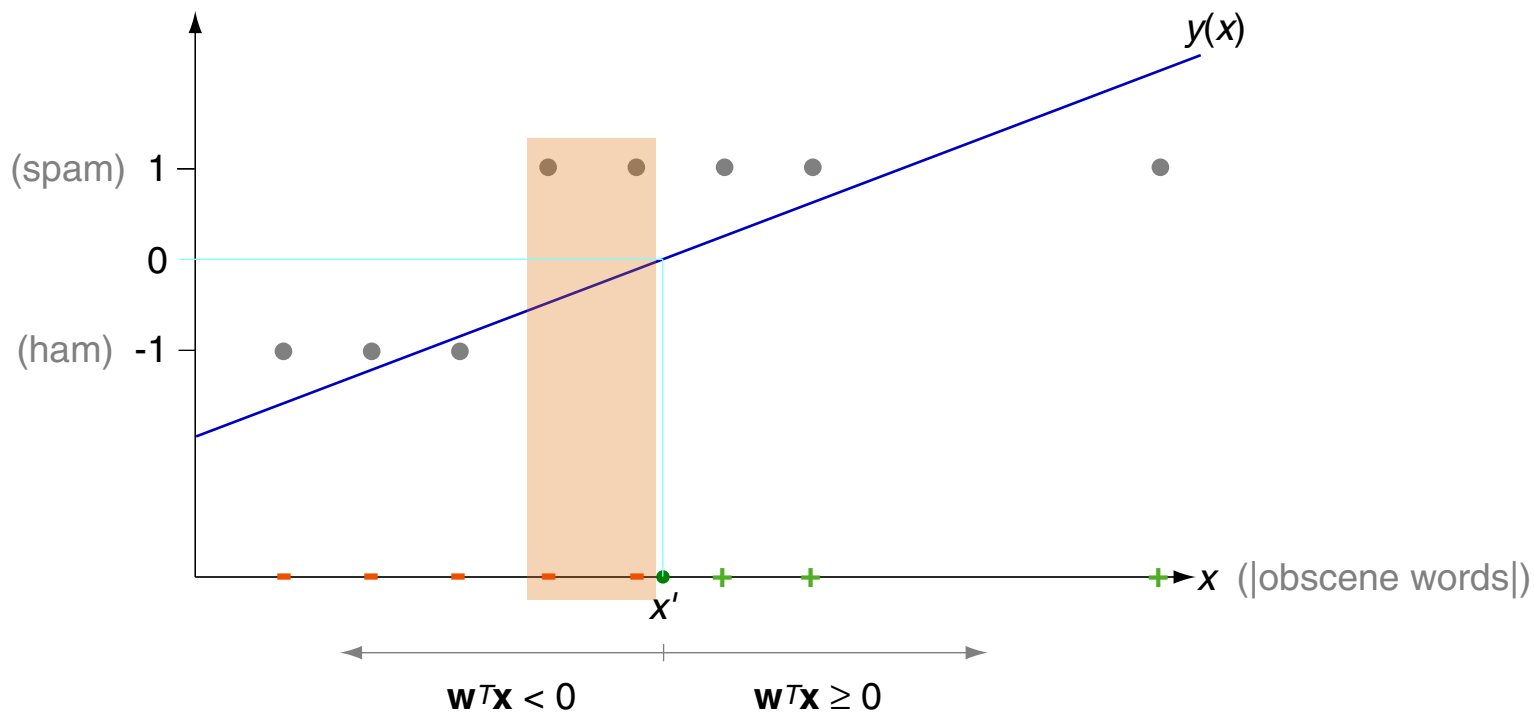
## Linear Regression (continued)



- Linear regression:  $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$
- Classification: Predict  $\begin{cases} \text{"spam"}, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ \text{"ham"}, & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$

# Logistic Regression

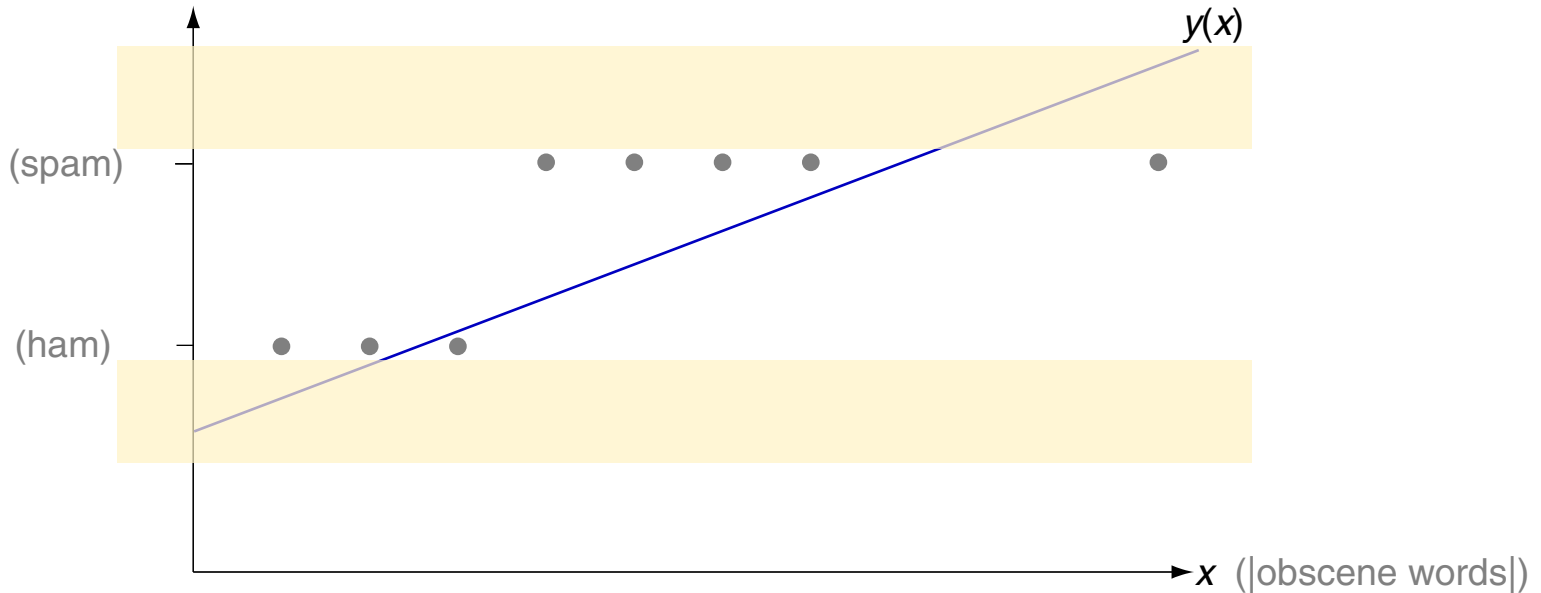
## Linear Regression (continued)



- Linear regression:  $y(x) \stackrel{(*)}{=} w^T x$
- Classification: Predict  $\begin{cases} \text{"spam"}, & \text{if } w^T x \geq 0 \\ \text{"ham"}, & \text{if } w^T x < 0 \end{cases}$

# Logistic Regression

## Linear Regression (continued)



Restrict the range of  $y(\mathbf{x})$  to reflect the two-class classification semantics:

$$-1 \leq y(\mathbf{x}) \leq 1 \quad \text{or} \quad 0 \leq y(\mathbf{x}) \leq 1$$

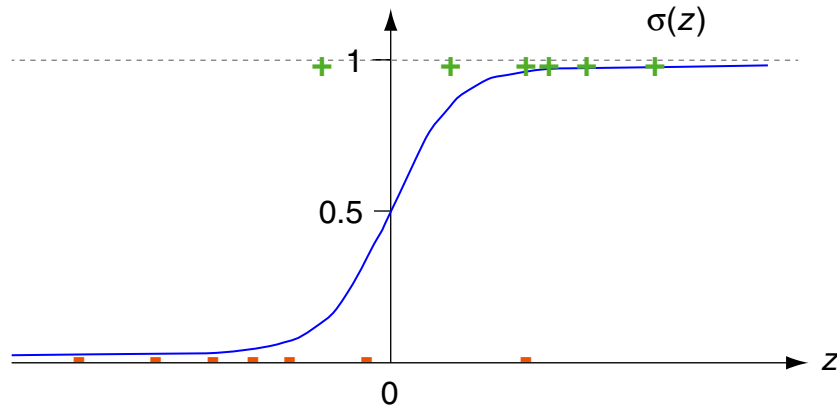


## Remarks:

- (★) Recap. We consider the feature vector  $\mathbf{x}$  in its extended form when used as operand in a scalar product with the weight vector,  $\mathbf{w}^T \mathbf{x}$ , and consequently, when noted as argument of the model function,  $y(\mathbf{x})$ . I.e.,  $\mathbf{x} = (1, x_1, \dots, x_p)^T \in \mathbb{R}^{p+1}$ , and  $x_0 = 1$ .

# Logistic Regression

## Sigmoid (Logistic) Function

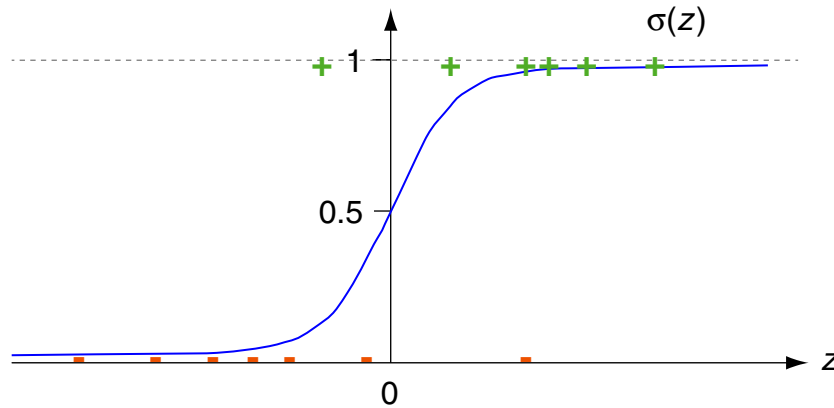


Sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

# Logistic Regression

## Sigmoid (Logistic) Function (continued)



Linear regression

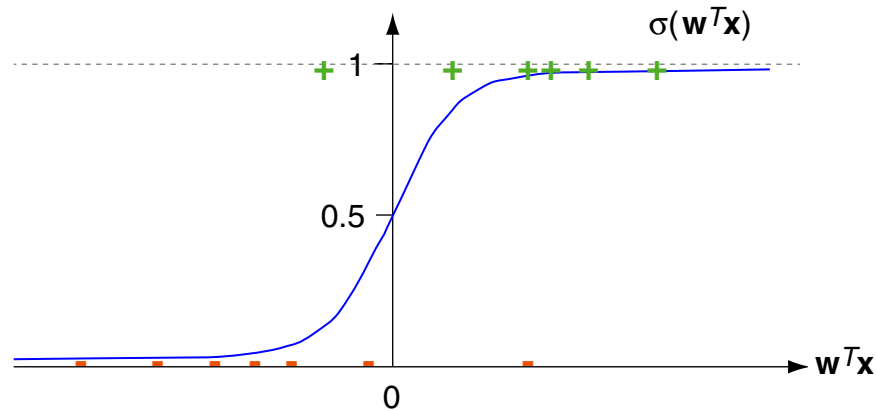
$$\mathbf{w}^T \mathbf{x}$$

Sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

# Logistic Regression

## Sigmoid (Logistic) Function (continued)



Linear regression

$$\mathbf{w}^T \mathbf{x}$$

Sigmoid function

$$\circ \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

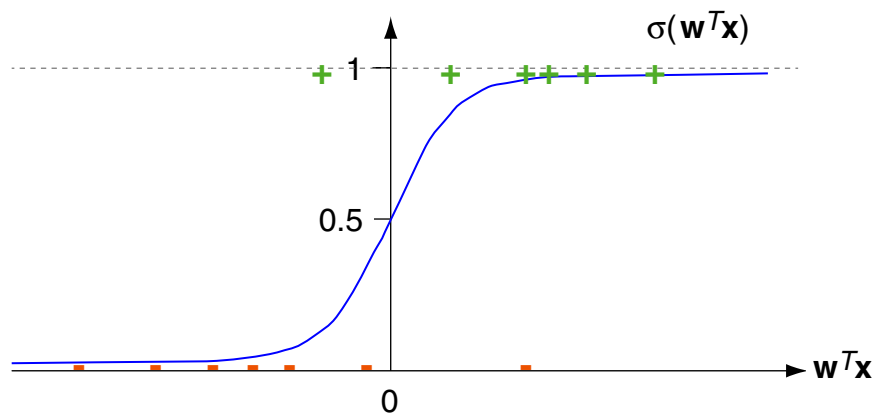
$\leadsto$

Logistic model function

$$y(\mathbf{x}) \equiv \sigma(\mathbf{w}^T \mathbf{x}) \stackrel{(*)}{=} \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

# Logistic Regression

## Sigmoid (Logistic) Function (continued)



Linear regression

$$\mathbf{w}^T \mathbf{x}$$

Sigmoid function

$$\circ \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

$\leadsto$

Logistic model function

$$y(\mathbf{x}) \equiv \sigma(\mathbf{w}^T \mathbf{x}) \stackrel{(*)}{=} \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

$$y(\mathbf{x}) : \mathbf{R}^{p+1} \rightarrow (0; 1)$$

# Logistic Regression

## Interpretation of the Logistic Model Function

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w}),$  “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w}),$  “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

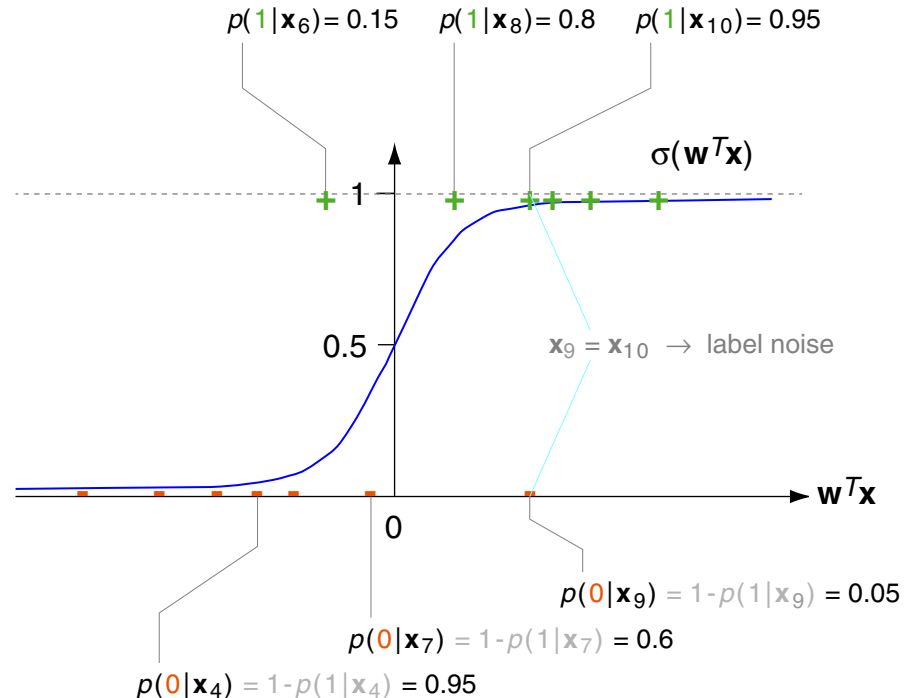
# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$ :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

The »+« and »-« are examples from  $D$ , projected on the  $z$ -axis, along with the probabilities as specified by  $\sigma(\mathbf{w}^T \mathbf{x})$ .



# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

Example (email spam classification) :

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ |\text{obscene words}| \end{pmatrix}, \quad \mathbf{x}_1 = \begin{pmatrix} 1 \\ 5 \end{pmatrix} \quad \text{and} \quad y(\mathbf{x}_1) = 0.67$$

$\leadsto$  67% chance that this email is spam.



# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated **probability for the event  $C=1$**  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

Example (email spam classification) :

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ |\text{obscene words}| \end{pmatrix}, \quad \mathbf{x}_1 = \begin{pmatrix} 1 \\ 5 \end{pmatrix} \quad \text{and} \quad y(\mathbf{x}_1) = 0.67$$

$\leadsto$  **67% chance that this email is spam.**

# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

Example (email spam classification) :

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ |\text{obscene words}| \end{pmatrix}, \quad \mathbf{x}_1 = \begin{pmatrix} 1 \\ 5 \end{pmatrix} \quad \text{and} \quad y(\mathbf{x}_1) = 0.67$$

$\leadsto$  67% chance that this email is spam.

Classification: Predict  $\begin{cases} 1, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \\ 0, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) < 0.5 \end{cases}$

# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

Example (email spam classification) :

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ |\text{obscene words}| \end{pmatrix}, \quad \mathbf{x}_1 = \begin{pmatrix} 1 \\ 5 \end{pmatrix} \quad \text{and} \quad y(\mathbf{x}_1) = 0.67$$

$\leadsto$  67% chance that this email is spam.

Classification: Predict  $\begin{cases} 1, & \text{if } \underline{\sigma(\mathbf{w}^T \mathbf{x})} \geq 0.5 & \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0 \\ 0, & \text{if } \underline{\sigma(\mathbf{w}^T \mathbf{x})} < 0.5 & \Leftrightarrow \mathbf{w}^T \mathbf{x} < 0 \end{cases}$

# Logistic Regression

## Interpretation of the Logistic Model Function (continued)

$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as the estimated probability for the event  $C=1$  :

- $y(\mathbf{x}) =: P(C=1 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(1 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=1$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”
- $1 - y(\mathbf{x}) =: P(C=0 \mid \mathbf{X}=\mathbf{x}; \mathbf{w}) = p(0 \mid \mathbf{x}; \mathbf{w})$ , “Probability for  $C=0$  given  $\mathbf{x}$ , parameterized by  $\mathbf{w}$ .”

Estimate optimum  $\mathbf{w}$  by maximizing the probability  $p(D; \mathbf{w})$  :

$$\mathbf{w}_{\text{ML}} = \underset{\mathbf{w} \in \mathbf{R}^{p+1}}{\operatorname{argmax}} p(D; \mathbf{w})$$

⋮ [derivation]

$$= \underset{\mathbf{w} \in \mathbf{R}^{p+1}}{\operatorname{argmin}} L_{\sigma}(\mathbf{w}) \quad \text{[RSS minimization]}$$

I.e., optimizing  $\mathbf{w}$  by maximizing  $p(D; \mathbf{w})$  is equivalent to optimizing  $\mathbf{w}$  by minimizing the logistic loss  $L_{\sigma}(\mathbf{w})$ .

## Remarks (probabilistic view to classification) :

- If  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  is interpreted as “probability for  $C=1$  given feature vector  $\mathbf{x}$ ”, then  $\mathbf{w}$  is the (unique) characterizing parameter vector of the hidden stochastic process that generates the observed data  $D$ .

Recap. As a consequence,  $\mathbf{w}$  is not the realization of a random variable—which would come along with a distribution—but an *exogenous parameter*, which is varied in order to find the maximum probability  $p(D; \mathbf{w})$  or the minimum loss  $L_\sigma(\mathbf{w})$ .

The fact that  $\mathbf{w}$  is an exogenous parameter and not a the realization of a random variable is reflected by the notation, which uses a  $\gg; \ll$  instead of a  $\gg | \ll$  in  $p()$ .

- The underlying probability space—which can be left implicit—looks as follows:

The sample space  $\Omega$  corresponds to a set  $O$  of real-world objects,  $P$  is a probability measure defined on  $\mathcal{P}(\Omega)$ . The classification task (experiment) suggests two types of random variables,  $\mathbf{X} : \Omega \rightarrow \mathbf{X}$  and  $C : \Omega \rightarrow \{0, 1\}$ .

See section Evaluating Effectiveness of part Machine Learning Basics for an illustration of the probabilistic view to classification, and section Probability Basics of part Bayesian Learning for a recap of concepts from probability theory.

- $\mathbf{X}$  and  $C$  denote (multivariate) random variables with ranges  $\mathbf{X}$  and  $C$  respectively.

$\mathbf{X}$  corresponds to a model formation function  $\alpha$  that returns for a real-world object  $o \in O$  its feature vector  $\mathbf{x}$ ,  $\mathbf{x} = \alpha(o)$ , and  $C$  corresponds to an ideal classifier  $\gamma$  that returns its class  $c$ ,  $c = \gamma(o)$ .

## Remarks (probabilistic view to classification) : (continued)

- ❑ Interpreting  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  as probability for the event  $C=1$  means that  $P(C=1 \mid \mathbf{X}=\mathbf{x})$  is defined as  $y(\mathbf{x})$ .
- ❑ The interpretation of  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$  as probability for the event  $C=1$  is not a mathematical consequence but a decision of the modeler. This decision is based on the advantageous properties of the sigmoid function, on practical considerations, and on heuristic simplifications of the real world.

Consider the following two aspects where an interpretation of  $\sigma(\mathbf{w}^T \mathbf{x})$  is questionable:

1. The sigmoid function implies that with  $\mathbf{w}^T \mathbf{x} \rightarrow +\infty$  we get  $y(\mathbf{x}) \rightarrow 1$  or  $P(C=1) \rightarrow 1$ . Likewise, with  $\mathbf{w}^T \mathbf{x} \rightarrow -\infty$  we get  $P(C=1) \rightarrow 0$ . Though such a strict monotonicity appears self-evident, it need not necessarily correspond to the observed behavior in a real world experiment.
2. The sigmoid function implies a smooth, virtually linear transition from low probability values (around 0.1) to high probability values (around 0.9) as its argument  $\mathbf{w}^T \mathbf{x}$  increases. This link between the continuous growth of  $\mathbf{w}^T \mathbf{x}$  and the continuous growth of probability values  $P(C=1)$  presumes a proportional connection between cause (in the form of  $\mathbf{X}=\mathbf{x}$ ) and effect (in the form of  $C=1$ ). Again, such a relation appears sensible but may not necessarily model the real world.

Remarks (derivation of  $L_\sigma(\mathbf{w})$ ):

- The most probable (= optimum) hypothesis in the space  $H$  of possible hypotheses,  $h_{\text{ML}}$ , can be estimated with the maximum likelihood principle:  $h_{\text{ML}} = \underset{h \in H}{\operatorname{argmax}} p(D; h)$ .
- Applied to logistic regression:  $\mathbf{w}_{\text{ML}} = \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} p(D; \mathbf{w})$ , where

$$\begin{aligned}\underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} p(D; \mathbf{w}) &= \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, c) \in D} p(\mathbf{x}, c; \mathbf{w}) = \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, c) \in D} \left( p(c \mid \mathbf{x}; \mathbf{w}) \cdot p(\mathbf{x}) \right) \\&= \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, c) \in D} p(\mathbf{x}) \cdot \prod_{(\mathbf{x}, c) \in D} p(c \mid \mathbf{x}; \mathbf{w}) \stackrel{(1)}{=} \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, c) \in D} p(c \mid \mathbf{x}; \mathbf{w}) \\&= \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, 1) \in D} \sigma(\mathbf{w}^T \mathbf{x}) \cdot \prod_{(\mathbf{x}, 0) \in D} (1 - \sigma(\mathbf{w}^T \mathbf{x})) \\&= \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \prod_{(\mathbf{x}, 1) \in D} y(\mathbf{x}) \cdot \prod_{(\mathbf{x}, 0) \in D} (1 - y(\mathbf{x})) \\&\stackrel{(2)}{=} \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \log \prod_{(\mathbf{x}, 1) \in D} y(\mathbf{x}) + \log \prod_{(\mathbf{x}, 0) \in D} (1 - y(\mathbf{x})) \\&= \underset{\mathbf{w} \in \mathbb{R}^{p+1}}{\operatorname{argmax}} \sum_{(\mathbf{x}, 1) \in D} \log(y(\mathbf{x})) + \sum_{(\mathbf{x}, 0) \in D} \log(1 - y(\mathbf{x})) \\&= \hookrightarrow \text{p. 24}\end{aligned}$$

Remarks (derivation of  $L_\sigma(\mathbf{w})$ ): (continued)

$$\begin{aligned}
 &= \operatorname{argmax}_{\mathbf{w} \in \mathbf{R}^{p+1}} \sum_{(\mathbf{x}, c) \in D} \left( c \cdot \log(y(\mathbf{x})) + (1 - c) \cdot \log(1 - y(\mathbf{x})) \right) \\
 &\stackrel{(3)}{=} \operatorname{argmin}_{\mathbf{w} \in \mathbf{R}^{p+1}} - \sum_{(\mathbf{x}, c) \in D} \left( c \cdot \log(y(\mathbf{x})) + (1 - c) \cdot \log(1 - y(\mathbf{x})) \right) \\
 &\stackrel{(4)}{=} \operatorname{argmin}_{\mathbf{w} \in \mathbf{R}^{p+1}} \sum_{(\mathbf{x}, c) \in D} \overbrace{\left( -c \cdot \log(y(\mathbf{x})) - (1 - c) \cdot \log(1 - y(\mathbf{x})) \right)}^{l_\sigma :=} \\
 &=: \operatorname{argmin}_{\mathbf{w} \in \mathbf{R}^{p+1}} \sum_{(\mathbf{x}, c) \in D} l_\sigma(c, y(\mathbf{x})) = \operatorname{argmin}_{\mathbf{w} \in \mathbf{R}^{p+1}} L_\sigma(\mathbf{w})
 \end{aligned}$$

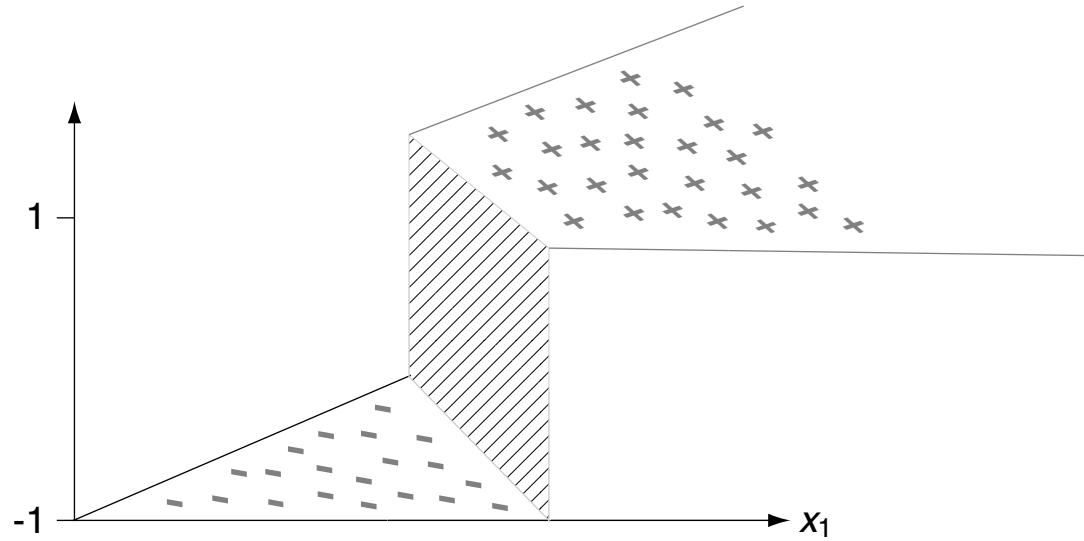
## □ Hints:

- (1)  $\prod_{(\mathbf{x}, c) \in D} p(\mathbf{x})$  is constant with respect to the variation of  $\mathbf{w}$ .
- (2)  $\operatorname{argmax}_x f(x) = \operatorname{argmax}_x g \circ f(x)$  (similarly for  $\operatorname{argmin}$ ) if  $g(z)$  is a strictly monotonically increasing function. Here,  $\log(z)$  is in the role of  $g(z)$ . Conversely, if  $g(z)$  is a strictly monotonically decreasing function, then  $\operatorname{argmax}_x f(x) = \operatorname{argmin}_x g \circ f(x)$ .
- (3) The maximization problem (the  $\operatorname{argmax}$ -expression) can be reformulated as minimization problem, i.e., as an  $\operatorname{argmin}$ -expression. See in this regard the second part of Hint (2).
- (4) The argument of the  $\operatorname{argmin}$ -expression quantifies  $L_\sigma(\mathbf{w})$ , the global logistic loss related to some  $\mathbf{w}$ , and, analogously, the pointwise logistic loss,  $l_\sigma(c, y(\mathbf{x}))$ .



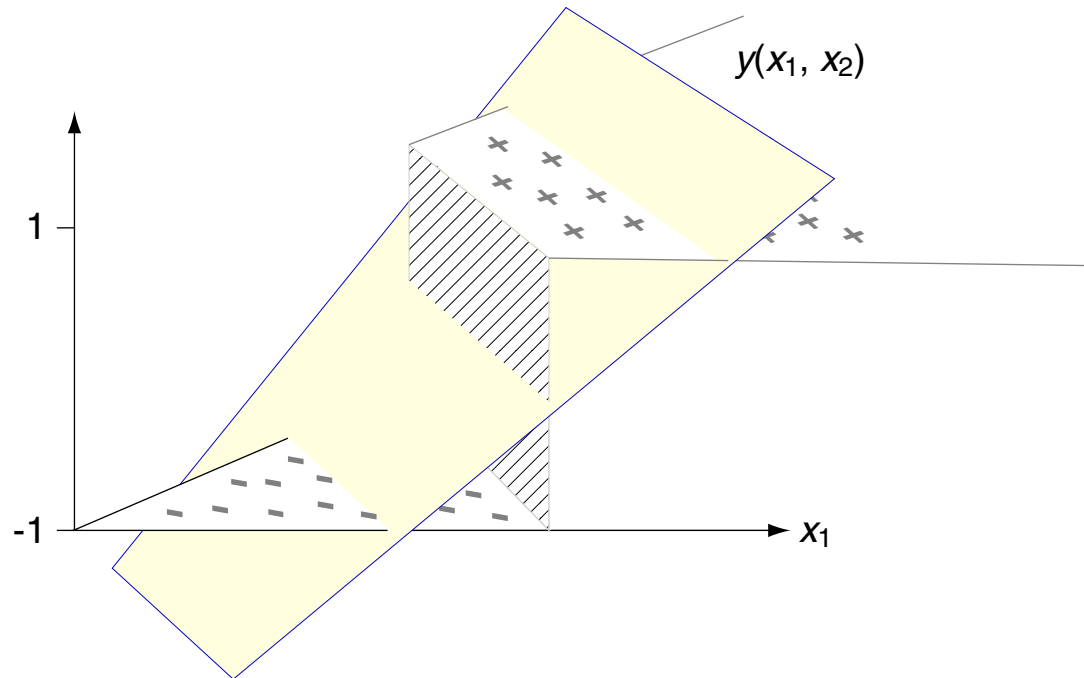
# Logistic Regression

Recap. Linear Regression for Classification (illustrated for  $p = 2$ )



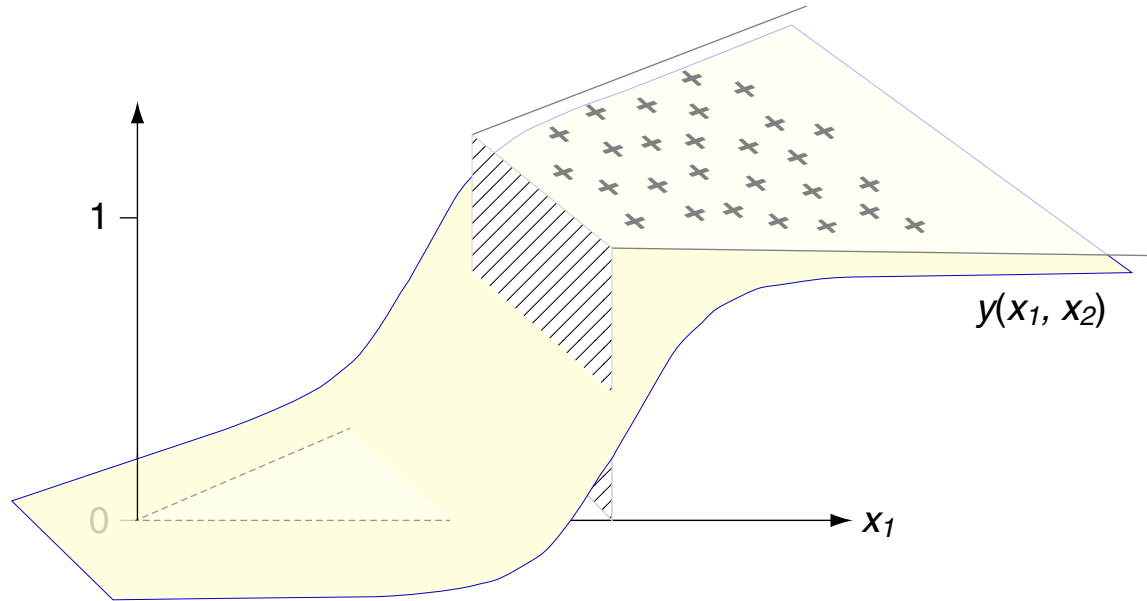
# Logistic Regression

Recap. Linear Regression for Classification (illustrated for  $p = 2$ ) (continued)



# Logistic Regression

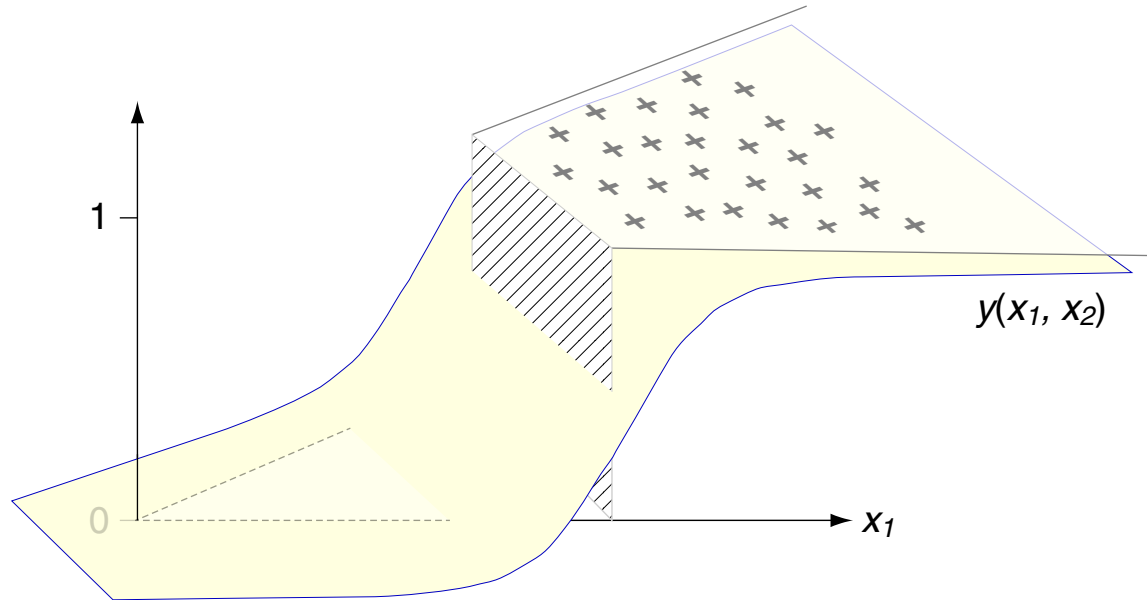
Logistic Regression for Classification (illustrated for  $p = 2$ )



# Logistic Regression

Logistic Regression for Classification (illustrated for  $p = 2$ ) (continued)

Use logistic regression to learn  $w$  from  $D$ , where  $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ .

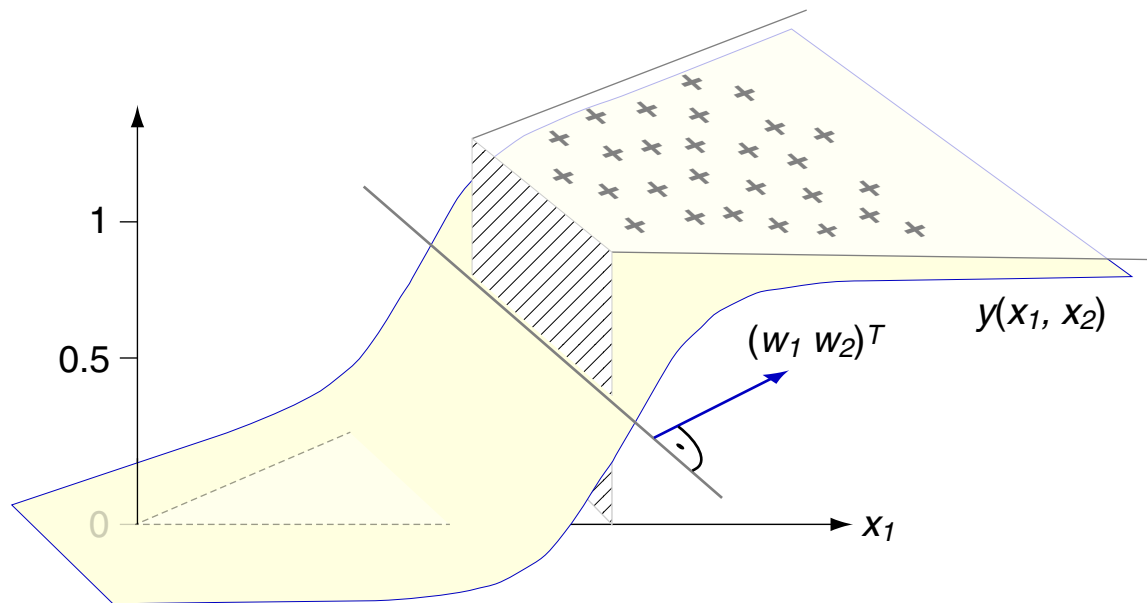


$$y(x_1, x_2) = \frac{1}{1 + e^{-(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)}}$$

# Logistic Regression

Logistic Regression for Classification (illustrated for  $p = 2$ ) (continued)

Use logistic regression to learn  $w$  from  $D$ , where  $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ .

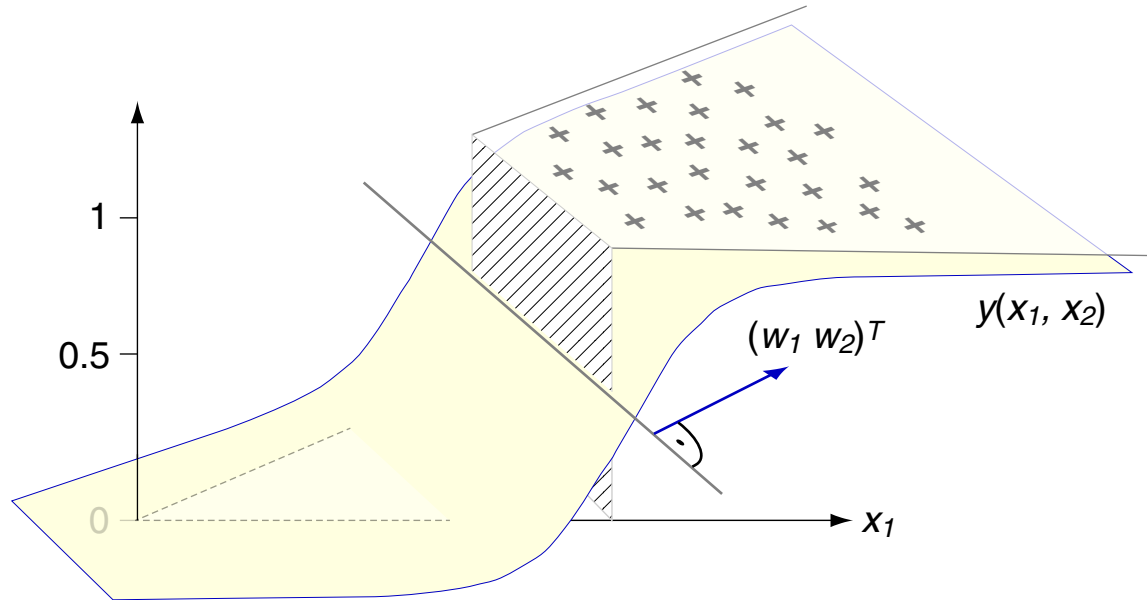


$$y(x_1, x_2) = \frac{1}{1 + e^{-(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)}}$$

# Logistic Regression

## Logistic Regression for Classification (illustrated for $p = 2$ ) (continued)

Use logistic regression to learn  $\mathbf{w}$  from  $D$ , where  $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ .

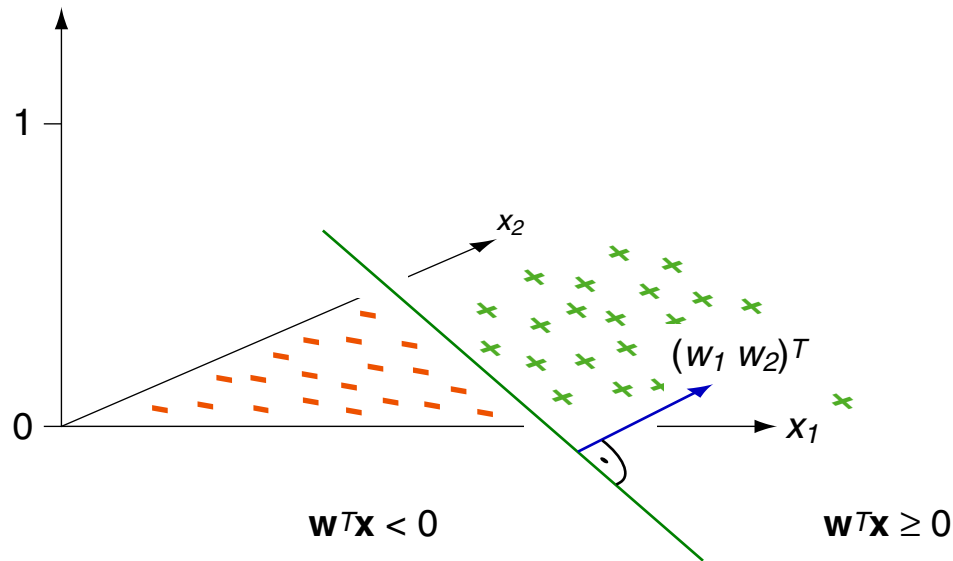


Classification: Predict  $\begin{cases} 1, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \\ 0, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) < 0.5 \end{cases}$

# Logistic Regression

## Logistic Regression for Classification (illustrated for $p = 2$ ) (continued)

Use logistic regression to learn  $\mathbf{w}$  from  $D$ , where  $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ .

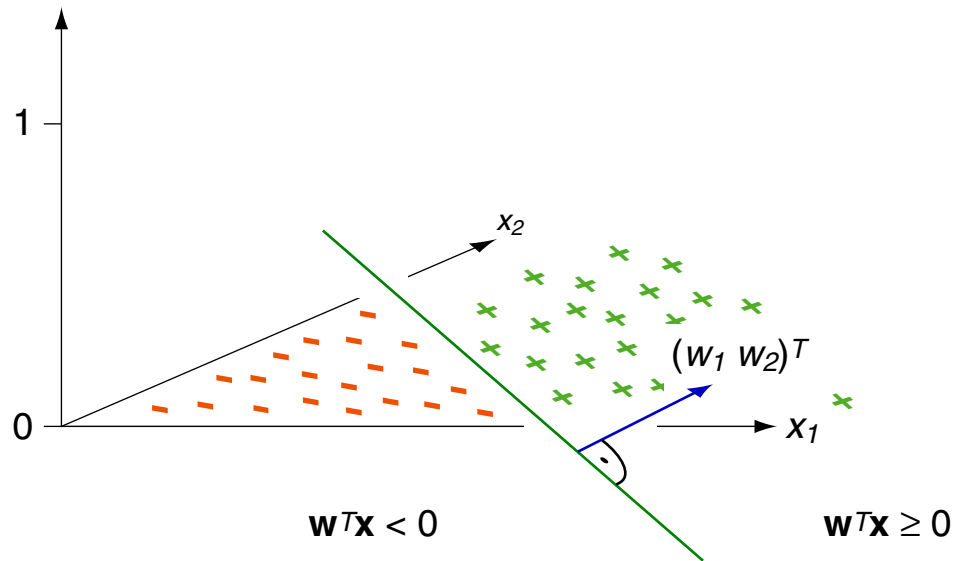


Classification: Predict  $\begin{cases} 1, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \\ 0, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) < 0.5 \end{cases}$

# Logistic Regression

## Logistic Regression for Classification (illustrated for $p = 2$ ) (continued)

Use logistic regression to learn  $\mathbf{w}$  from  $D$ , where  $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ .



Classification: Predict 
$$\begin{cases} 1, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 & \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0 \\ 0, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) < 0.5 & \Leftrightarrow \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$



# Logistic Regression

## The $\text{BGD}_\sigma$ Algorithm

[algorithms: LMS  $\text{BGD}_\sigma$  PT BGD IGD]

Algorithm:  $\text{BGD}_\sigma$  Batch Gradient Descent.

Input:  $D$  Multiset of examples  $(\mathbf{x}, c)$  with  $\mathbf{x} \in \mathbb{R}^p$ ,  $c \in \{0, 1\}$ .

$\eta$  Learning rate, a small positive constant.

Output:  $\mathbf{w}$  Weight vector from  $\mathbb{R}^{p+1}$ . (= hypothesis)

$\text{BGD}_\sigma(D, \eta)$

1. *initialize\_random\_weights*( $\mathbf{w}$ ),  $t = 0$
2. **REPEAT**
3.    $t = t + 1$ ,  $\Delta \mathbf{w} = 0$
4.   **FOREACH**  $(\mathbf{x}, c) \in D$  **DO**
5.      $y(\mathbf{x}) \stackrel{(*)}{=} \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$
6.      $\delta = c - y(\mathbf{x})$
7.      $\Delta \mathbf{w} \stackrel{(*)}{=} \Delta \mathbf{w} + \eta \cdot \delta \cdot \mathbf{x}$    //  $-\delta \cdot \mathbf{x}$  is the derivative of  $l_\sigma(c, y(\mathbf{x}))$  wrt.  $\mathbf{w}$ .
8.   **ENDDO**
9.    $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
10. **UNTIL**(*convergence*( $D, y(), t$ ))
11. *return*( $\mathbf{w}$ )

# Logistic Regression

## The $\text{BGD}_\sigma$ Algorithm

[algorithms: LMS  $\text{BGD}_\sigma$  PT BGD IGD]

Algorithm:  $\text{BGD}_\sigma$  Batch Gradient Descent.

Input:  $D$  Multiset of examples  $(\mathbf{x}, c)$  with  $\mathbf{x} \in \mathbb{R}^p$ ,  $c \in \{0, 1\}$ .

$\eta$  Learning rate, a small positive constant.

Output:  $\mathbf{w}$  Weight vector from  $\mathbb{R}^{p+1}$ . (= hypothesis)

$\text{BGD}_\sigma(D, \eta)$

1. *initialize\_random\_weights*( $\mathbf{w}$ ),  $t = 0$
2. **REPEAT**
3.  $t = t + 1$ ,  $\Delta \mathbf{w} = 0$
4. **FOREACH**  $(\mathbf{x}, c) \in D$  **DO**
5.  Model function evaluation.
6.  Calculation of residual.
7.  Calculation of derivative of the loss, accumulate for  $D$ .
8. **ENDDO**
9.  Parameter vector update  $\hat{=}$  one gradient step down.
10. **UNTIL**(*convergence*( $D, y(), t$ ))
11. *return*( $\mathbf{w}$ )

## Remarks (BGD<sub>σ</sub> Algorithm) :

- ❑ The BGD<sub>σ</sub> Algorithm is an iterative method to estimate  $\mathbf{w}_{\text{ML}}$  in the model function  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$ . There is no direct method (such as the normal equations in linear regression) to tackle the optimization problem.
- ❑ The BGD<sub>σ</sub> Algorithm exploits the derivative of the pointwise logistic loss  $l_\sigma(c, y(\mathbf{x}))$  with respect to  $\mathbf{w}$ , which is  $-\delta \cdot \mathbf{x} = -(c - y(\mathbf{x})) \cdot \mathbf{x} = -(c - \sigma(\mathbf{w}^T \mathbf{x})) \cdot \mathbf{x}$ . The derivation of this term, as well as notes regarding the speed of convergence of the basic gradient descent are given in section [Gradient Descent in Detail](#) of part Linear Models.
- ❑ Each BGD<sub>σ</sub> iteration “REPEAT ... UNTIL”
  1. computes the direction of steepest loss descent as
$$-\nabla L_\sigma(\mathbf{w}_t) = \sum_{(\mathbf{x}, c) \in D} (c - y_t(\mathbf{x})) \cdot \mathbf{x}, \text{ and}$$
  2. updates  $\mathbf{w}_t$  by taking a step of length  $\eta$  in this direction.
- ❑ The output of the logistic model (function),  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$ , is a continuous variable in the domain  $(0, 1)$ , while the dependent variable that we use for training is a Bernoulli-distributed:  $c = 0 \vee c = 1$ .

## Remarks: (continued)

- Recap. A hypothesis is a proposed explanation for a phenomenon. [\[Wikipedia\]](#)

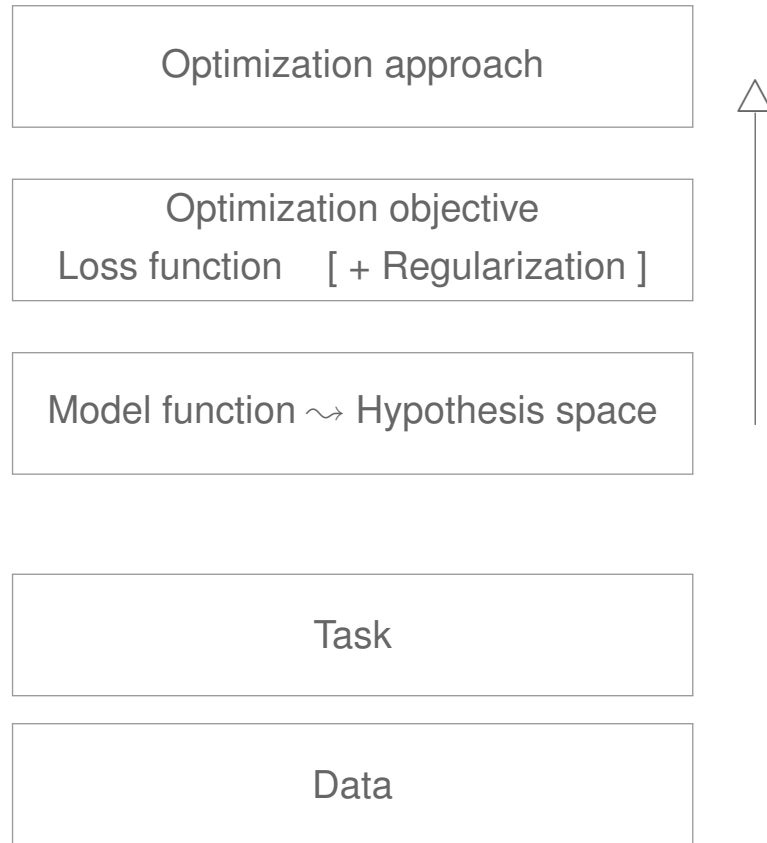
Here, hypothesis “explains” (= fits) the data  $D$ . Hence, a concrete model function  $y()$ ,  $\mathbf{y}()$ , or, if the function type is clear from the context, its parameters  $\mathbf{w}$  or  $\boldsymbol{\theta}$  are called “hypothesis”.

The variable name  $h$  (similarly:  $h_1$ ,  $h_2$ ,  $h_i$ ,  $h'$ , etc.) may be used to refer to a specific instance of a model function or its parameters.

- Recap. The function *convergence*() can analyze the global logistic loss,  $L_\sigma(\mathbf{w}_t)$ , or the norm of the loss gradient,  $\|\nabla L_\sigma(\mathbf{w}_t)\|$ , and compare it to a small positive bound  $\varepsilon$ . Consider in this regard the vectors of observed and computed classes,  $D|_c$  and  $y(D|_x)$  respectively. In addition, the function may check via  $t$  an upper bound on the number of iterations.

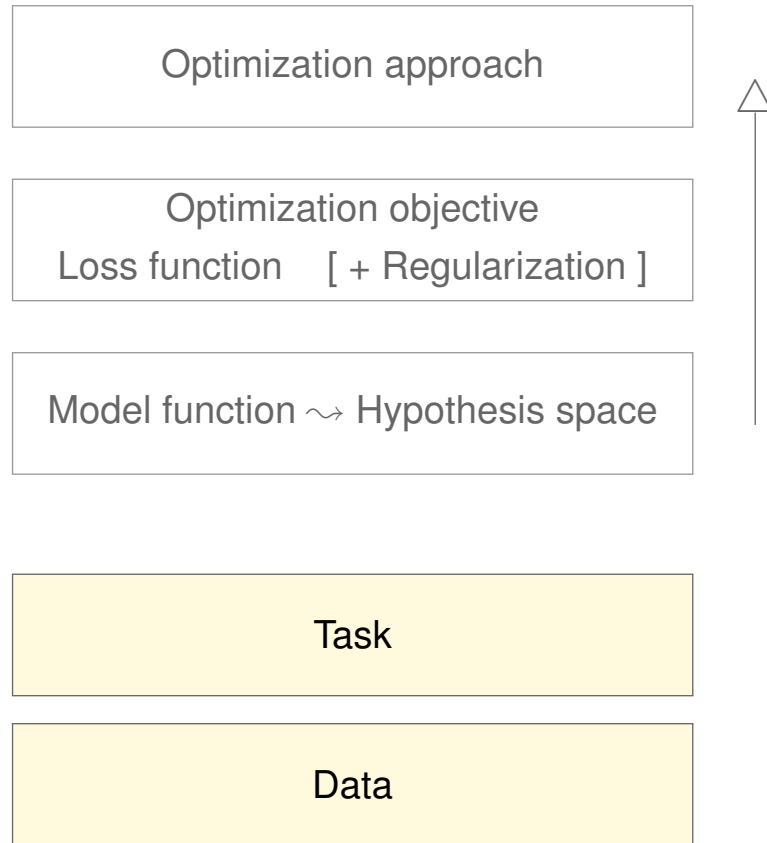
# Logistic Regression

## Machine Learning Stack of Logistic Regression



# Logistic Regression

## Machine Learning Stack of Logistic Regression (continued)

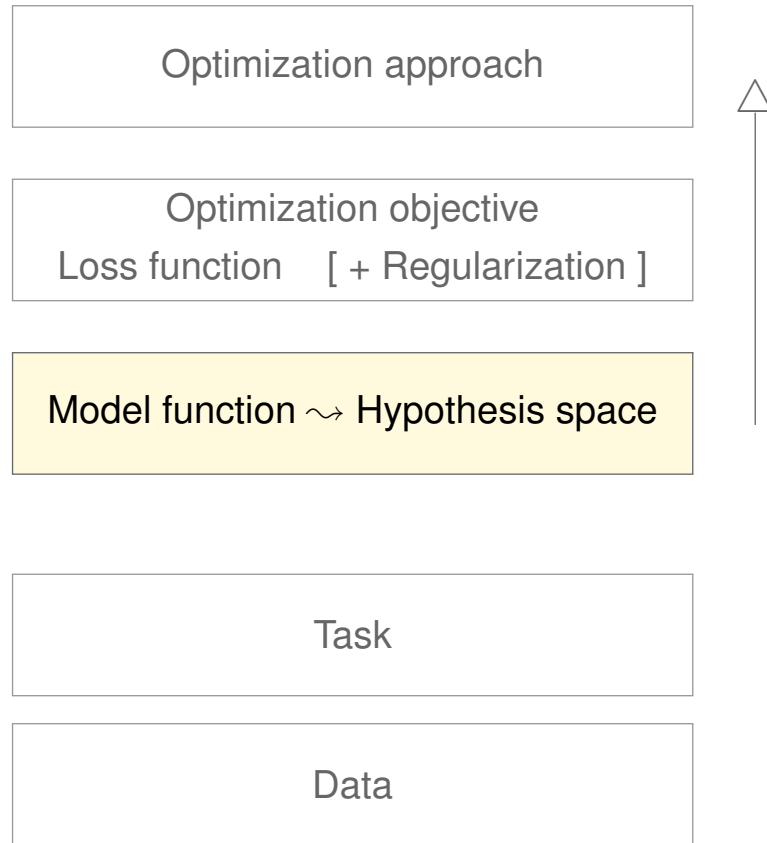


Binary classification

$$D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times \{0, 1\}$$

# Logistic Regression

## Machine Learning Stack of Logistic Regression (continued)



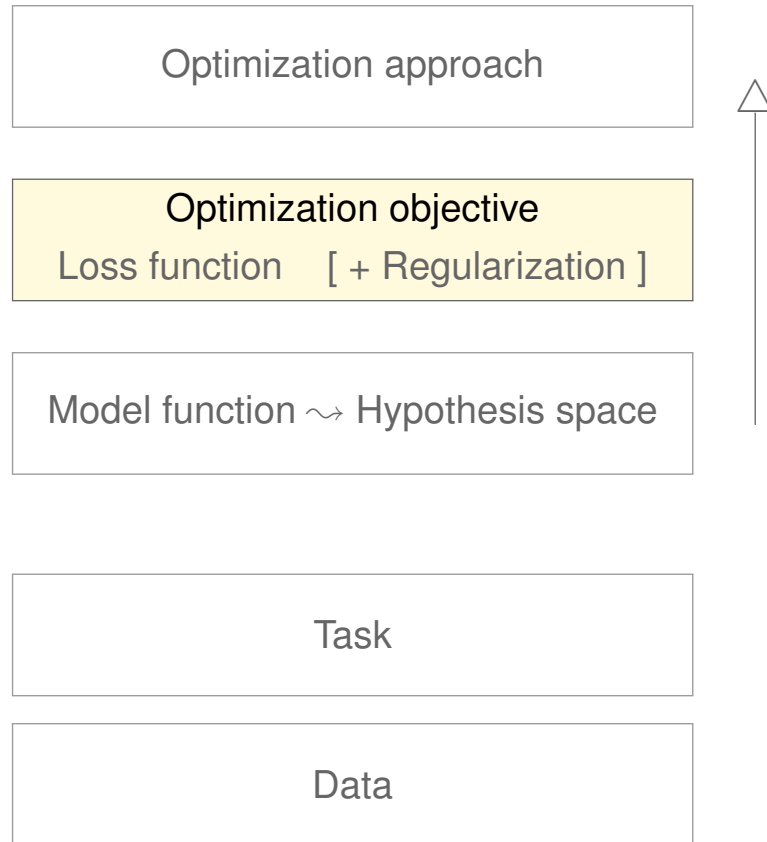
- ❑ Hypothesis space:  $\mathbf{w} \in \mathbf{R}^{p+1}$
- ❑ Logistic model:  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$

Binary classification

$$D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times \{0, 1\}$$

# Logistic Regression

## Machine Learning Stack of Logistic Regression (continued)



- ❑ Objective: minimize logistic loss  $L_\sigma(\mathbf{w})$
- ❑ Regularization: none
- ❑ Loss:  $l_\sigma(c, y(\mathbf{x})) = -c \cdot \log(y(\mathbf{x})) - (1-c) \cdot \log(1-y(\mathbf{x}))$
- ❑ Hypothesis space:  $\mathbf{w} \in \mathbb{R}^{p+1}$
- ❑ Logistic model:  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$

Binary classification

$$D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times \{0, 1\}$$



# Logistic Regression

## Machine Learning Stack of Logistic Regression (continued)

Optimization approach

Optimization objective  
Loss function [ + Regularization ]

Model function  $\leadsto$  Hypothesis space

Task

Data



BGD, Newton-Raphson, BFGS, Conjugate GD

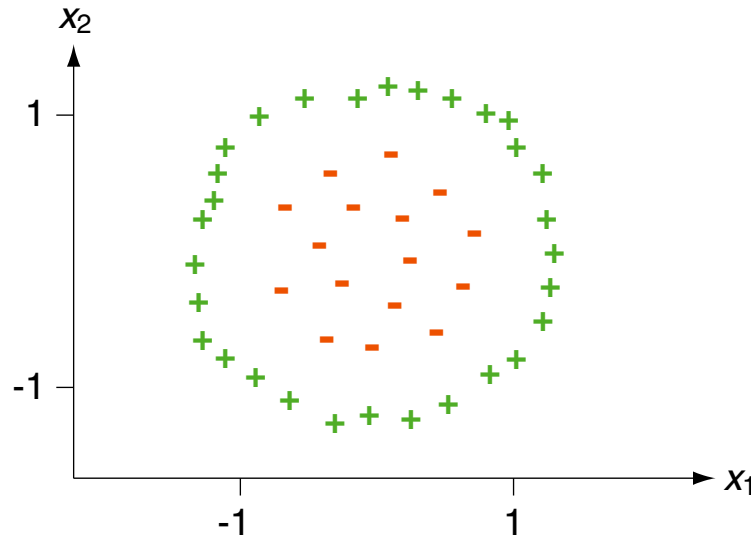
- ❑ Objective: minimize logistic loss  $L_\sigma(\mathbf{w})$
- ❑ Regularization: none
- ❑ Loss:  $l_\sigma(c, y(\mathbf{x})) = -c \cdot \log(y(\mathbf{x})) - (1-c) \cdot \log(1-y(\mathbf{x}))$
- ❑ Hypothesis space:  $\mathbf{w} \in \mathbb{R}^{p+1}$
- ❑ Logistic model:  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$

Binary classification

$$D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_n, c_n)\} \subseteq X \times \{0, 1\}$$

# Logistic Regression

## Non-Linear Decision Boundaries [linear regression]

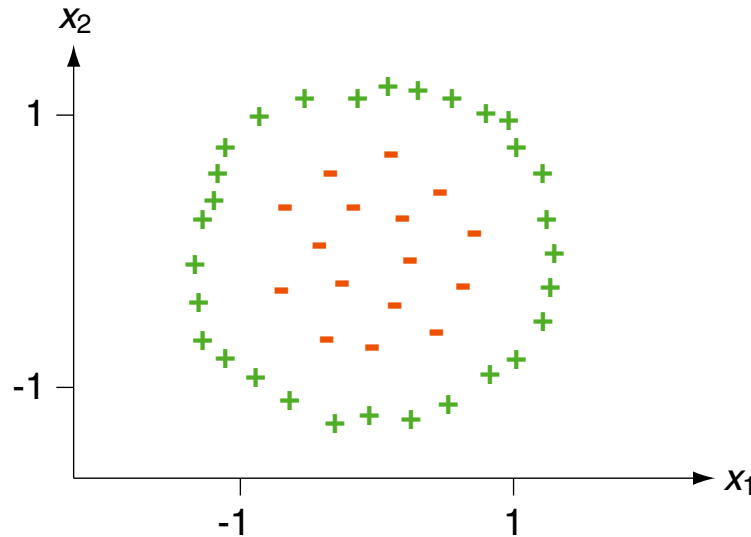


Higher order polynomial terms in the features (linear in the parameters):

$$y(\mathbf{x}) = \sigma(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_1^2 + w_4 \cdot x_2^2)$$

# Logistic Regression

## Non-Linear Decision Boundaries (continued) [linear regression]

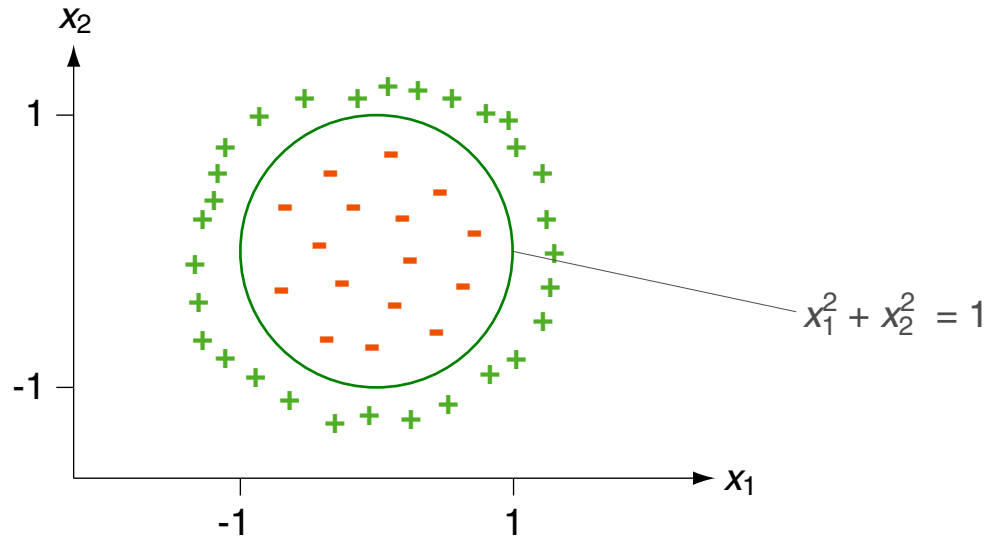


Higher order polynomial terms in the features (linear in the parameters):

$$y(\mathbf{x}) = \sigma(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_1^2 + w_4 \cdot x_2^2)$$

# Logistic Regression

## Non-Linear Decision Boundaries (continued) [linear regression]

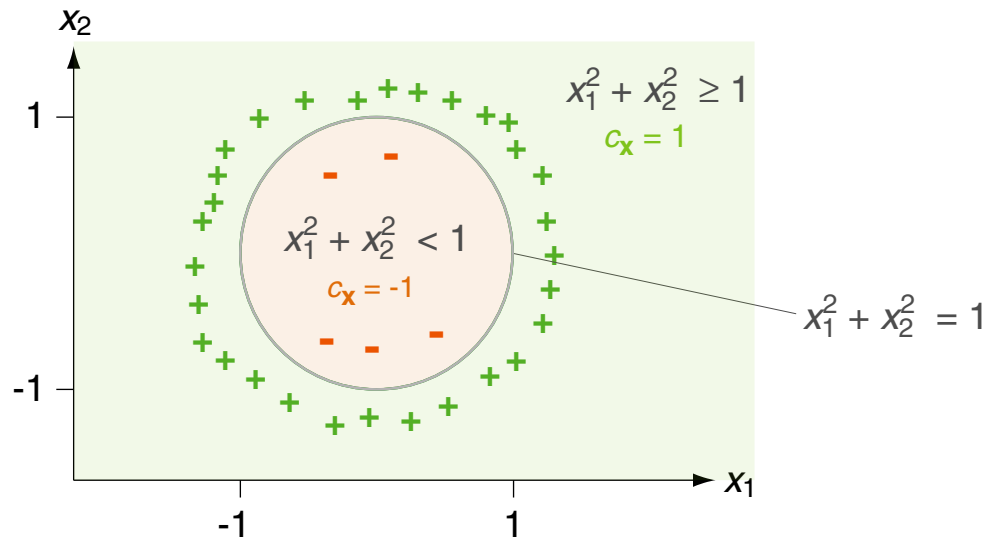


Higher order polynomial terms in the features (linear in the parameters):

$$\text{with } \mathbf{w} = \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} \quad \rightsquigarrow \quad y(\mathbf{x}) = \frac{1}{1 + e^{-(-1 + x_1^2 + x_2^2)}}$$

# Logistic Regression

## Non-Linear Decision Boundaries (continued) [linear regression]



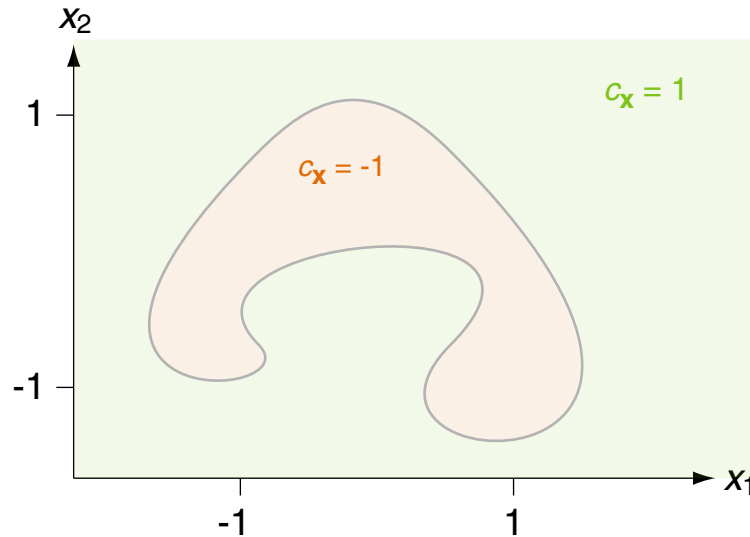
Higher order polynomial terms in the features (linear in the parameters):

$$\text{with } \mathbf{w} = \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix} \quad \rightsquigarrow \quad y(\mathbf{x}) = \frac{1}{1 + e^{-(-1 + x_1^2 + x_2^2)}}$$

Classification: Predict  $\begin{cases} 1, & \text{if } x_1^2 + x_2^2 \geq 1 & \Leftrightarrow \mathbf{w}^T \mathbf{x} \geq 0 \\ 0, & \text{if } x_1^2 + x_2^2 < 1 & \Leftrightarrow \mathbf{w}^T \mathbf{x} < 0 \end{cases}$

# Logistic Regression

## Non-Linear Decision Boundaries (continued) [linear regression]



More complex polynomials entail more complex decision boundaries:

$$y(\mathbf{x}) = \sigma(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_1^2 + w_4 \cdot x_1^2 \cdot x_2 + w_5 \cdot x_1^2 \cdot x_2^2 + \dots)$$

## Remarks:

- Under logistic regression the structure of a hypothesis, i.e., the forms of possible decision boundaries, is identical to the hypothesis structure under linear regression. Similarly, the respective hypothesis spaces are the same.

Hence, the expressiveness, i.e., the complexity of classification problems that can be tackled (or, the effectiveness at which classification problems can be decided) is identical for the two regression approaches.

- Linear regression and logistic regression differ in the way how the model function parameters,  $\mathbf{w}$ , are determined. In both cases the optimum  $\mathbf{w}$  is the result of a loss minimization problem.

Recall that “loss” means “interpretation of residuals.” Linear regression and logistic regression differ with respect to this interpretation: While the former simply squares the residuals, this way putting a high weight onto outliers, the latter models an increasing confidence in class membership probability with increasing hyperplane distance. This different interpretation will usually lead to a different parameter vector  $\mathbf{w}$ , i.e., a different hyperplane.

# Chapter ML:III

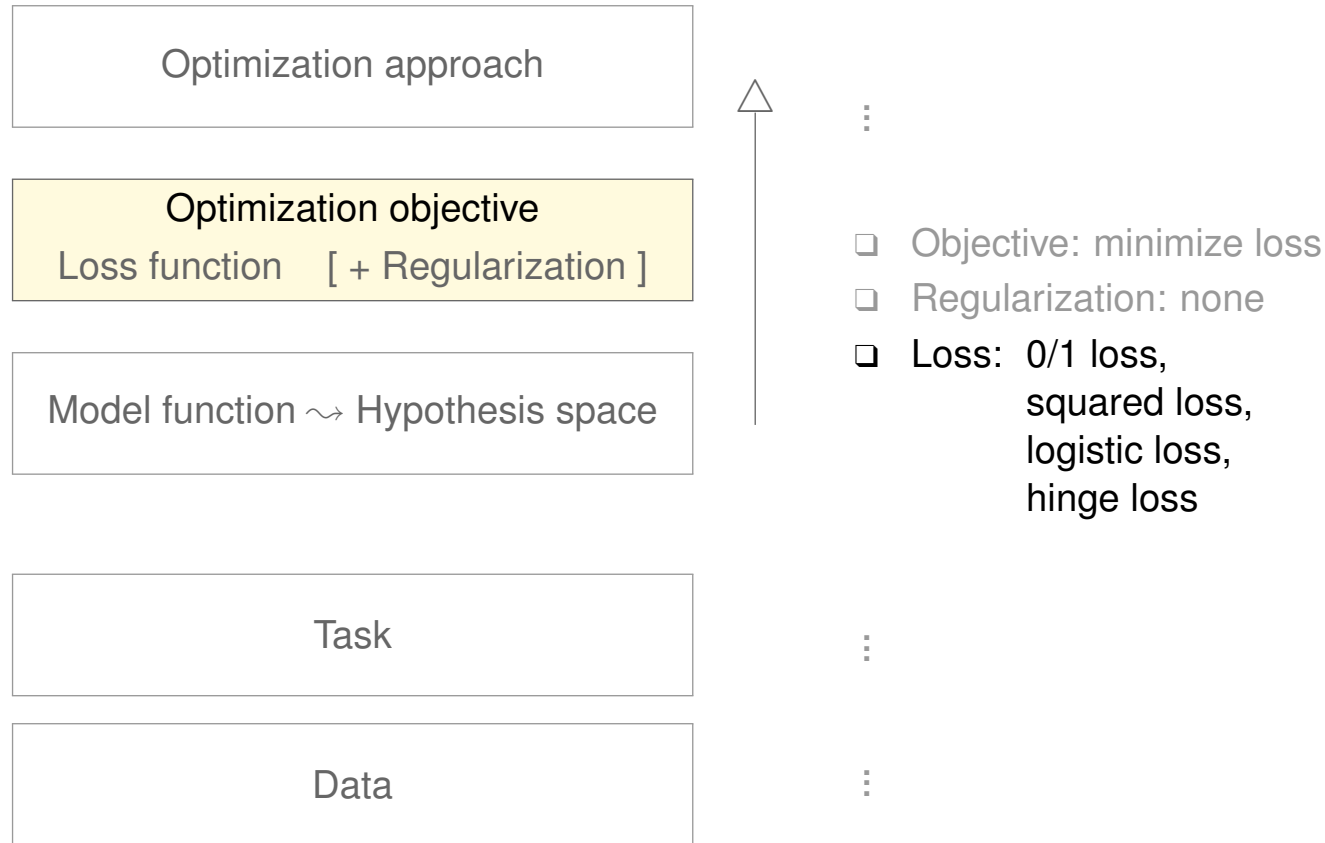
## III. Linear Models

- ❑ Logistic Regression
- ❑ Loss Computation in Detail
- ❑ Overfitting
- ❑ Regularization
- ❑ Gradient Descent in Detail



# Loss Computation in Detail

## Loss Computation in the Machine Learning Stack [ML stack: logistic regression]



## Remarks (loss function) :

- ❑ Given a hypothesis  $\mathbf{w}$ , its (global) loss,  $L(\mathbf{w})$ , tells us something about the effectiveness of  $\mathbf{w}$ . When used as sole criterion (e.g., no regularization is applied) we select from two hypotheses that with the smaller loss. I.e., the most effective hypothesis is found by loss minimization. Conversely, we call a function, whose minimization determines the most effective hypothesis, a loss function.
- ❑ Loss functions can be distinguished with respect to the problem class they are typically applied to: regression versus classification. Keep in mind that this distinction is not unique since loss functions with continuous output are applied to classification problems as well.
- ❑ Furthermore, we distinguish the
  1. *pointwise loss*  $l(c, y(\mathbf{x}))$ , which is computed for a single  $\mathbf{x}$ , and the
  2. *global loss*  $L(\mathbf{w})$ , which accumulates the pointwise losses of all  $\mathbf{x} \in X$  for the weight vector  $\mathbf{w}$  used in a specific  $y()$ :

$$L(\mathbf{w}) = \sum_{(\mathbf{x}, c) \in D} l(c, y(\mathbf{x}))$$

The pointwise loss is also called per-example loss. [p.268, Goodfellow/Bengio/Courville 2016]

- ❑ Instead of “loss” (function, computation) also the terms “error” (function, computation), “cost” (function, computation), or “performance” (function, computation) are used, usually with the same semantics as introduced here. We will use the term “error” for classification problems and the term “loss” for both classification and regression problems.

Remarks (different roles of loss functions) :

- Observe that loss functions are employed at two places (in two roles) in an optimization approach:
  1. For the fitting of the data (i.e., the parameter update during regression / optimization / hyperplane search), where a new position of the hyperplane is computed.  
Example: Lines 6+7 in the [BGD <sub>\$\sigma\$</sub>  Algorithm](#), which is illustrated [here](#) over the hyperplane distance (red curve in the figure at the bottom right).
  2. For the evaluation of the effectiveness of a hypothesis, where the proportion of correctly and misclassified examples is analyzed.  
Example: Line 10 in the [BGD <sub>\$\sigma\$</sub>  Algorithm](#).  
General: section [Evaluating Effectiveness](#) of part Machine Learning Basics.

Typically, (1) fitting (optimization) and (2) effectiveness evaluation are done with different loss functions. E.g., logistic regression uses  $L_\sigma$  and  $L_{0/1}$  for fitting and evaluation respectively. However, linear regression (not classification) uses RSS (the  $L_2$  loss) for both fitting and evaluation. The basic perceptron learning algorithm uses the misclassification information (the  $L_{0/1}$  loss) for both fitting and evaluation.

# Loss Computation in Detail

## Linear Regression

- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

# Loss Computation in Detail

## Linear Regression (continued)

- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y(\cdot)$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

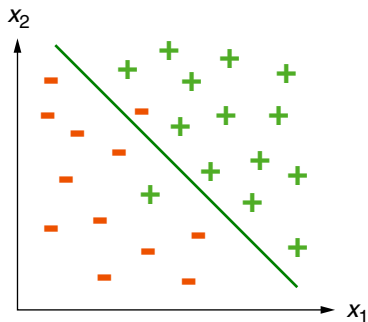
# Loss Computation in Detail

## Linear Regression (continued)

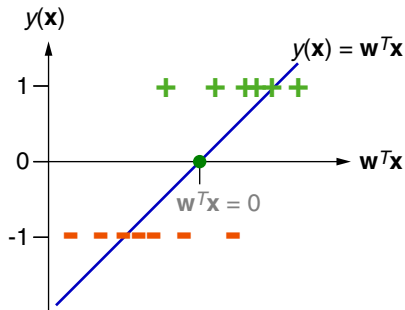
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

Illustration for a particular  $\mathbf{w}$ :

Input space:



$y()$  over hyperplane distance:

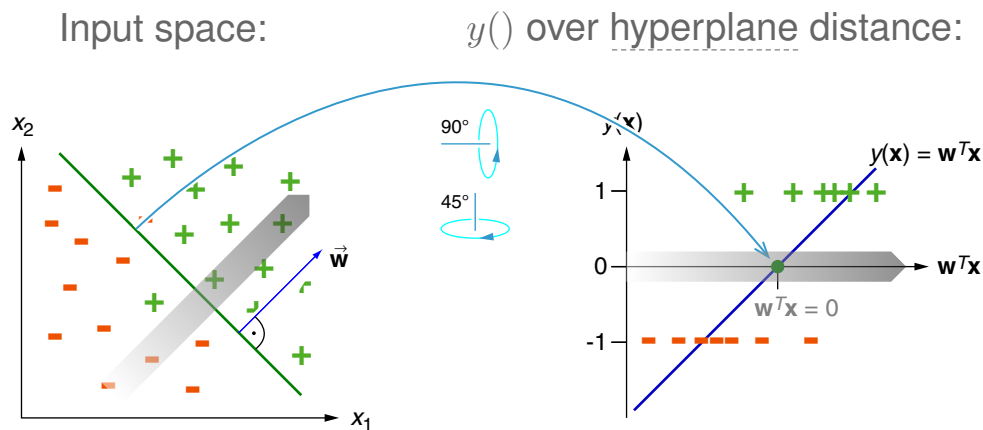


# Loss Computation in Detail

## Linear Regression (continued)

- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

Illustration for a particular  $\mathbf{w}$ :



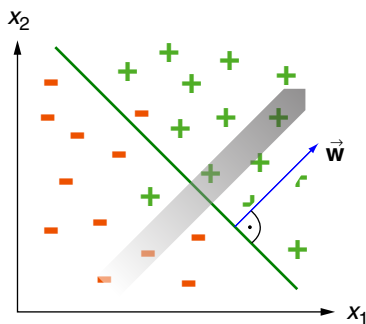
# Loss Computation in Detail

## Linear Regression (continued)

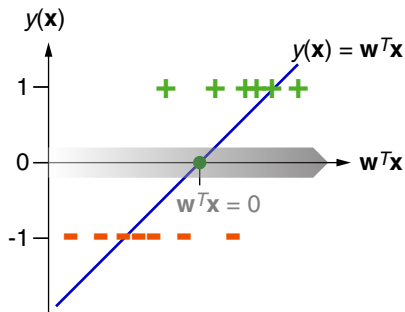
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

Illustration for a particular  $\mathbf{w}$ :

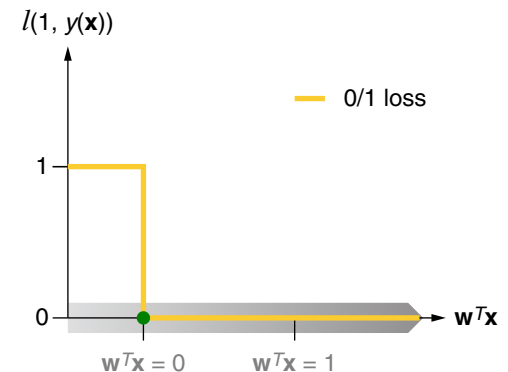
Input space:



$y()$  over hyperplane distance:



Loss over hyperplane distance:





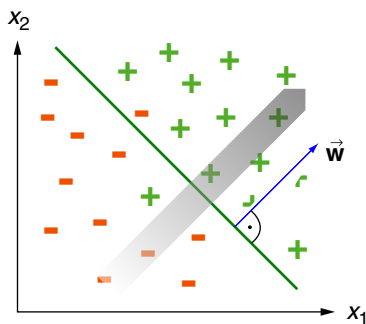
# Loss Computation in Detail

## Linear Regression (continued)

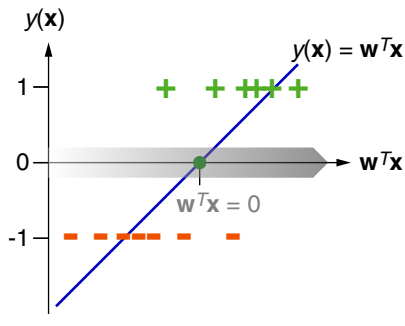
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

Illustration for a particular  $\mathbf{w}$ :

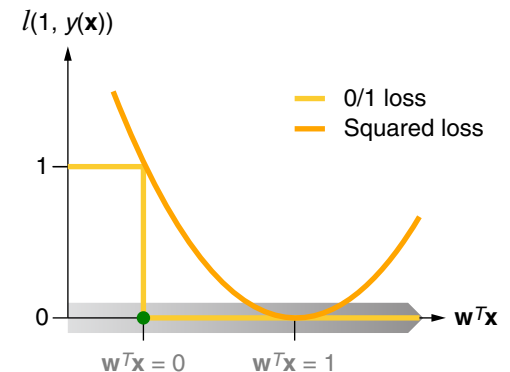
Input space:



$y()$  over hyperplane distance:



Loss over hyperplane distance:



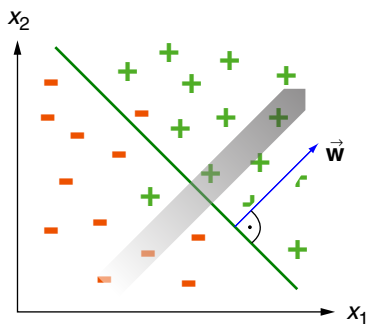
# Loss Computation in Detail

## Linear Regression (continued)

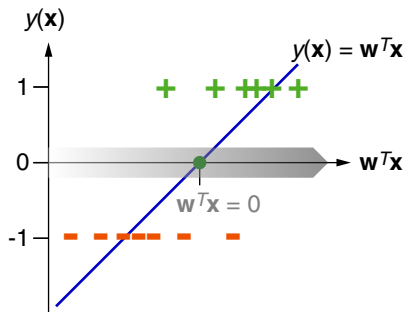
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y()$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \text{sign}(y(\mathbf{x}))) = \begin{cases} 0 & \text{if } c = \text{sign}(y(\mathbf{x})) \\ 1 & \text{otherwise} \end{cases}$
  - Squared loss.  $l_2(c, y(\mathbf{x})) = (c - y(\mathbf{x}))^2$

Illustration for a particular  $\mathbf{w}$ :

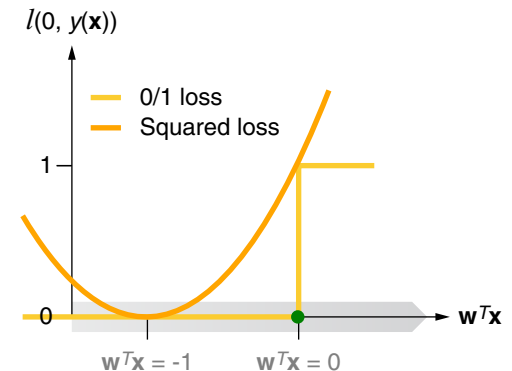
Input space:



$y()$  over hyperplane distance:



Loss over hyperplane distance:



## Remarks:

- ❑ The 0/1 loss computes the misclassification error. Recall in this regard the definition of the true misclassification rate.
- ❑ The pointwise squared loss computes the squared residual. The global squared loss,  $L_2(\mathbf{w}) = \sum_{(\mathbf{x}, c) \in D} l_2(c, y(\mathbf{x}))$ , hence computes the residual sum of squares (RSS) related to some  $\mathbf{w}$ .
- ❑  $I_{\neq}$  is an indicator function that returns 1 if its arguments are *unequal* (and 0 if its arguments are equal).
- ❑ Recap. We label  $y(0)$  with the “positive” class and define  $\text{sign}(0) = 1$  here.
- ❑ Regarding the illustration:  $\mathbf{w}^T \mathbf{x}$  is the hyperplane distance in relation to  $\|\mathbf{w}\|$ , the length of  $\mathbf{w}$ . By scaling  $\mathbf{w}$  such that  $\|\mathbf{w}\| = 1$  the hyperplane distance  $\mathbf{w}^T \mathbf{x}$  becomes normalized and is also called “geometric distance”.

# Loss Computation in Detail

## Logistic Regression

- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y(\cdot)$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \lfloor y(\mathbf{x}) + 0.5 \rfloor)$  [decision rule]
  - Logistic loss.  $l_{\sigma}(c, y(\mathbf{x})) = \begin{cases} -\log(y(\mathbf{x})) & \text{if } c = 1 \\ -\log(1 - y(\mathbf{x})) & \text{if } c = 0 \end{cases}$  [derivation]

# Loss Computation in Detail

## Logistic Regression (continued)

- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y(\cdot)$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \lfloor y(\mathbf{x}) + 0.5 \rfloor)$  [decision rule]
  - Logistic loss.  $l_{\sigma}(c, y(\mathbf{x})) = \begin{cases} -\log(y(\mathbf{x})) & \text{if } c = 1 \\ -\log(1 - y(\mathbf{x})) & \text{if } c = 0 \end{cases}$  [derivation]

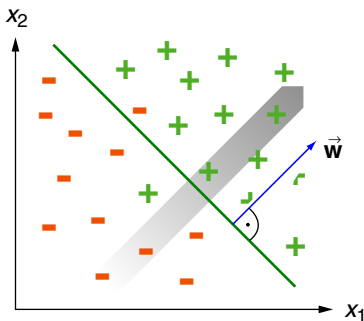
# Loss Computation in Detail

## Logistic Regression (continued)

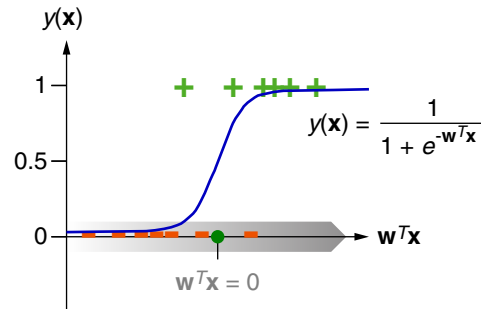
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y(\cdot)$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \lfloor y(\mathbf{x}) + 0.5 \rfloor)$  [decision rule]
  - Logistic loss.  $l_{\sigma}(c, y(\mathbf{x})) = \begin{cases} -\log(y(\mathbf{x})) & \text{if } c = 1 \\ -\log(1 - y(\mathbf{x})) & \text{if } c = 0 \end{cases}$  [derivation]

Illustration for a particular  $\mathbf{w}$ :

Input space:



$y(\cdot)$  over hyperplane distance:



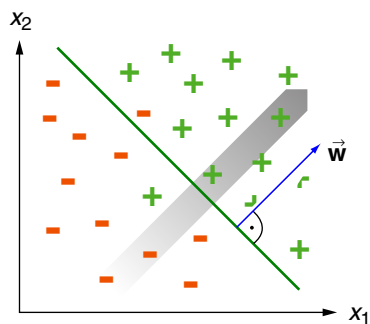
# Loss Computation in Detail

## Logistic Regression (continued)

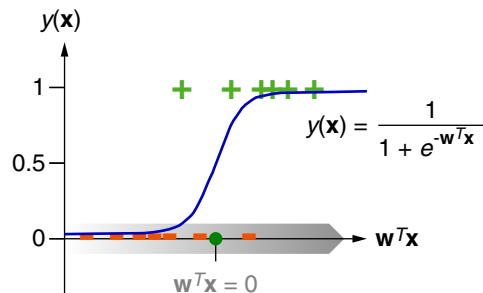
- The pointwise loss,  $l(c, y(\mathbf{x}))$ , quantifies the error introduced by some  $\mathbf{x}$ . The loss depends on a hypothesis  $y(\cdot)$  and the true class,  $c$ , of  $\mathbf{x}$ .
- For  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$  we define the following pointwise loss functions:
  - 0/1 loss.  $l_{0/1}(c, y(\mathbf{x})) = I_{\neq}(c, \lfloor y(\mathbf{x}) + 0.5 \rfloor)$  [decision rule]
  - Logistic loss.  $l_{\sigma}(c, y(\mathbf{x})) = \begin{cases} -\log(y(\mathbf{x})) & \text{if } c = 1 \\ -\log(1 - y(\mathbf{x})) & \text{if } c = 0 \end{cases}$  [derivation]

Illustration for a particular  $\mathbf{w}$ :

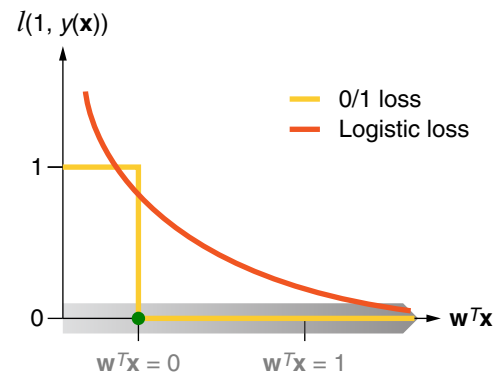
Input space:



$y(\cdot)$  over hyperplane distance:



Loss over hyperplane distance:



## Remarks:

- As before, the 0/1 loss computes the misclassification error.
- The pointwise logistic loss can be rewritten by combining the two cases algebraically:

$$l_{\sigma}(c, y(\mathbf{x})) = -c \cdot \log(y(\mathbf{x})) - (1 - c) \cdot \log(1 - y(\mathbf{x}))$$

$L_{\sigma}(\mathbf{w}) = \sum_{(\mathbf{x}, c) \in D} l_{\sigma}(c, y(\mathbf{x}))$  computes the global logistic loss related to some  $\mathbf{w}$ .

- Recall from the [derivation](#) of the logistic loss  $L_{\sigma}(\mathbf{w})$  that its minimization determines  $\mathbf{w}_{\text{ML}}$ , the most probable hypothesis in  $\mathbb{R}^{p+1}$  under the logistic regression model.
- [Recap](#).  $I_{\neq}$  is an indicator function that returns 1 if its arguments are *unequal* (and 0 if its arguments are equal).
- [Recap](#). We label  $y(0)$  with the “positive” class, which is consistent with the definition of the [floor function](#), since  $\lfloor y(0) + 0.5 \rfloor = 1$ .
- Q. Why not use mean squared error (MSE) as a loss function for logistic regression?  
A. MSE is not convex as loss function for logistic regression. [\[ai.stackexchange\]](#) [\[Varma 2018\]](#)  
[\[Nabert 2021\]](#) [\[DeLuca 2023\]](#)