# COMPUTATIONAL INTELLIGENCE

(INTRODUCTION TO MACHINE LEARNING) SS16

Lecture 3:

- Classification with Logistic Regression
- Advanced optimization techniques
- Underfitting & Overfitting
- Model selection (Training- & Validation- & Testset)

# CLASSIFICATION WITH LOGISTIC REGRESSION

# Logistic Regression



- **Classification** and not regression
- Classification = recognition

Mogees and vibration classification: https://www.youtube.com/watch?v=xv4hII-_h10
Action recognition: https://www.youtube.com/watch?v=ajswsWVWQvY

# Logistic Regression

- „The" default **classification** model
  - Binary classification
  - Extensions to multi-class later in the course

- Simple classification algorithm
  - **Convex** cost - unique local optimum
  - Gradient descent
  - No more parameter than with linear regre

- Interpretability of parameters

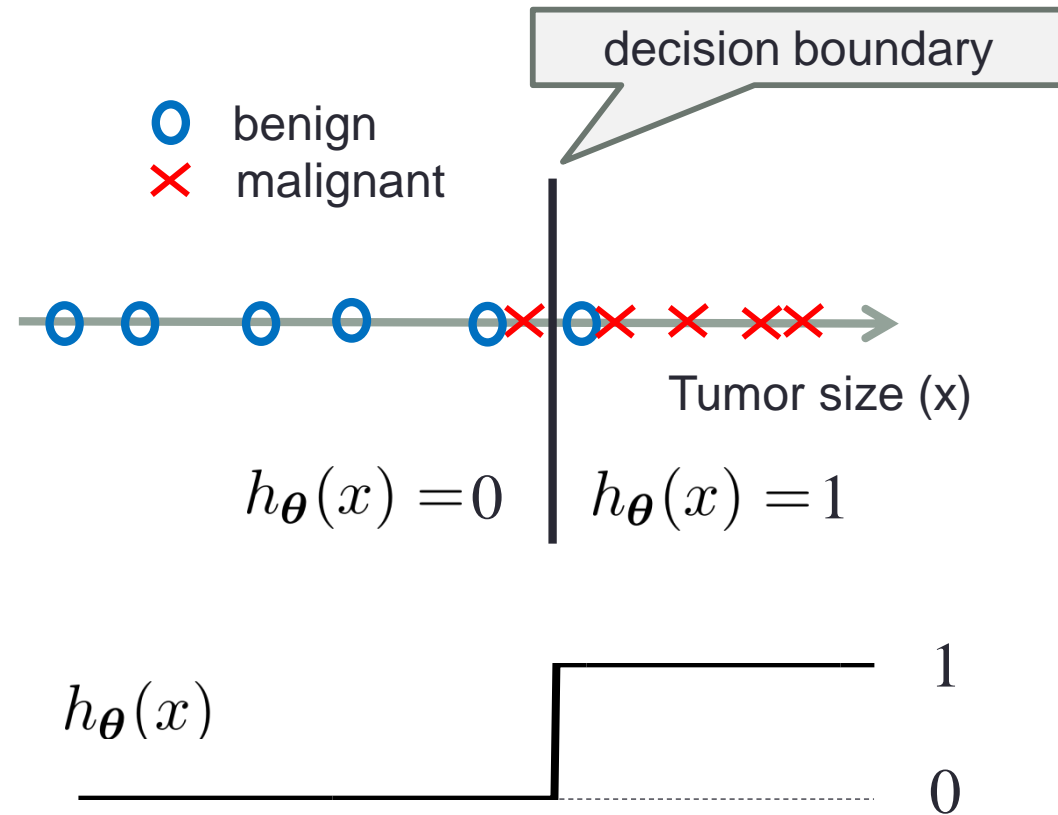- Fast evaluation of hypothesis for making predictions

# LOGISTIC REGRESSION

Hypothesis

# Example (step function hypothesis)

„labeled data"

| $i$ | Tumor size (mm) | Malignant ? |
|-----|-----------------|-------------|
|     | x               | y           |
| 1   | 2.3             | **0 (N)**   |
| 2   | 5.1             | **1 (Y)**   |
| 3   | 1.4             | **0 (N)**   |
| 4   | 6.3             | **1 (Y)**   |
| 5   | 5.3             | **1 (Y)**   |
|     | ...             | ...         |

↑
labels

decision boundary

O benign
✕ malignant

Tumor size (x)

$h_{\boldsymbol{\theta}}(x) = 0$  |  $h_{\boldsymbol{\theta}}(x) = 1$

$h_{\boldsymbol{\theta}}(x)$    1

0

# Example (logistic function hypothesis)

„labeled data"

| $i$ | Tumor size (mm) | Malignant ? |
|-----|-----------------|-------------|
|     | x               | y           |
| 1   | 2.3             | 0 (N)       |
| 2   | 5.1             | 1 (Y)       |
| 3   | 1.4             | 0 (N)       |
| 4   | 6.3             | 1 (Y)       |
| 5   | 5.3             | 1 (Y)       |
|     | …               | …           |

↑
labels

decision boundary

O benign
× malignant

Tumor size (x)

$h_{\boldsymbol{\theta}}(x) < 0.5$  |  $h_{\boldsymbol{\theta}}(x) \geq 0.5$

p=0.2, class 0

$h_{\boldsymbol{\theta}}(x)$

?

1
0.5
0

**Hypothesis**: Tumor is malignant with **probability p**

Classification:  if $p < 0.5$: 0
                      if $p \geq 0.5$: 1

# Example (logistic function hypothesis)

„labeled data"

| $i$ | Tumor size (mm) | Malignant ? |
|-----|-----------------|-------------|
| | x | y |
| 1 | 2.3 | 0 (N) |
| 2 | 5.1 | 1 (Y) |
| 3 | 1.4 | 0 (N) |
| 4 | 6.3 | 1 (Y) |
| 5 | 5.3 | 1 (Y) |
| | … | … |

↑
labels

decision boundary

O benign
× malignant

Tumor size (x)

$h_{\boldsymbol{\theta}}(x) < 0.5$   $h_{\boldsymbol{\theta}}(x) \geq 0.5$

p=0.001, class 0

$h_{\boldsymbol{\theta}}(x)$

1
0.5
0

**Hypothesis**: Tumor is malignant with **probability p**

Classification:  if $p < 0.5$: 0
                 if $p \geq 0.5$: 1

# Logistic (Sigmoid) function

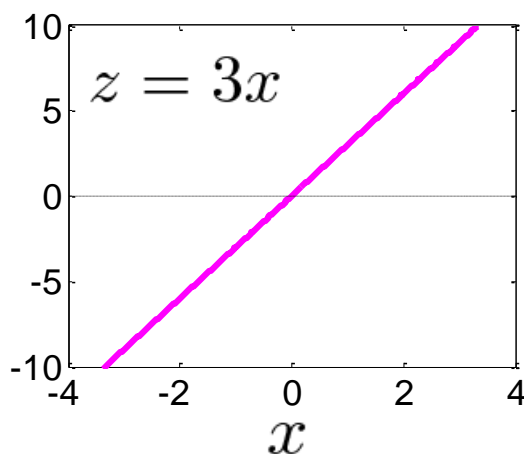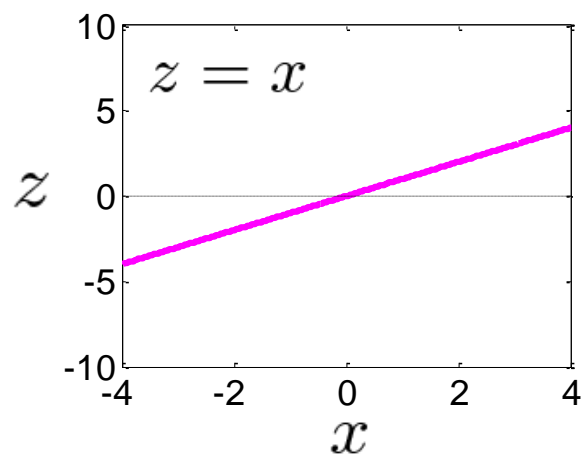$$\sigma(-3) \approx 0.05 \quad \sigma(0) = 0.5 \quad \sigma(3) \approx 0.95$$

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

- Advantages over step function for classification:
  - Differentiable → (**gradient** descent)
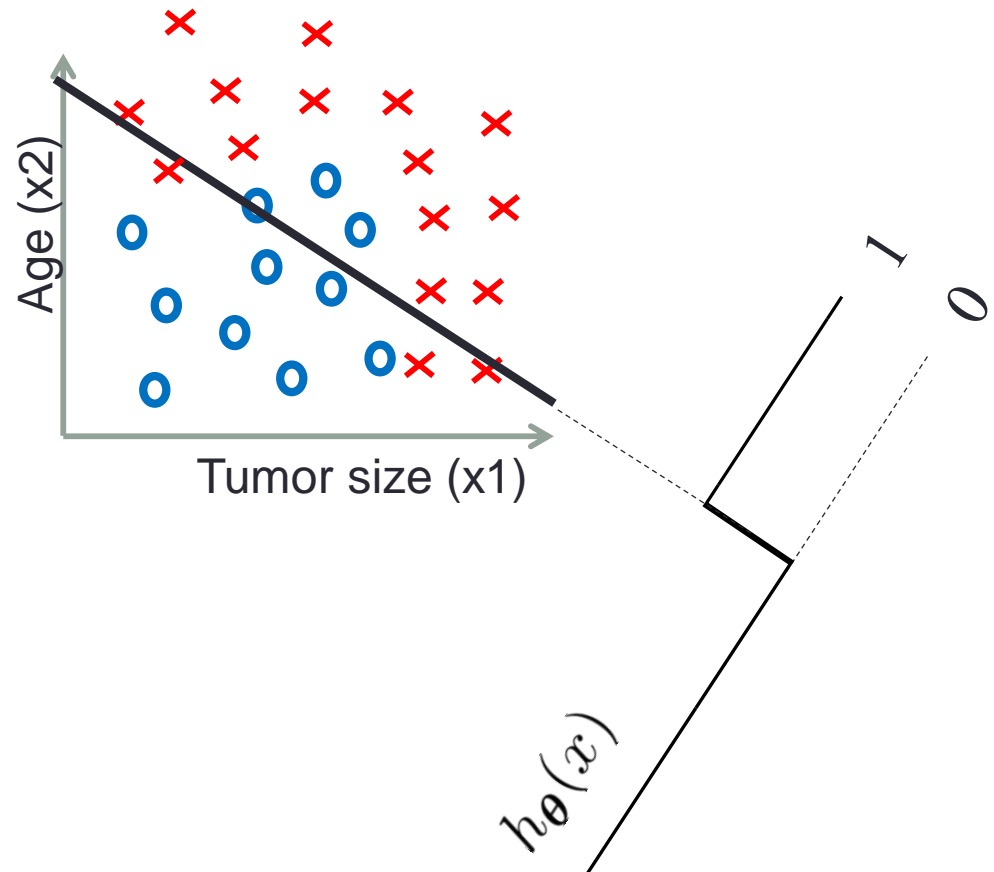  - Contains additional information (how certain is the prediction?)

# Logistic regression hypothesis (one input)

$$h_{\boldsymbol{\theta}}(x) = \sigma(z) = \sigma(\theta_0 + \theta_1 \cdot x)$$
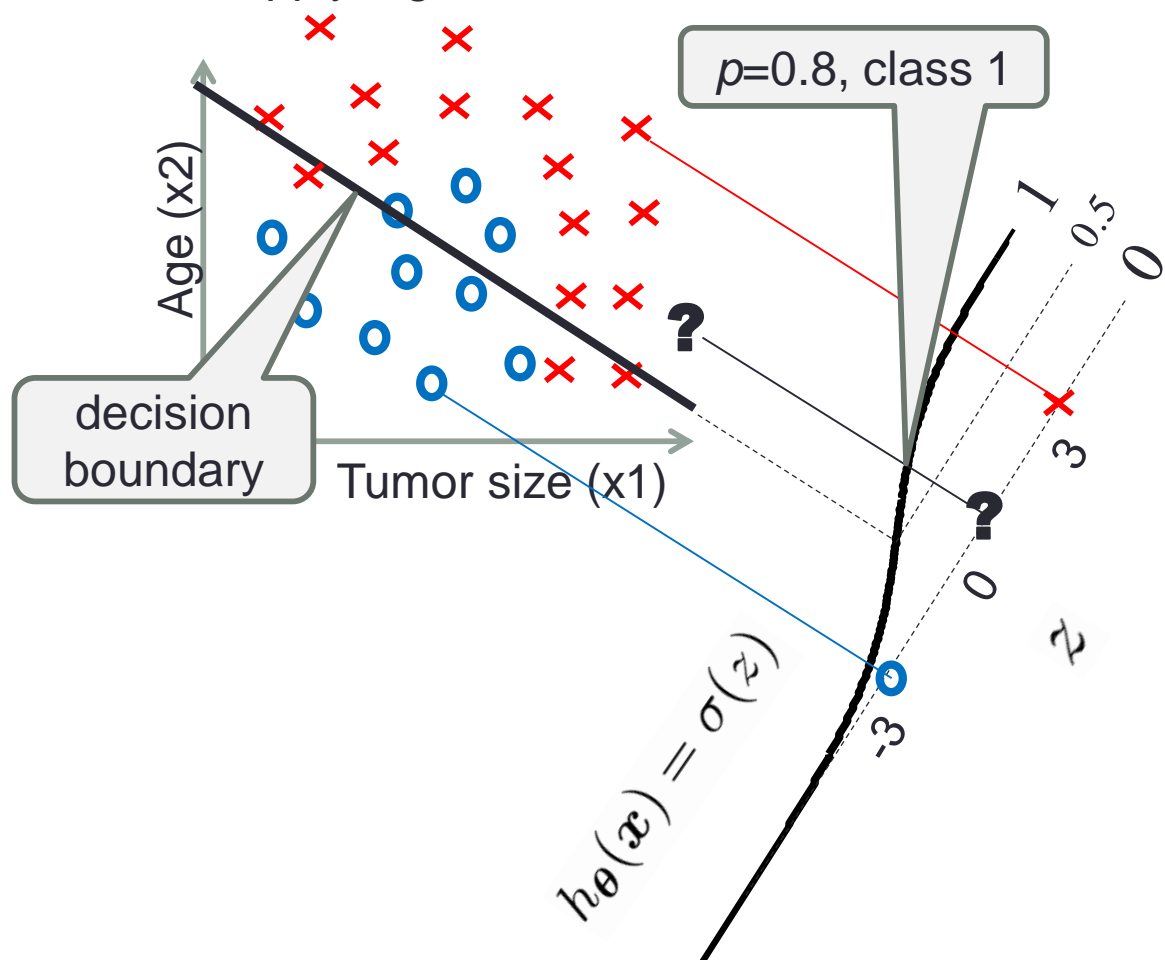
# Classification with multiple inputs

| $i$ | Tumor size (mm) x1 | Age x2 | Malignant? y |
|---|---|---|---|
| 1 | 2.3 | 25 | 0 (N) |
| 2 | 5.1 | 62 | 1 (Y) |
| 3 | 1.4 | 47 | 0 (N) |
| 4 | 6.3 | 39 | 1 (Y) |
| 5 | 5.3 | 72 | 1 (Y) |
| | ... | | ... |

# Multiple inputs and logistic hypothesis

| $i$ | Tumor size (mm) | Age | Malignant? |
|---|---|---|---|
| | x1 | x2 | y |
| 1 | 2.3 | 25 | 0 (N) |
| 2 | 5.1 | 62 | 1 (Y) |
| 3 | 1.4 | 47 | 0 (N) |
| 4 | 6.3 | 39 | 1 (Y) |
| 5 | 5.3 | 72 | 1 (Y) |
| | … | | … |

1. Reduce point in high-dimensional space to a scalar $z$
2. Apply logistic function

p=0.8, class 1

Age (x2)

decision boundary

Tumor size (x1)

$h_\theta(x) = \sigma(z)$

1

0.5

0

3

0

-3

$z$

?

?

# Classification with multiple inputs

| $i$ | Tumor size (mm) | Age | Malignant? |
|---|---|---|---|
| | x1 | x2 | y |
| 1 | 2.3 | 25 | 0 (N) |
| 2 | 5.1 | 62 | 1 (Y) |
| 3 | 1.4 | 47 | 0 (N) |
| 4 | 6.3 | 39 | 1 (Y) |
| 5 | 5.3 | 72 | 1 (Y) |
| | … | | … |

1. Reduce point in high-dimensional space to a scalar $z$
2. Apply logistic function



$p$=0.999, class 1

Age (x2)

Tumor size (x1)

$h_\theta(x) = \sigma(z)$

# Logistic regression hypothesis

1. Reduce high-dimensional input $\boldsymbol{x}$ to a scalar

$$z = \boldsymbol{x}^T \boldsymbol{\theta}$$
$$= \theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n$$

2. Apply logistic function

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T \boldsymbol{\theta})$$
$$= \sigma(\theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n)$$

3. Interpret output $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ as probability and predict class:

$$\text{Class} = \begin{cases} 0 & \text{if } h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5 \\ 1 & \text{if } h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5 \end{cases}$$

# LOGISTIC REGRESSION

Cost function

# Logistic regression cost function

- How well does the hypothesis $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T \boldsymbol{\theta})$ fit the data?



Prediction: 0.98
Actual value y: 1

Prediction: 0.6
Actual value y: 0

Prediction: 0.1
Actual value y: 0

# Logistic regression cost function

- **Probabilistic model**: y is 1 with probability:
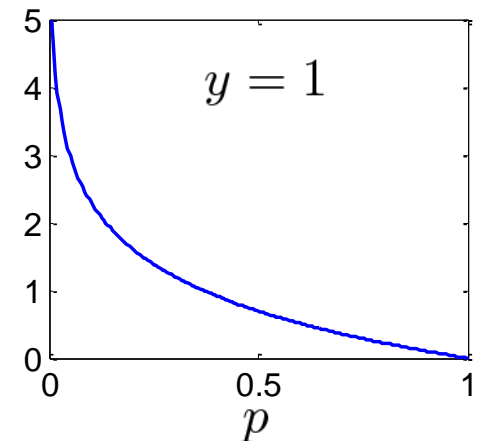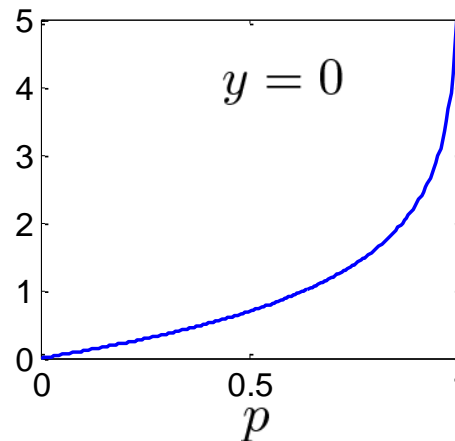
$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T\boldsymbol{\theta})$$

# Logistic regression cost function

- **Probabilistic model**: y is 1 with probability p(x,y=1) = $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T\boldsymbol{\theta})$

The parameters should maximize the **likelihood** of the data

$$\max_{\theta} \log p(X = (x_1 \ldots x_n), y = (y_1, \ldots y_n)|\theta)$$

If data points are independants $p(x_i, y_i, x_j, y_j|\theta) = p(x_j, y_j|\theta).p(x_i, y_i|\theta)$

$$\max_{\theta} \sum_i \log p(x_i, y_i|\theta)$$

Separating positive and negative examples

$$\max_{\theta} \sum_{y_i=1} \log p(x_i, 1|\theta) + \sum_{y_i=0} \log p(x_i, 0|\theta)$$

$\sigma(x^T\theta)$   $1 - \sigma(x^T\theta)$

# Logistic regression cost function

- How well does the hypothesis $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T \boldsymbol{\theta})$ fit the data?
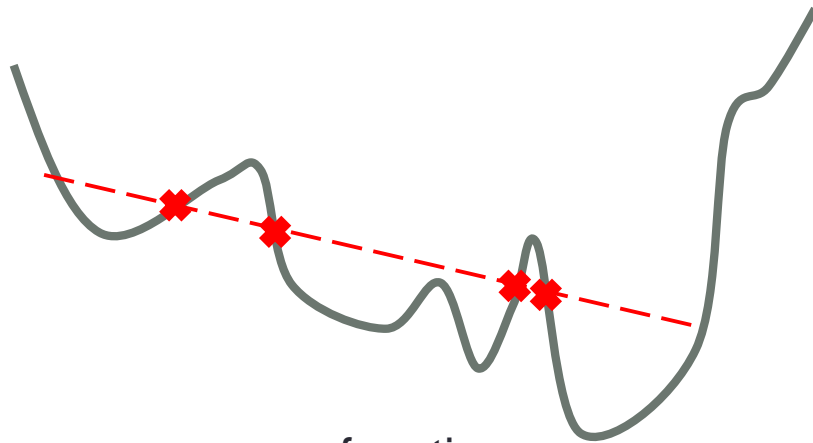- „Cost" for predicting probability $p$ when the real value is $y$:

$$\mathrm{Cost}(p,\ y) = \begin{cases} -\log(1-p) & \text{if } y = 0 \ , \\ -\log(p) & \text{if } y = 1 \ . \end{cases}$$
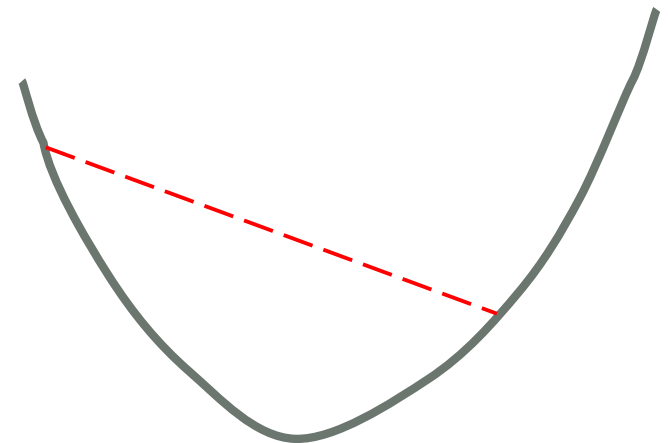


- Mean over all training examples: $\quad J(\boldsymbol{\theta}) = \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} \mathrm{Cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}),\ y^{(i)})$

# Multiple inputs and logistic hypothesis

- How well does the hypothesis $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T \boldsymbol{\theta})$ fit the data?



$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(z)$$

*Prediction: 0.1*
*Actual value y: 0*

*Prediction: 0.6*
*Actual value y: 0*

*Prediction: 0.98*
*Actual value y: 1*

# Comparison cost functions

**Linear regression**



$$\text{Cost}(h,\ y) = (h-y)^2$$

**Logistic regression**



$$\text{Cost}(p,\ y) = \begin{cases} -\log(1-p) & \text{if } y = 0 \ , \\ -\log(p) & \text{if } y = 1 \ . \end{cases}$$

Mean over all training examples:

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}),\ y^{(i)})$$

# Why not mean squared error (MSE) again?

- **MSE** with logistic hypothesis is **non-convex** (many local minima)

- Logistic regression **is convex** (unique minimum)

- Cost function can be derived from statistical principles ("**maximum likelihood**")

non-convex function                                    convex function

# LOGISTIC REGRESSION

Learning from data

# Minimizing the cost via gradient descent

- Gradient descent

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

*(simultaneous update for j=0…n)*



- Gradient of logistic regression cost:

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \left( \underbrace{h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) - y^{(i)}}_{\text{„error“}} \right) \cdot \underbrace{x_j^{(i)}}_{\text{„input“}}$$

(for j=0: $x_0^{(i)} = 1$ )

# Linear to non-linear features



$$x_1 = \text{Tumor Size}, \ x_2 = \text{Age}$$
$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(-10 + 2 \cdot x_1 + 0.05 \cdot x_2)$$

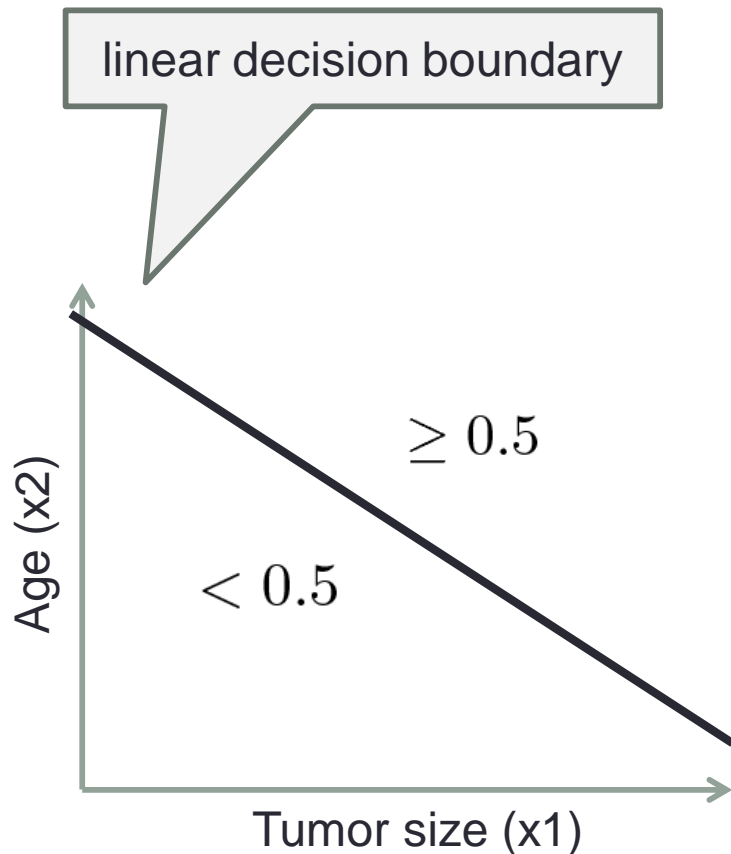# Linear to non-linear features



$$\phi_1 = \text{Tumor Size}, \ \phi_2 = \text{Age}, \ \phi_3 = \text{Tumor Size}^2,$$

$$\phi_4 = \text{Age}^2, \ \phi_5 = \text{Tumor Size} \cdot \text{Age}, \ \dots$$

$$h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \sigma(-3 + 1.2 \cdot \phi_1 + 0.07 \cdot \phi_2 - 0.9 \cdot \phi_3 + \dots)$$

# Decision boundaries

linear decision boundary

non-linear decision boundary

$\geq 0.5$

$< 0.5$

Age (x2)

Tumor size (x1)

$\geq 0.5$

$< 0.5$

Age (x2)

Tumor size (x1)

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(-10 + 2 \cdot x_1 + 0.05 \cdot x_2) \qquad h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \sigma(-3 + 1.2 \cdot \phi_1 + 0.07 \cdot \phi_2 - 0.9 \cdot \phi_3 + \ldots)$$

*Decision boundary is a property of hypothesis, not of data!*

# Linear vs. Logistic Regression

## Linear Regression

- Regression
- Hypothesis $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{\theta}$
- Cost for one training example:

$$\text{Cost}(h,\ y) = (h - y)^2$$

- Gradient

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \frac{2}{m} \sum_{i=1}^{m} \underbrace{\left( h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) - y^{(i)} \right)}_{\text{„error"}} \cdot \underbrace{x_j^{(i)}}_{\text{„input"}}$$

- Analytical:

$$\boldsymbol{\theta}^* = \left( \boldsymbol{X}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}^T \boldsymbol{y}$$

## Logistic Regression

- Binary classification (!)
- Hypothesis $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{x}^T \boldsymbol{\theta})$
- Cost for one training example:

$$\text{Cost}(p,\ y) = \begin{cases} -\log(1-p) & \text{if } y = 0 \ , \\ -\log(p) & \text{if } y = 1 \ . \end{cases}$$

- Gradient

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\left( h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) - y^{(i)} \right)}_{\text{„error"}} \cdot \underbrace{x_j^{(i)}}_{\text{„input"}}$$

- No analytical solution!

# GRADIENT DESCENT TRICKS, AND MORE ADVANCED OPTIMIZATION TECHNIQUES

For linear regression, logistic regression, ….
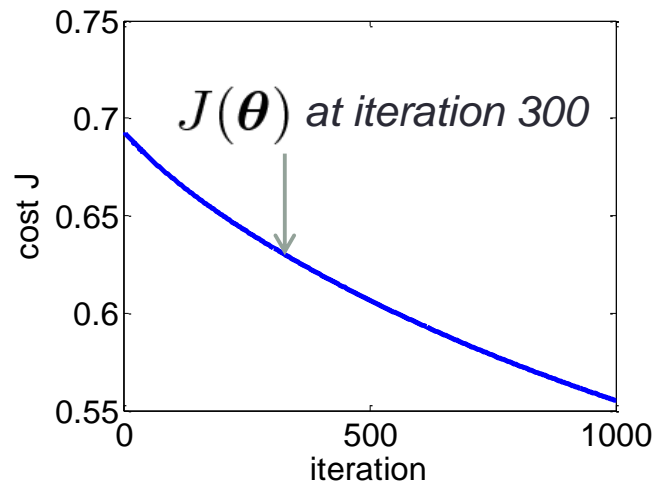
# GD trick #1: feature scaling

- Feature scaling and mean normalization
  - Bring all features into a similar range

  - E.g.: shift and scale each feature to have mean 0 and variance 1

Mean of unscaled feature

$$x_j \leftarrow \frac{x_j - \mu_j}{\sigma_j} \qquad \phi_j \leftarrow \frac{\phi_j - \mu_j}{\sigma_j}$$

(when using non-linear features)

Standard deviation of unscaled feature

  - Do not apply to constant feature $x_0 / \phi_0$ !

- Typically leads to much faster convergence

# GD trick #2: monitoring convergence

- Diagnose typical issues with Gradient Descent:



… slow convergence
(increase learning rate?)

…oscillations
(decrease learning rate)

…divergence
(decrease learning rate)

# GD trick #3: adaptive learning rate

- At each iteration
  - Compare cost function value $J(\boldsymbol{\theta})$ before and after Gradient Descent update

- If cost increased:
  - Reject update (go back to previous parameters)
  - Multiply learning rate $\eta$ by **0.7** (for example)

- If cost decreased:
  - Multiply learning rate $\eta$ by **1.02** (for example)

*Often eliminates slow convergence and divergence issues*
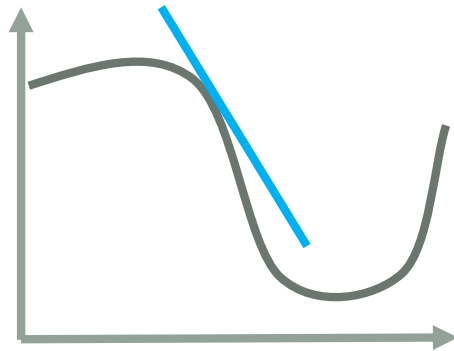
# A black box view of gradient descent

- Write code for computing the cost and its gradient



```
Code to compute J(θ)
```

```
Code to compute   ∂/∂θⱼ J(θ)
```

Gradient descent → $\boldsymbol{\theta}^*$

*Local/global minimum?*

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

Learning rate $\eta$

- Stop when $||\nabla J|| < 10^{-4}$

# More advanced optimization methods

- Gradient methods = order 1 ⟶ Newton methods = order 2



- Need **Hessian matrix** or approximations
- Avoid choosing a learning rate
- Conjugate gradient, BFGS, L-BFGS, …

- Tricky to implement (numerical stability, etc.)
  - Use available toolbox / library implementations!
    `scipy.optimize.minimize`
  - **Only use** when fighting for performance

# EVALUATION OF HYPOTHESIS
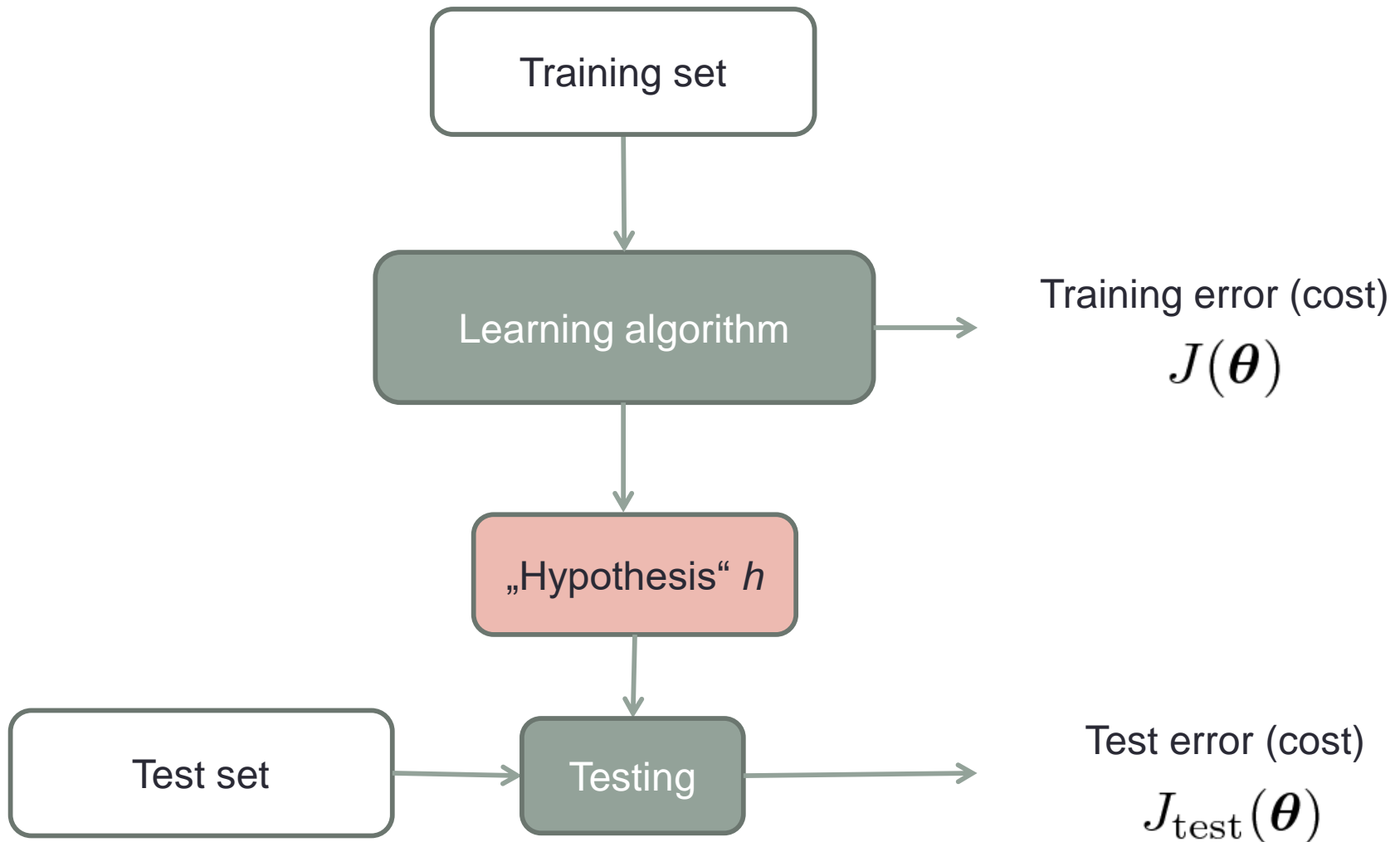
Training and test set

# Training and Test set

- **Training set**: used by learning algorithm to fit parameters and find a hypothesis.

- **Test set**: independent data set, used after learning to estimate the performance of the hypothesis on **new (unseen) test examples.**

*Regression example*



$$\langle \boldsymbol{x}^{(1)}, y^{(1)} \rangle \qquad \langle \boldsymbol{x}_{\text{test}}^{(1)}, y_{\text{test}}^{(1)} \rangle$$

$$\langle \boldsymbol{x}^{(2)}, y^{(2)} \rangle \qquad \langle \boldsymbol{x}_{\text{test}}^{(2)}, y_{\text{test}}^{(2)} \rangle$$

$$\vdots \qquad\qquad \vdots$$

$$\langle \boldsymbol{x}^{(m)}, y^{(m)} \rangle \qquad \langle \boldsymbol{x}_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})} \rangle$$

E.g. 70% randomly chosen examples from dataset are training examples, the remaining 30% are test examples. Must be **disjoint subsets**!

# Training and Test set workflow

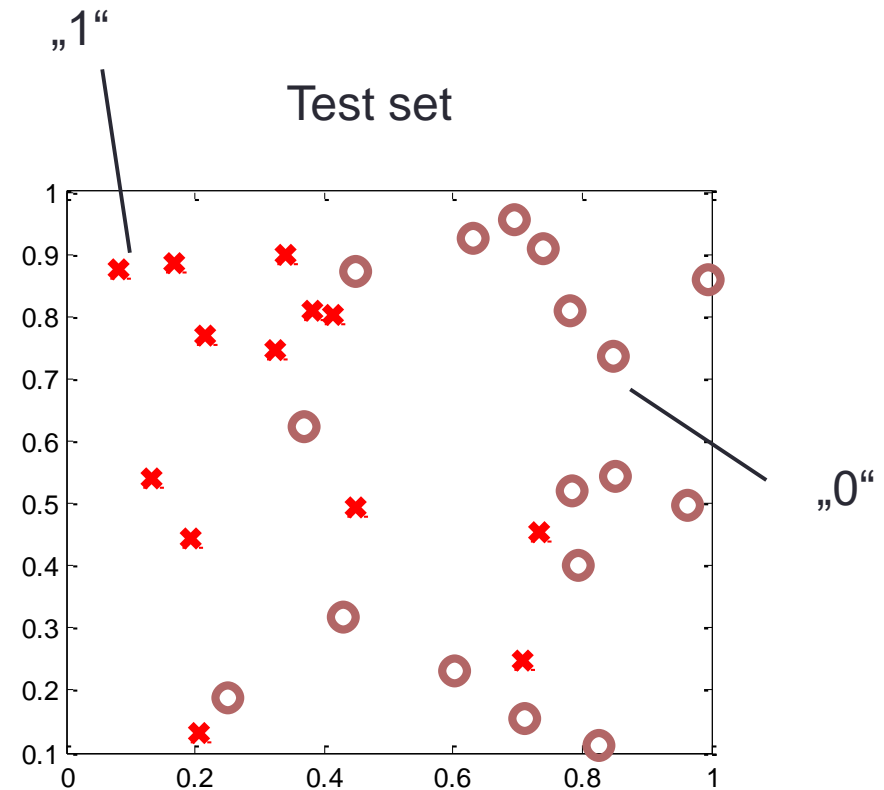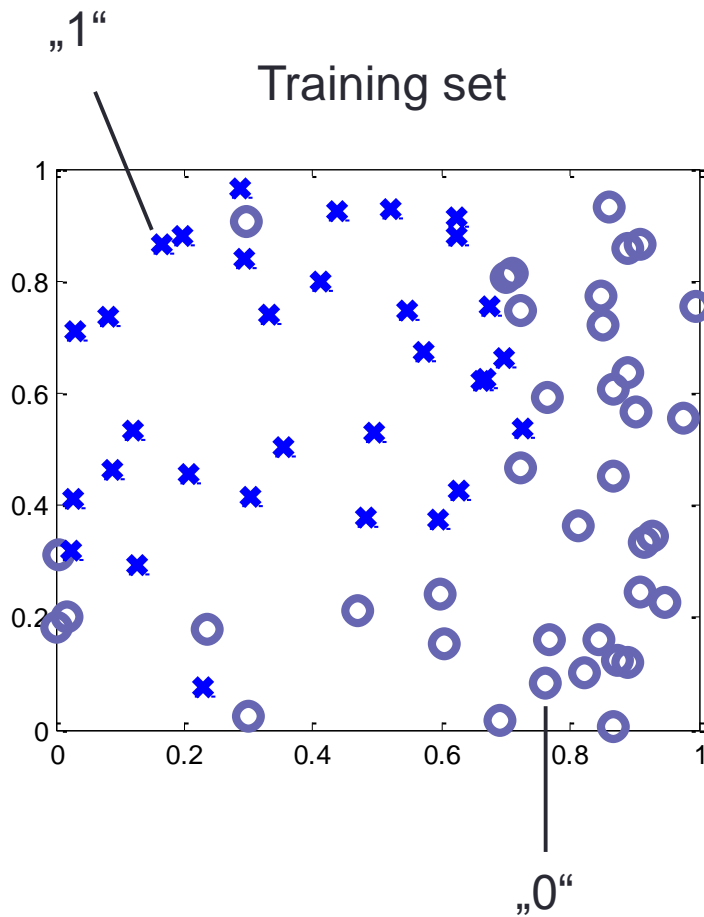# Linear regression training vs. test error



MSEtrain=0.01

MSEtest=2.02

polynomial degree 14

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$
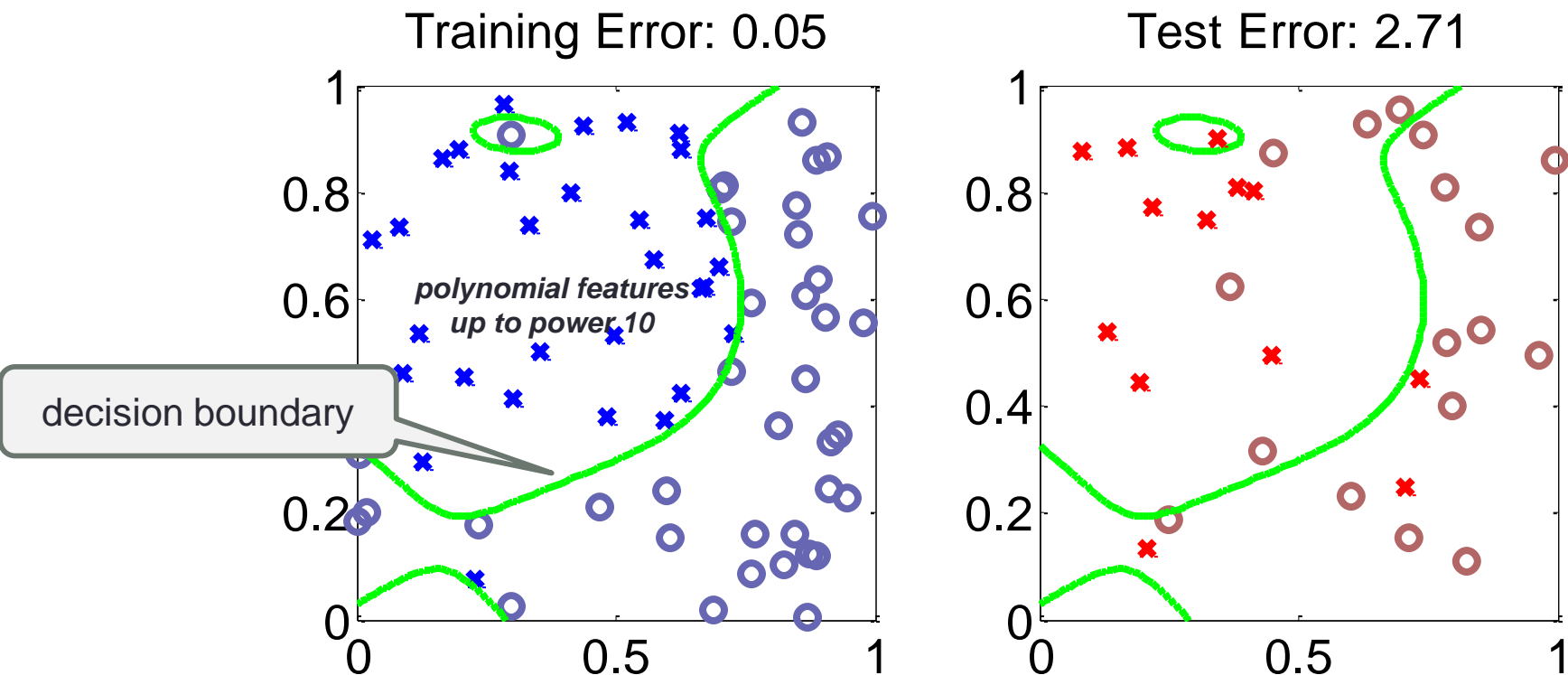
$$J_{\text{test}}(\boldsymbol{\theta}) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}_{\text{test}}^{(i)} \right) - y_{\text{test}}^{(i)} \right)^2$$

# Classification Training / Test set

# Logistic regression training vs. test error



Training Error: 0.05 — Test Error: 2.71

polynomial features up to power 10

decision boundary

$$J(\boldsymbol{\theta}) = -\frac{1}{m}\sum_{i=1}^{m}\mathrm{Cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}),\ y^{(i)})$$

$$J_{\mathrm{test}}(\boldsymbol{\theta}) = -\frac{1}{m_{\mathrm{test}}}\sum_{i=1}^{m_{\mathrm{test}}}\mathrm{Cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}_{\mathrm{test}}^{(i)}),\ y_{\mathrm{test}}^{(i)})$$

# UNDERFITTING AND OVERFITTING

# Polynomial regression under-/overfitting

# Polynomial regression under-/overfitting



MSEtrain=0.01, MSEtest=0.01

Legend:
- training
- test
- polynomial fit degree 5
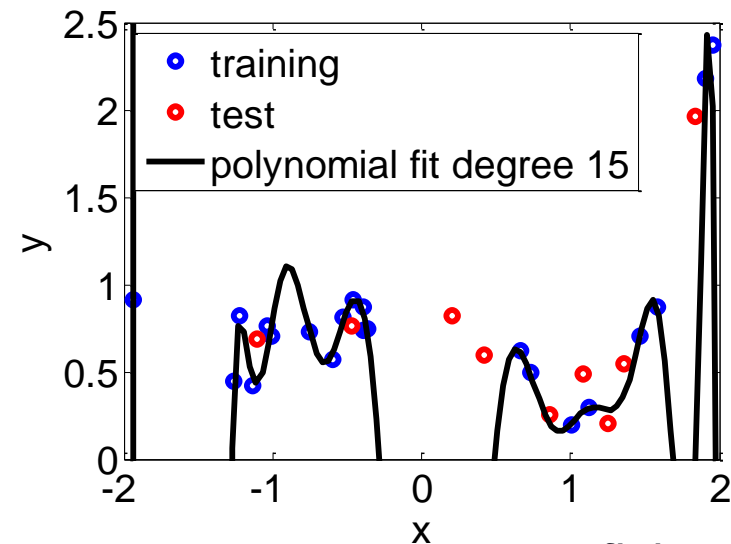
# Polynomial regression under-/overfitting

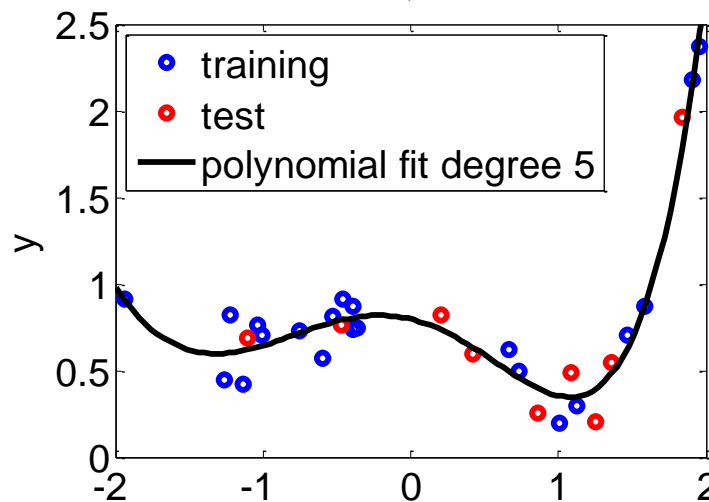# Polynomial regression under-/overfitting



MSEtrain=0.22, MSEtest=0.29
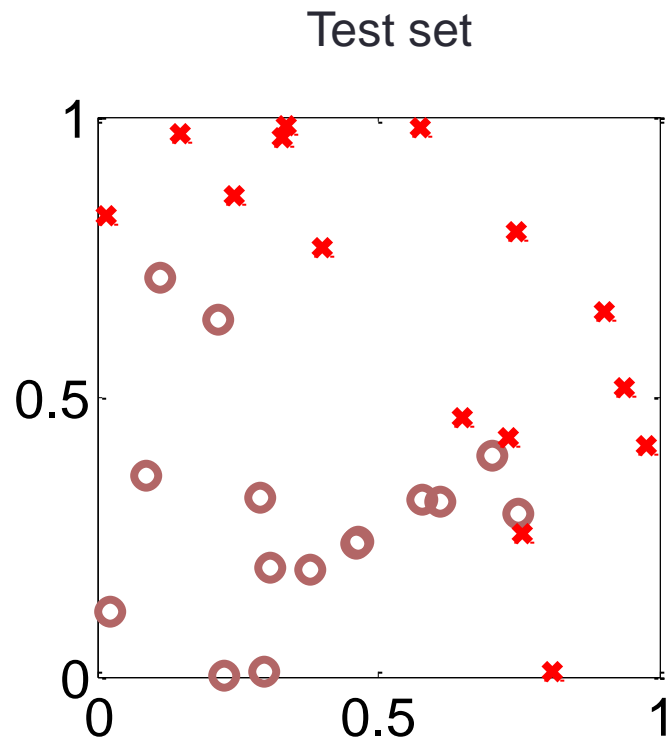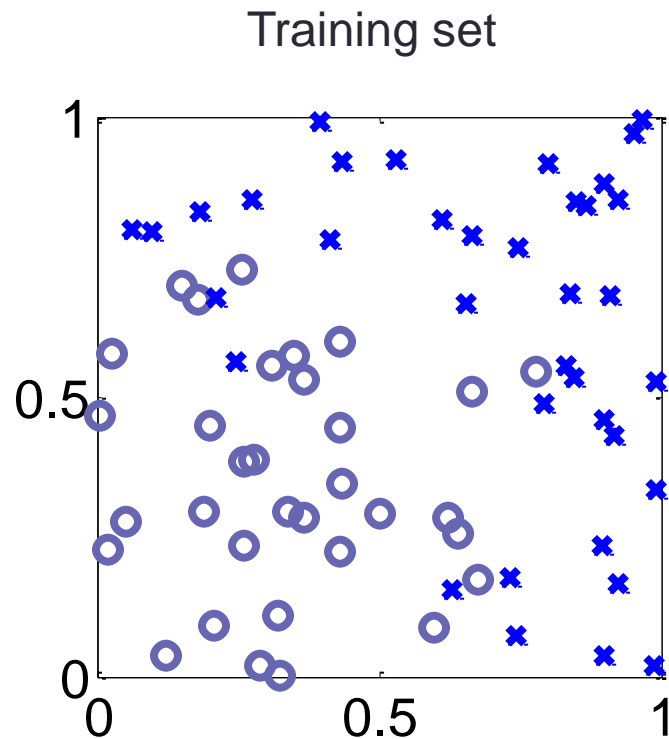
*underfitting*

MSEtrain=0.00, MSEtest=3.17

*overfitting*

MSEtrain=0.01, MSEtest=0.01

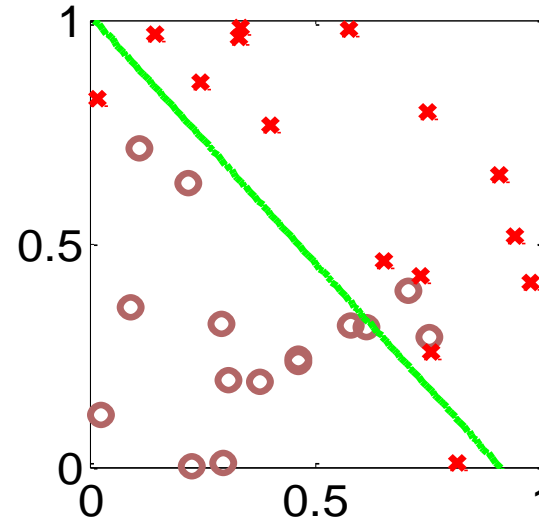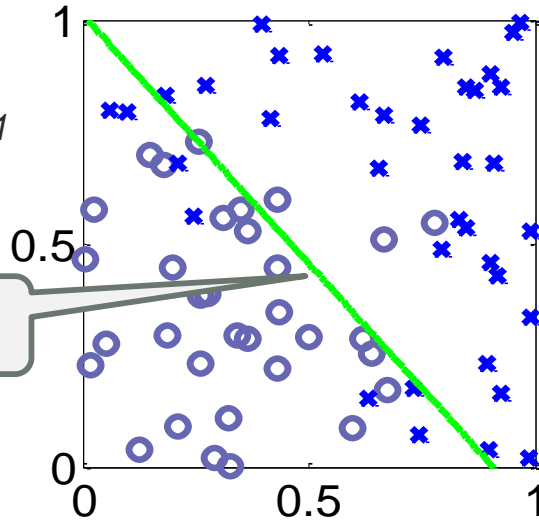*„just right"*

# Logistic regression with polynomial terms



Training set        Test set

# Logistic regression with polynomial terms



Terms up to power 1

Training Error: 0.33

Test Error: 0.32

decision boundary

Terms up to power 2

Training Error: 0.19

Test Error: 0.19

# Logistic regression with polynomial terms

*Terms up to power 20*



Training Error: 0.17

Test Error: 0.44

Training Error: 0.33    Test Error: 0.32

*underfitting*

Training Error: 0.17    Test Error: 0.44

*overfitting*

Training Error: 0.19    Test Error: 0.19

*„just right"*

# Under-/ and Overfitting

Training error (cost)

Test error (cost)

*underfitting*

*overfitting*

*„just right"*

Model complexity
*(e.g. degree of polynomial terms)*

# Under- and Overfitting

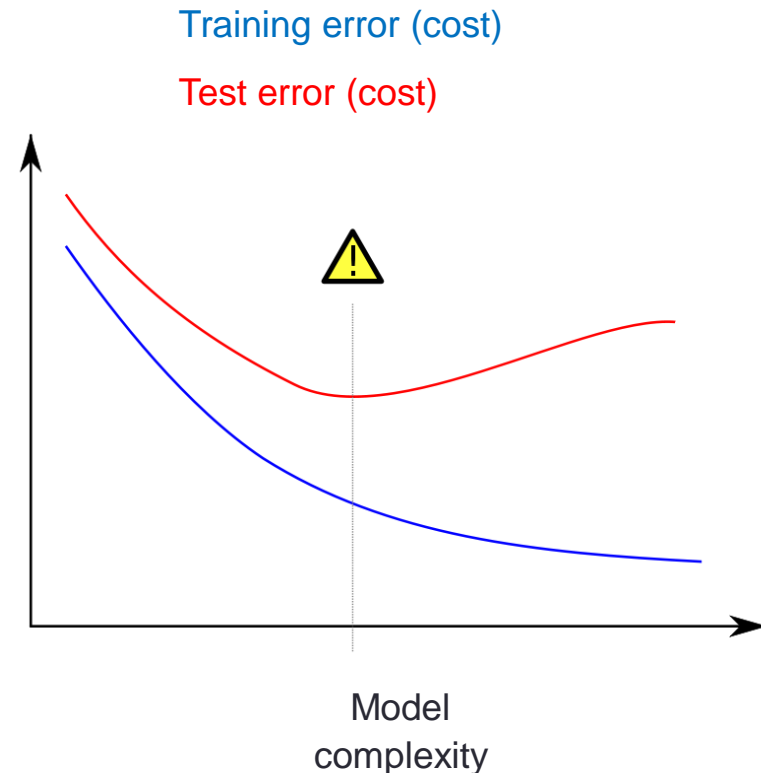- Underfitting:
  - Model is **too simple** (often: too few parameters)
  - **High training error, high test error**

- Overfitting
  - Model is **too complex** (often: too many parameters relative to number of training examples)
  - Low training error, **high test error**

- In between:
  - Model has „right" complexity
  - Moderate training error
  - **Lowest test error**

Training error (cost)

Test error (cost)



Model complexity

# How to deal with overfitting

- Use **model selection** to automatically select the right model complexity

- Use **regularization** to keep parameters small *(other lecture…)*

- Collect more data
  (often not possible or inefficient)

- Manually throw out features which are unlikely to contribute
  (often hard to guess which ones, potentially throwing out the wrong ones)

- Change features vectors, use pre-processing
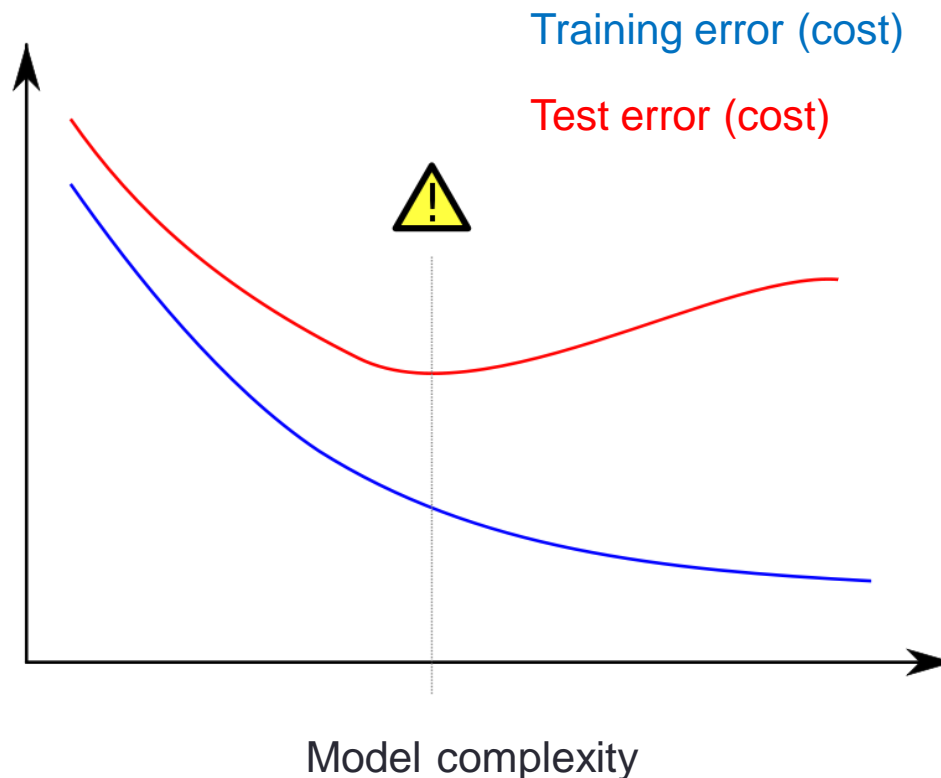  (often not possible or inefficient, time consuming)

# MODEL SELECTION

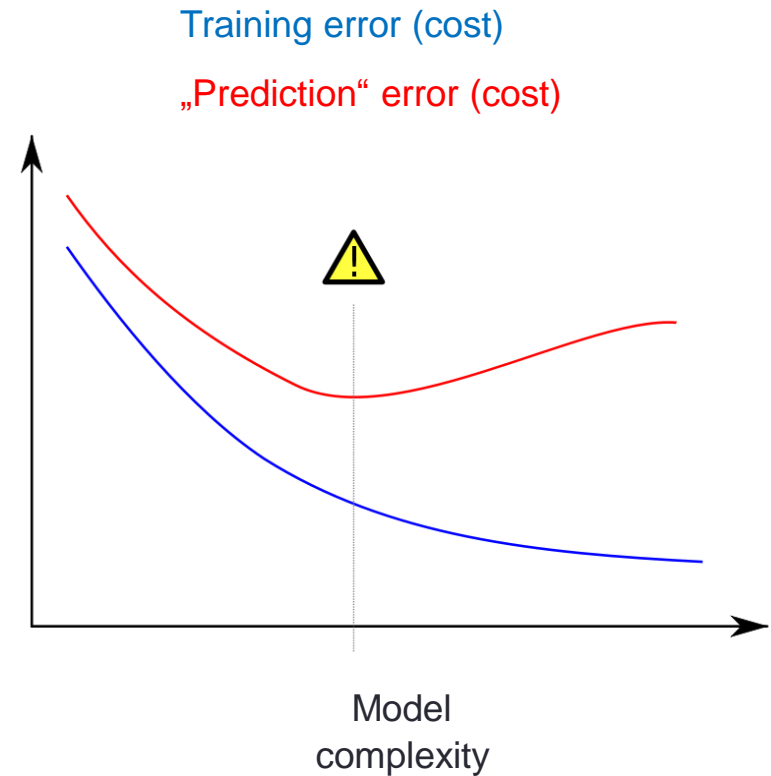Training, Validation and Test sets

# Model selection

- Selection of learning algorithm and „hyperparameters" (model complexity) that are **most suitable** for a given learning problem



Training error (cost)

Test error (cost)

Model complexity

# Idea

- Try out different learning algorithms/variants
  - Vary degree of polynomial
  - Try different sets of features
  - …


- Select variant with best predictive performance

Training error (cost)

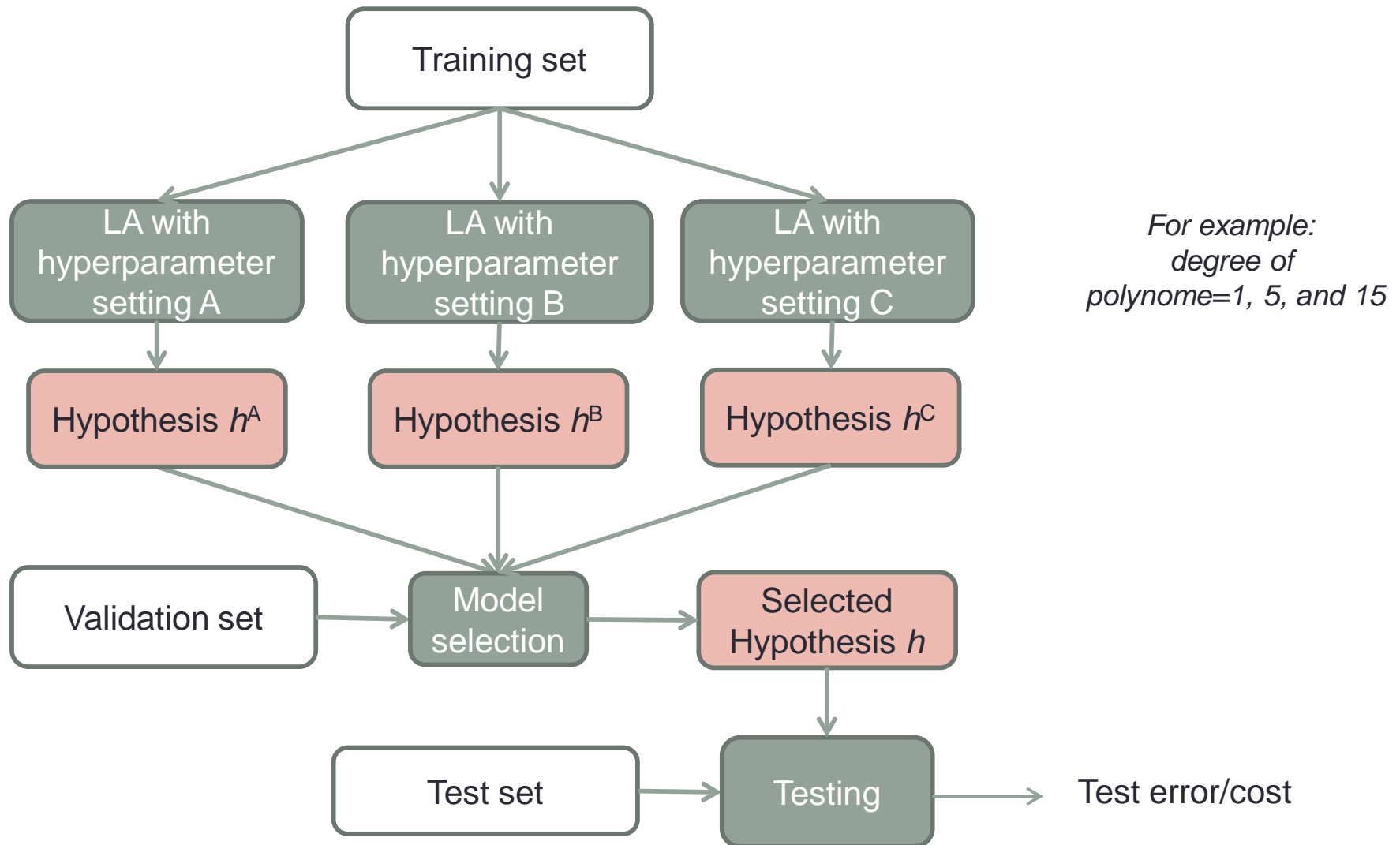„Prediction" error (cost)



Model complexity

# Training, Validation, Test set

- **Training set**: used by learning algorithm **to fit parameters** and find a hypothesis for each learning algorithm/variant.

- **Validation set**: used to estimate predictive performance of each learning algorithm/variant. The hypothesis with **lowest validation error** (cost) is selected.

- **Test set**: independent data set, used after learning and model selection to estimate the performance of the final (selected) hypothesis on **new (unseen) test examples**.
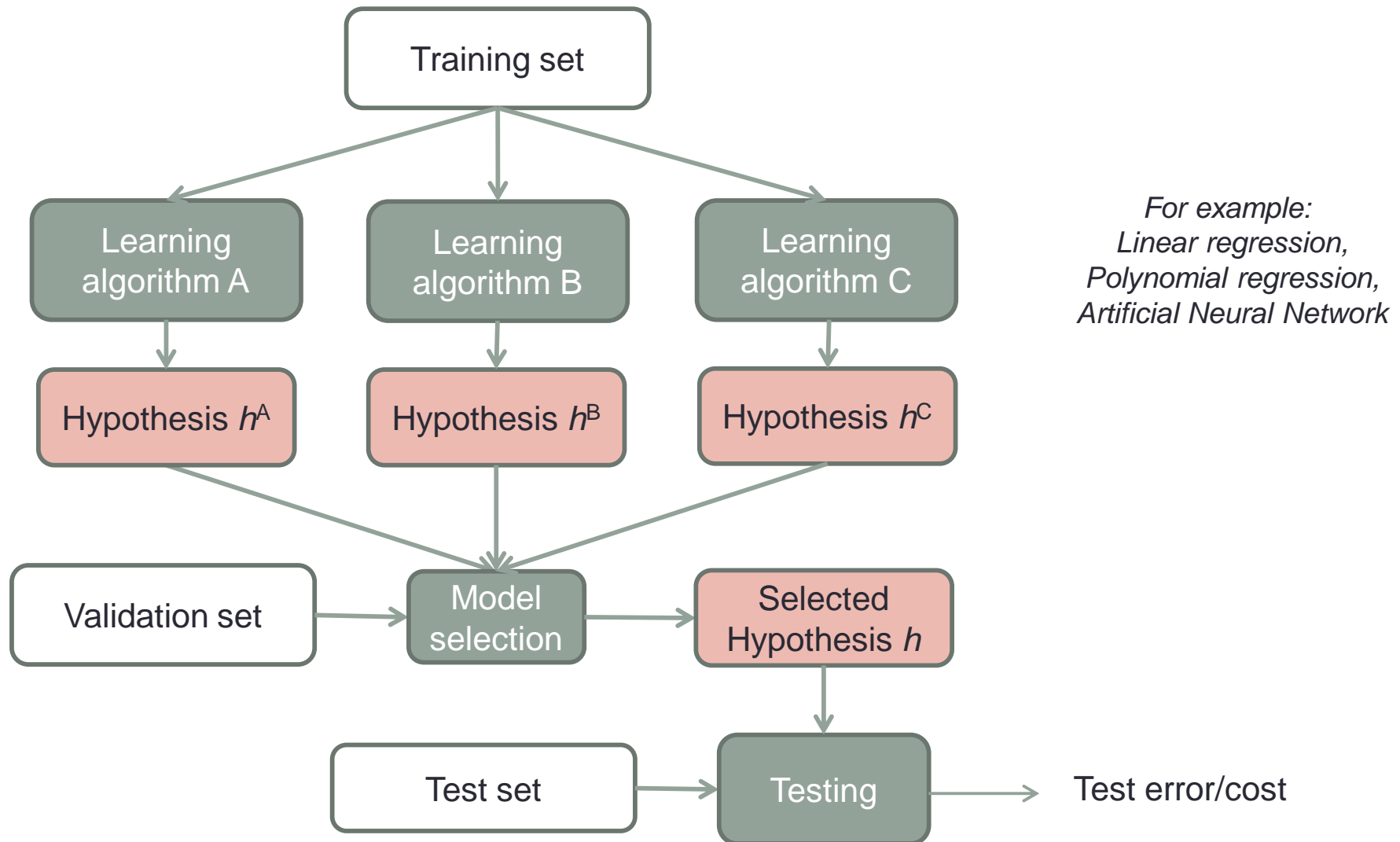
E.g. 60/20/20 % randomly chosen examples from dataset. Must be disjoint subsets!

# Training/Validation/Test set workflow

# Training/Validation/Test set workflow



*For example:
Linear regression,
Polynomial regression,
Artificial Neural Network*

# Some questions…

- Logistic regression is a method for … regression/classification?
- Logistic regression hypothesis?
- What's the cost function used for logistic regression?
- Is it convex or non-convex?
- What does „adaptive learning rate" mean in the context of gradient descent?
- How to evaluate a hypothesis?
- What is under-/overfitting?
- What is model selection?
- What are training, validation and test sets?
- How does model selection work (procedure)?

# What is next?

- Neural Networks:

  - Perceptron

  - Feedforward Neural Network

  - Backpropagation