# *** Applied Machine Learning Fundamentals ***
## Logistic Regression

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Winter term 2021/2022

# Lecture Overview

| | |
|---|---|
| **Unit I** | Machine Learning Introduction |
| **Unit II** | Mathematical Foundations |
| **Unit III** | Bayesian Decision Theory |
| **Unit IV** | Probability Density Estimation |
| **Unit V** | Regression |
| **Unit VI** | **Classification I** |
| **Unit VII** | Evaluation |
| **Unit VIII** | Classification II |
| **Unit IX** | Clustering |
| **Unit X** | Dimensionality Reduction |

# Agenda for this Unit

**Section:**

**Introduction**

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

What is logistic Regression?
Why you should not use linear Regression

# What is logistic Regression?

- Learning algorithm for **classification** *(despite the name...)*
- In its standard form it's applicable to **binary classification problems only**, but you can use techniques like:
    - **One-vs-One (OVO)**
    - **One-vs-Rest (OVR)**
- **Class labels:**
    - The 'positive class' $\oplus$ is encoded as **1**
    - The 'negative class' $\ominus$ as **0**
- **Probabilistic interpretation:** The output of the algorithm is between 0 and 1 *(probability of the instance belonging to the positive class)*

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

What is logistic Regression?
Why you should not use linear Regression

# Why you should not use linear Regression...

**Introduction**
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

What is logistic Regression?
Why you should not use linear Regression

# Why you should not use linear Regression...

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up
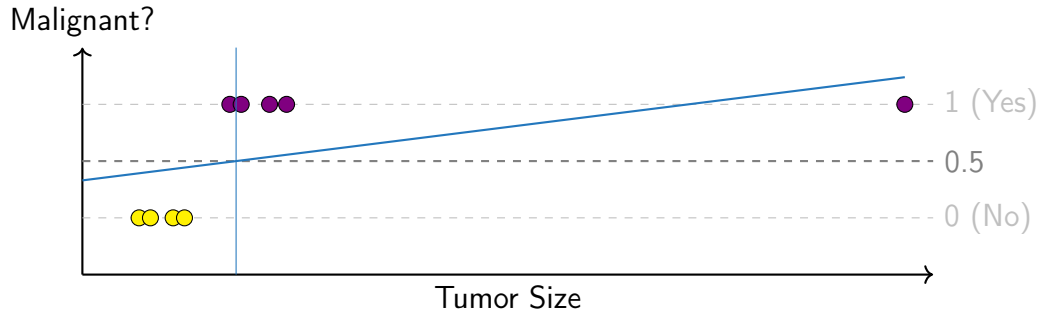
What is logistic Regression?
Why you should not use linear Regression

# Why you should not use linear Regression... (Ctd.)

- Linear regression: $h_{\theta}(\boldsymbol{x}) = \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$

- By putting a **threshold** at 0.5, we can turn linear regression into a classifier
  - If $h_{\theta}(\boldsymbol{x}) \geqslant 0.5$, predict $y = 1$
  - If $h_{\theta}(\boldsymbol{x}) < 0.5$, predict $y = 0$

- **Problems:**
  1. **Outliers heavily affect the decision boundary**
  2. Furthermore, we only want $0 \leqslant h_{\theta}(\boldsymbol{x}) \leqslant 1$, linear regression can output values $h_{\theta}(\boldsymbol{x}) \ll 0$ or $h_{\theta}(\boldsymbol{x}) \gg 1$

- We need a better strategy!

# Section:
## Model Architecture

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

**Sigmoid Function**
Probabilistic Interpretation
Model Training
Decision Boundary

## Logistic Regression Model

- Remember that we want: $0 \leqslant h_{\boldsymbol{\theta}}(\boldsymbol{x}) \leqslant 1$

- **Solution**: Logistic / Sigmoid function:

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

- We plug $\boldsymbol{\theta}^{\intercal}\boldsymbol{x}$ into the sigmoid function:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g(\boldsymbol{\theta}^{\intercal}\boldsymbol{x}) = \frac{1}{1 + e^{-(\boldsymbol{\theta}^{\intercal}\boldsymbol{x})}} \tag{2}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
Decision Boundary

# Logistic/Sigmoid Function



- $g(z)$ is symmetric around $z = 0$
- $0 \leqslant g(z) \leqslant 1$ holds true

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

**Sigmoid Function**
Probabilistic Interpretation
Model Training
Decision Boundary

# Where does the Sigmoid come from?

$$p(\mathcal{C}_1|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\boldsymbol{x})} = \frac{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{\sum_j p(\boldsymbol{x}, \mathcal{C}_j)} = \frac{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{\sum_j p(\boldsymbol{x}|\mathcal{C}_j)p(\mathcal{C}_j)}$$

$$= \frac{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\boldsymbol{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

$$= \frac{1}{1 + p(\boldsymbol{x}|\mathcal{C}_2)p(\mathcal{C}_2)/(p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1))}$$

$$= \frac{1}{1 + \exp\{-z\}} = g(z) \qquad\qquad \longrightarrow \textbf{logistic sigmoid}$$

$$z = \log \frac{p(\boldsymbol{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\boldsymbol{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \qquad\qquad \longrightarrow \textbf{log odds}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
**Probabilistic Interpretation**
Model Training
Decision Boundary

# Interpretation of Hypothesis Output

- $h_{\theta}(\boldsymbol{x})$ is interpreted as the probability of instance $\boldsymbol{x}$ belonging to class $y = 1$
- **Example:**

$$\boldsymbol{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ tumorSize \end{bmatrix} \tag{3}$$

- If $h_{\theta}(\boldsymbol{x}) = 0.7$, we have to tell the patient that there is a **70 % chance** of the tumor being malignant $\Rightarrow p(y = 1 | \boldsymbol{x}, \boldsymbol{\theta})$
- **Binary case:** $p(y = 0 | \boldsymbol{x}, \boldsymbol{\theta}) = 1 - p(y = 1 | \boldsymbol{x}, \boldsymbol{\theta})$

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
Decision Boundary

## Training Setup

- We have a labeled training set ($\Rightarrow$ **supervised learning**):

$$\mathcal{D} = \left\{(\boldsymbol{x}^{(1)}, y^{(1)}), (\boldsymbol{x}^{(2)}, y^{(2)}), \ldots, (\boldsymbol{x}^{(n)}, y^{(n)})\right\} = \left\{(\boldsymbol{x}^{(i)}, y^{(i)})\right\}_{i=1}^{n} \quad (4)$$

- Each $\boldsymbol{x}$ is a vector of features:

$$\boldsymbol{x} = \begin{bmatrix} x_0 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^{m+1} \quad \text{and} \quad x_0 = 1 \quad \text{and} \quad y \in \{0, 1\} \quad (5)$$

- **How to choose the parameters $\theta$?**

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

# Logistic Regression Cost Function

- Gradient descent is performed in order to find the parameters $\boldsymbol{\theta}$
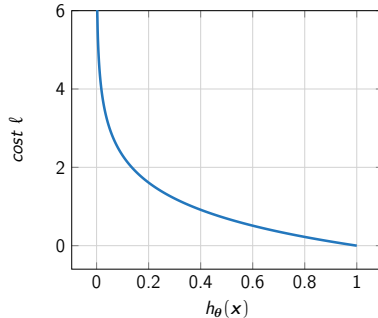- To this end, a cost function is needed:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \ell(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}), y^{(i)}) \tag{6}$$

- The cost function $\ell(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y)$ is defined as follows:
  *(square loss would be **non-convex**...)*

$$\ell(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases} \tag{7}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

# Logistic Regression Cost Function (Ctd.)

$y = 1$:

$y = 0$:

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

# Logistic Regression Cost Function (Ctd.)

- $\ell(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y)$ can be written in a more compact form:

$$\ell(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = -y \log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) - (1 - y) \log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \tag{8}$$

  - If $y = 1$, we get: $-\log(h_{\boldsymbol{\theta}}(\boldsymbol{x}))$
  - If $y = 0$, we get: $-\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}))$

- This gives the **(binary) cross entropy** cost function $\mathcal{J}(\boldsymbol{\theta})$:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \left[ -y^{(i)} \log(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})) \right] \tag{9}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

## Derivation of Cross Entropy

- The likelihood function can be written in the form:

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^{n} h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})^{y^{(i)}} \cdot (1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}))^{1-y^{(i)}} \quad (10)$$

- The cost function is then given by the **negative log-likelihood**:

$$\mathcal{J}(\boldsymbol{\theta}) = -\log \mathcal{L}(\boldsymbol{\theta}) \quad (11)$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

# Derivative of the Sigmoid Function

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

$$\frac{\mathrm{d}}{\mathrm{d}z}g(z) = \frac{0 \cdot (1 + e^{-z}) - (-e^{-z})}{(1 + e^{-z})^2}$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{(1 - 1) + e^{-z}}{(1 + e^{-z})^2} = \frac{1 + e^{-z}}{(1 + e^{-z})^2} - \frac{1}{(1 + e^{-z})^2}$$

$$= \frac{1}{1 + e^{-z}} \left[ 1 - \frac{1}{1 + e^{-z}} \right]$$

$$= \boxed{g(z)(1 - g(z))}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

# Derivation of the Gradient based on a single Example $(\boldsymbol{x}, y)$

$$\frac{\partial}{\partial \theta_j} \mathcal{J}(\boldsymbol{\theta}) = -\frac{\partial}{\partial \theta_j} y \log(g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})) - \frac{\partial}{\partial \theta_j}(1-y)\log(1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})) \qquad \textbf{(derivative of sum terms)}$$

$$= \left[ -\frac{y}{g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})} + \frac{1-y}{1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})} \right] \frac{\partial}{\partial \theta_j} g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}) \qquad \textbf{(derivative of log function)}$$

$$= \left[ -\frac{y}{g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})} + \frac{1-y}{1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})} \right] g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})(1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}))\frac{\partial}{\partial \theta_j} \boldsymbol{\theta}^\mathsf{T} \boldsymbol{x} \qquad \textbf{(chain rule)}$$

$$= \left[ \frac{g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}) - y}{g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})(1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}))} \right] g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})(1 - g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}))x_j \qquad \textbf{(algebraic manipulation)}$$

$$= (g(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x}) - y)x_j = \boxed{(h_{\boldsymbol{\theta}}(\boldsymbol{x}) - y)x_j} \qquad \textbf{(cancelling terms)}$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
**Model Training**
Decision Boundary

## Gradient Descent

- The goal is to minimize $\mathcal{J}(\boldsymbol{\theta})$: $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta})$

- Repeat until convergence {
  $$\boldsymbol{\theta}^{(t+1)} \longleftarrow \boldsymbol{\theta}^{(t)} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}^{(t)}) \quad \textit{// simultaneously update all } \theta_j$$
  }

- The gradient $\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta})$ is given by:

$$\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) \mathbf{x}^{(i)} \tag{12}$$

**Algorithm looks identical to linear regression, but $h_{\theta}(x)$ is different!**

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
**Decision Boundary**

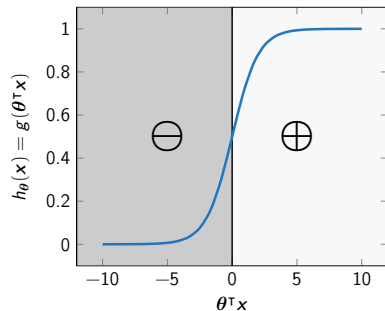## Decision Boundary

- **We have to set a threshold**
- Setting the threshold to 0.5 means:
  - Predict the positive class, if

    $$h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geqslant 0.5 \Leftrightarrow \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x} \geqslant 0$$

  - Predict the negative class, if

    $$h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5 \Leftrightarrow \boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x} < 0$$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
**Decision Boundary**

# Decision Boundary (Ctd.)

- Suppose we have the following hypothesis:
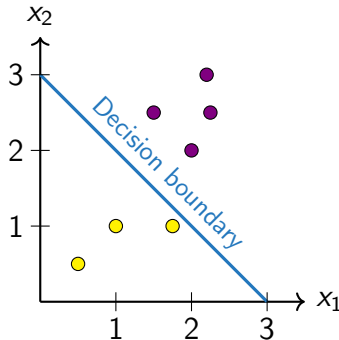
$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

- Using gradient descent we obtained the following coefficients:

$$\theta_0 = -3 \qquad \theta_1 = 1 \qquad \theta_2 = 1$$

- Predict $y = 1$, if $-3 + x_1 + x_2 \geqslant 0$

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
**Decision Boundary**
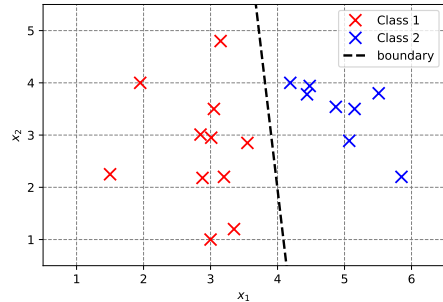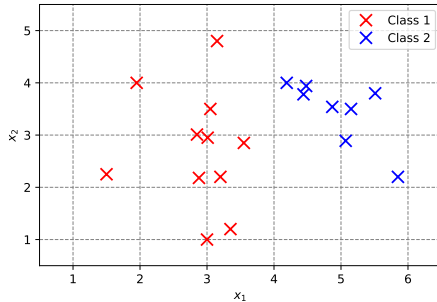
# Decision Boundary (Ctd.)



- Predict $y = 1$, if $-3 + x_1 + x_2 \geqslant 0$
- The decision boundary satisfies $-3 + x_1 + x_2 = 0$
- If $x_2 = 0$, then $x_1 = 3$ and vice versa

**Logistic regression is not a maximum-margin classifier (although the cost function can be adjusted to get that $\Rightarrow$ Hinge loss)**

Introduction
**Model Architecture**
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
**Decision Boundary**

# Example: Decision Boundary



## Where is the sigmoid function?

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

Sigmoid Function
Probabilistic Interpretation
Model Training
Decision Boundary

# Example: Logistic Function

Section:
Non-linear Data

Introduction
Model Architecture
**Non-linear Data**
Multi-Class Classification
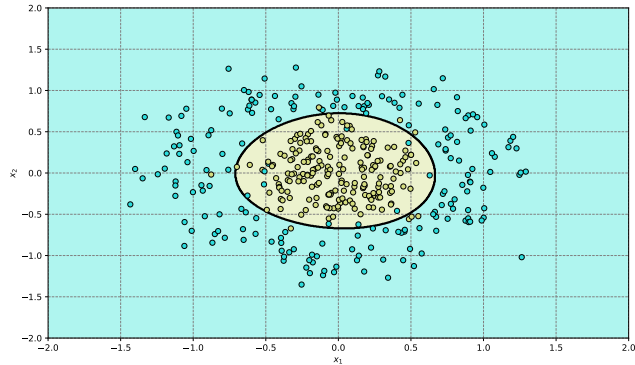Wrap-Up

Feature Mapping
Regularization

# Non-Linear Decision Boundaries

- **Feature mapping** can be used to obtain non-linear decision boundaries
- **Example:**
  - Imagine a circular data set
  - Using the features...

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

  - ...the algorithm could e.g. choose: $\boldsymbol{\theta} = \begin{bmatrix} -1, 0, 0, 1, 1 \end{bmatrix}^{\mathsf{T}}$
  - So we would get: $x_1^2 + x_2^2 = 1 \Rightarrow$ **equation of a unit circle**

Introduction
Model Architecture
**Non-linear Data**
Multi-Class Classification
Wrap-Up

Feature Mapping
Regularization

# Example: Non-Linear Decision Boundary

Introduction
Model Architecture
**Non-linear Data**
Multi-Class Classification
Wrap-Up

Feature Mapping
Regularization

# It is still linear!

**Basis function classification**

Introduction
Model Architecture
**Non-linear Data**
Multi-Class Classification
Wrap-Up

Feature Mapping
Regularization

# Logistic Regression with Regularization

- We should apply regularization for non-linear decision boundaries:

$$\frac{1}{n} \sum_{i=1}^{n} \left[ -y^{(i)} \log(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2 \qquad (13)$$

- The last term prevents the parameters $\theta_j$ from becoming too large
- $\lambda \geqslant 0$ controls the degree of regularization
- This leads to smoother decision boundaries

**Section:**
**Multi-Class Classification**

DHBW
Duale Hochschule
Baden-Württemberg

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

**Multiple Classes**
Multinomial Logistic Regression
One-vs-Rest (OVR)
One-vs-One (OVO)

# Multi-Class Classification

- In its basic form logistic regression can handle two classes only

- **What if there are more than two classes?**

- Two approaches:

  1. Change the algorithm so that it can deal with more classes
     ($\rightarrow$ **Multinomial Logistic Regression** / **Softmax Regression**)

  2. Transform the problem into several binary problems.
     Two common techniques are:
     - **One-vs-Rest (OvR)** $\rightarrow$ One-against-All
     - **One-vs-One (OvO)** $\rightarrow$ Pairwise classification

- Let's examine these approaches a bit closer

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
**Multinomial Logistic Regression**
One-vs-Rest (OVR)
One-vs-One (OVO)

# Multinomial Logistic Regression Introduction

- The logistic regression model has to be changed in order to deal with multiple classes

- The sigmoid function is replaced by the **Softmax** function:

$$\boldsymbol{g} : \mathbb{R}^{\kappa} \rightarrow \mathbb{R}^{\kappa} \qquad \boldsymbol{z} \mapsto \boldsymbol{g}(\boldsymbol{z}) \qquad g_k(\boldsymbol{z}) = \frac{e^{z_k}}{\sum_{n=1}^{\kappa} e^{z_n}} \qquad (14)$$

- $\kappa$ is the number of possible outcomes / classes

- The softmax function returns a **probability distribution over the possible outcomes**, i.e. $\sum_{k=1}^{\kappa} g_k(\boldsymbol{z}) = 1$

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
**Multinomial Logistic Regression**
One-vs-Rest (OVR)
One-vs-One (OVO)

# Multinomial Logistic Regression Introduction (Ctd.)

- $z = \begin{pmatrix} \theta_1^\mathsf{T} x & \theta_2^\mathsf{T} x & \dots & \theta_\kappa^\mathsf{T} x \end{pmatrix}^\mathsf{T}$ is the vector of logits

- This means we learn a separate set of parameters $\theta_k$ for each possible class

- All parameter vectors $\theta_k$ are stacked into a single matrix $\Theta$:

$$\Theta = \begin{pmatrix} | & | & \dots & | \\ \theta_1 & \theta_2 & \dots & \theta_\kappa \\ | & | & \dots & | \end{pmatrix} \tag{15}$$

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
**Multinomial Logistic Regression**
One-vs-Rest (OVR)
One-vs-One (OVO)

# Derivative of the Cross-Entropy Function

$$\mathcal{J}(\boldsymbol{\Theta}) = -\sum_{k=1}^{\kappa} y_k \log(g_k(\boldsymbol{z})) \qquad \text{with } \boldsymbol{z} = \begin{pmatrix} \boldsymbol{\theta}_1^{\mathsf{T}} \boldsymbol{x} \\ \boldsymbol{\theta}_2^{\mathsf{T}} \boldsymbol{x} \\ \vdots \\ \boldsymbol{\theta}_\kappa^{\mathsf{T}} \boldsymbol{x} \end{pmatrix}$$

$$\frac{\partial}{\partial \theta_{ij}} \mathcal{J}(\boldsymbol{\Theta}) = -\sum_{k=1}^{\kappa} y_k \frac{\partial \log(g_k(\boldsymbol{z}))}{\partial g_k(\boldsymbol{z})} \cdot \frac{\partial g_k(\boldsymbol{z})}{\partial z_i} \cdot \frac{\partial z_i}{\partial \theta_{ij}} \qquad \longrightarrow \text{chain rule}$$

$$= -\sum_{k=1}^{\kappa} y_k \frac{1}{g_k(\boldsymbol{z})} \cdot \frac{\partial g_k(\boldsymbol{z})}{\partial z_i} \cdot \frac{\partial z_i}{\partial \theta_{ij}} \qquad \longrightarrow \frac{\mathrm{d}}{\mathrm{d}x} \log(x) = \frac{1}{x}$$

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
**Multinomial Logistic Regression**
One-vs-Rest (OVR)
One-vs-One (OVO)

# Derivative of the Softmax Function

$$g_k(\boldsymbol{z}) = \frac{e^{z_k}}{\sum_{n=1}^{\kappa} e^{z_n}} \qquad\qquad \frac{\partial}{\partial z_i} g_k(\boldsymbol{z}) = \begin{cases} g_i(\boldsymbol{z})(1 - g_i(\boldsymbol{z})) & \text{if } i = k \\ -g_k(\boldsymbol{z}) g_i(\boldsymbol{z}) & \text{if } i \neq k \end{cases}$$

**Case ❶:** $i = k$

$$\frac{\partial}{\partial z_i} g_k(\boldsymbol{z}) = \frac{e^{z_k} \sum_{n=1}^{\kappa} e^{z_n} - e^{z_k} e^{z_i}}{(\sum_{n=1}^{\kappa} e^{z_n})^2}$$

$$= \frac{e^{z_k}}{\sum_{n=1}^{\kappa} e^{z_n}} \left[ 1 - \frac{e^{z_i}}{\sum_{n=1}^{\kappa} e^{z_n}} \right]$$

$$= g_k(\boldsymbol{z})(1 - g_i(\boldsymbol{z}))$$

$$= g_k(\boldsymbol{z})(1 - g_k(\boldsymbol{z}))$$

**Case ❷:** $i \neq k$:

$$\frac{\partial}{\partial z_i} g_k(\boldsymbol{z}) = \frac{0 \cdot \sum_{n=1}^{\kappa} e^{z_n} - e^{z_k} e^{z_i}}{(\sum_{n=1}^{\kappa} e^{z_n})^2}$$

$$= -\frac{e^{z_k}}{\sum_{n=1}^{\kappa} e^{z_n}} \frac{e^{z_i}}{\sum_{n=1}^{\kappa} e^{z_n}}$$

$$= -g_k(\boldsymbol{z}) g_i(\boldsymbol{z})$$

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
**Multinomial Logistic Regression**
One-vs-Rest (OVR)
One-vs-One (OVO)

# Derivative of the Cross-Entropy Function (Ctd.)

$$\frac{\partial}{\partial \theta_{ij}} \mathcal{J}(\boldsymbol{\Theta}) = -\sum_{k=1}^{\kappa} \frac{y_k}{g_k(\boldsymbol{z})} \cdot \frac{\partial g_k(\boldsymbol{z})}{\partial z_i} \cdot \frac{\partial z_i}{\partial \theta_{ij}} \qquad \longrightarrow \text{see slide 33}$$
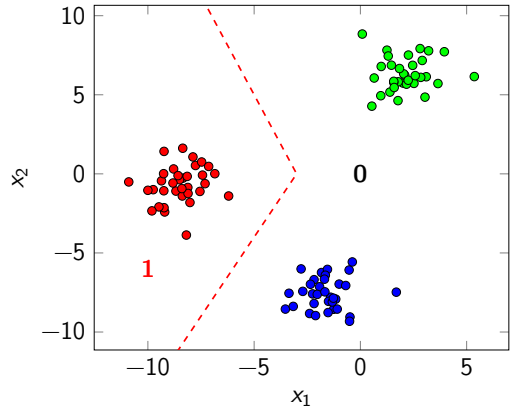
$$= \left[ -\frac{y_k}{g_k(\boldsymbol{z})} g_k(\boldsymbol{z})(1 - g_k(\boldsymbol{z})) + \sum_{\substack{k=1 \\ k \neq i}}^{\kappa} \frac{y_k}{g_k(\boldsymbol{z})} g_k(\boldsymbol{z}) g_i(\boldsymbol{z}) \right] \frac{\partial z_i}{\partial \theta_{ij}} \qquad \longrightarrow \text{separate cases}$$

$$= \left[ -y_k + y_k g_k(\boldsymbol{z}) + \sum_{\substack{k=1 \\ k \neq i}}^{\kappa} y_k g_i(\boldsymbol{z}) \right] \frac{\partial z_i}{\partial \theta_{ij}} = \left[ -y_k + \sum_{k=1}^{\kappa} y_k g_i(\boldsymbol{z}) \right] \frac{\partial z_i}{\partial \theta_{ij}} \qquad \longrightarrow \text{cancel terms}$$
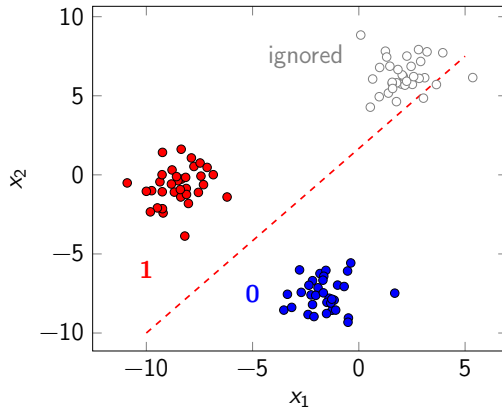
$$= (-y_k + g_k(\boldsymbol{z})) x_j = \boxed{(g_k(\boldsymbol{z}) - y_k) x_j}$$

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
Multinomial Logistic Regression
**One-vs-Rest (OVR)**
One-vs-One (OVO)

# Multi-Class Classification: One-vs-Rest (OVR)

- **Train one classifier per class** (expert for that class)

- We get $|\mathcal{C}|$ classifiers

- The $k$-th classifier learns to distinguish the $k$-th class from all the others

- Set the labels of examples from class $k$ to **1**, all the others to **0**

Introduction
Model Architecture
Non-linear Data
**Multi-Class Classification**
Wrap-Up

Multiple Classes
Multinomial Logistic Regression
One-vs-Rest (OVR)
**One-vs-One (OVO)**

# Multi-Class Classification: One-vs-One (OVO)



- **Train one classifier for each pair of classes**

- We get $\binom{|\mathcal{C}|}{2}$ classifiers

- Ignore all other examples that do not belong to either of the two classes

- **Voting**: Count how often each class wins; the class with the highest score is predicted

Section:

# Wrap-Up

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
**Wrap-Up**

**Summary**
Self-Test Questions
Lecture Outlook
Recommended Literature and further Reading
Meme of the Day

# Summary

- **Logistic regression is used for classification (!!!)**
- It is used for **binary classification problems** (generalizations exist)
- **Output**: Probability of instance belonging to positive class
- Apply a **threshold** to get the classification
- The algorithm minimizes the **cross entropy cost function**
- There is **no closed-form solution** (unlike for linear regression)
- **Basis functions** can be used for non-linear data
- **Multi-class classification**: One-vs-Rest, One-vs-One

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
**Wrap-Up**

Summary
**Self-Test Questions**
Lecture Outlook
Recommended Literature and further Reading
Meme of the Day

## Self-Test Questions

1. Why should you not use linear regression for classification?
2. State the formula for the logistic function.
3. Why do we use cross entropy instead of the squared error?
4. Does logistic regression find the best-separating hyper-plane?
5. What techniques do you know for multi-class classification problems?

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
**Wrap-Up**

Summary
Self-Test Questions
**Lecture Outlook**
Recommended Literature and further Reading
Meme of the Day

# What's next…?

| | |
|---|---|
| **Unit I** | Machine Learning Introduction |
| **Unit II** | Mathematical Foundations |
| **Unit III** | Bayesian Decision Theory |
| **Unit IV** | Probability Density Estimation |
| **Unit V** | Regression |
| **Unit VI** | **Classification I** |
| **Unit VII** | Evaluation |
| **Unit VIII** | Classification II |
| **Unit IX** | Clustering |
| **Unit X** | Dimensionality Reduction |

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
**Wrap-Up**

Summary
Self-Test Questions
Lecture Outlook
**Recommended Literature and further Reading**
Meme of the Day

# Recommended Literature and further Reading I

📕 **[1] Pattern Recognition and Machine Learning**
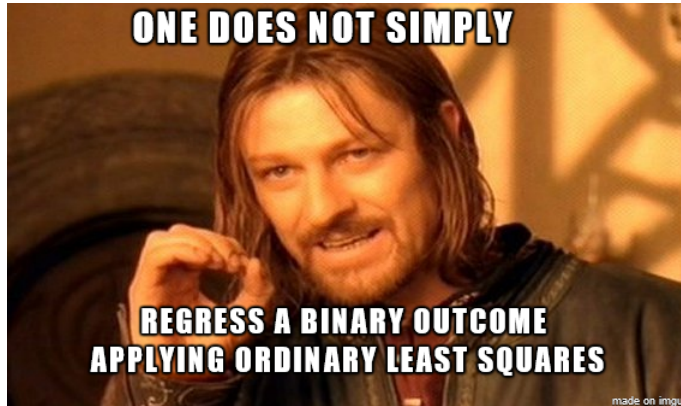*Christopher Bishop. Springer. 2006.*
→ <u>Link</u>, cf. chapter 4.3.2

📕 **[2] Machine Learning: A Probabilistic Perspective**
*Kevin Murphy. MIT Press. 2012.*
→ <u>Link</u>, cf. chapter 8

Introduction
Model Architecture
Non-linear Data
Multi-Class Classification
Wrap-Up

Summary
Self-Test Questions
Lecture Outlook
Recommended Literature and further Reading
Meme of the Day

# Meme of the Day

# Thank you very much for the attention!

**Topic:** *** Applied Machine Learning Fundamentals *** Logistic Regression
**Term:** Winter term 2021/2022

**Contact:**
Daniel Wehner, M.Sc.
SAP SE / DHBW Mannheim
daniel.wehner@sap.com

## Do you have any questions?