# Data Science and Machine Learning in Python

Stephan Weyers

Fachhochschule
Dortmund
University of Applied Sciences and Arts

# Topics covered in the online lectures

**Part 1: Data Science**

| | Date | Topics covered |
|---|---|---|
| 1 | Apr 13th | Course introduction<br>Data Science motivation<br>How to use Jupyter Notebook<br>Python types and lists<br>Loops, if/else, functions |
| 2 | Apr 20th | Python tuples, lists, dictionaries<br>Functions<br>Numpy basics, operations<br>Image processing |
| 3 | Apr 27th | Pandas Series, DataFrame<br>Pandas basic operations<br>Import/export files |
| 4 | May 4th | Principles of data visualization<br>Data cleaning and preparation<br>Join, combine and reshape data |
| 5 | May 11th | Volkswohl Bund dataset<br>Data visualization in Python<br>How to write Data Science reports<br>Data aggregation and grouping |

**Part 2: Machine Learning**

| | Date | Topics covered |
|---|---|---|
| 6 | Jun 1st | Introduction to supervised learning<br>Classification and regression<br>scikit-learn<br>k-Nearest Neighbors<br>Linear regression (ridge and lasso) |
| 7 | Jun 8th | Linear classification models<br>Decision trees<br>Random forests and gradient boosting<br>Support vector machines<br>Neural networks |
| 8 | Jun 15th | Introduction to unsupervised learning<br>Preprocessing and scaling<br>Dimensionality reduction<br>Principal component analysis |
| 9 | Jun 22nd | k-means clustering<br>Hierarchical clustering<br>DBSCAN |
| 10 | Jun 29th | Representing data<br>Engineering features |
| 11 | Jul 6th | Model evaluation and improvement<br>Text data analysis |

# Agenda for online lecture 6

| Session | Topic | Mode | Materials used | Minutes | End |
|---|---|---|---|---|---|
| 14:30-16:00 | Organizational questions | Q&A | | 10 | 14:40 |
| | Supermarket exercise | Team work in break-out rooms | Lecture 06a notebook | 40 | 15:20 |
| | k-Nearest Neighbors | Lecture / Q&A | Lecture slides | 15 | 15:35 |
| | Linear regression | Lecture / Q&A | Lecture slides | 20 | 15:55 |
| 16:10-17:40 | Regression toy data | Lecture / Q&A | Lecture 06b notebook | 20 | 16:30 |
| | California housing data | Team work in break-out rooms | Lecture 06c notebook | 45 | 17:15 |
| | Recap W02 / W03 | Lecture / Q&A | | 20 | 17:35 |
| 17:50-19:20 | Teams W04 | Lecture / Q&A | | 10 | 18:00 |
| | Happiness data | Team work in break-out rooms | Lecture 06d notebook | 80 | 19:20 |

# Types of problems

**Supervised Approaches**

- Labeled data
- Target values known

**Classification**

- Predict category

**Regression**

- Predict numeric value

**Unsupervised Approaches**

- Unlabeled data
- No target value provided

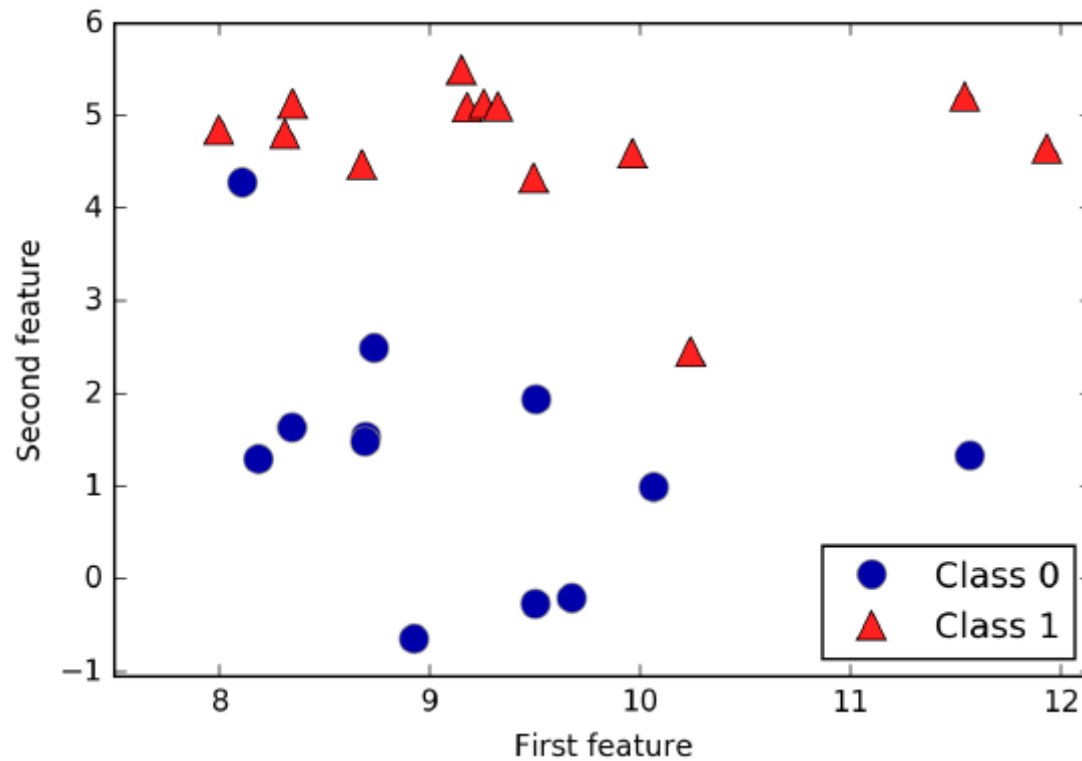**Cluster Analysis**

- Organize similar cases into segments
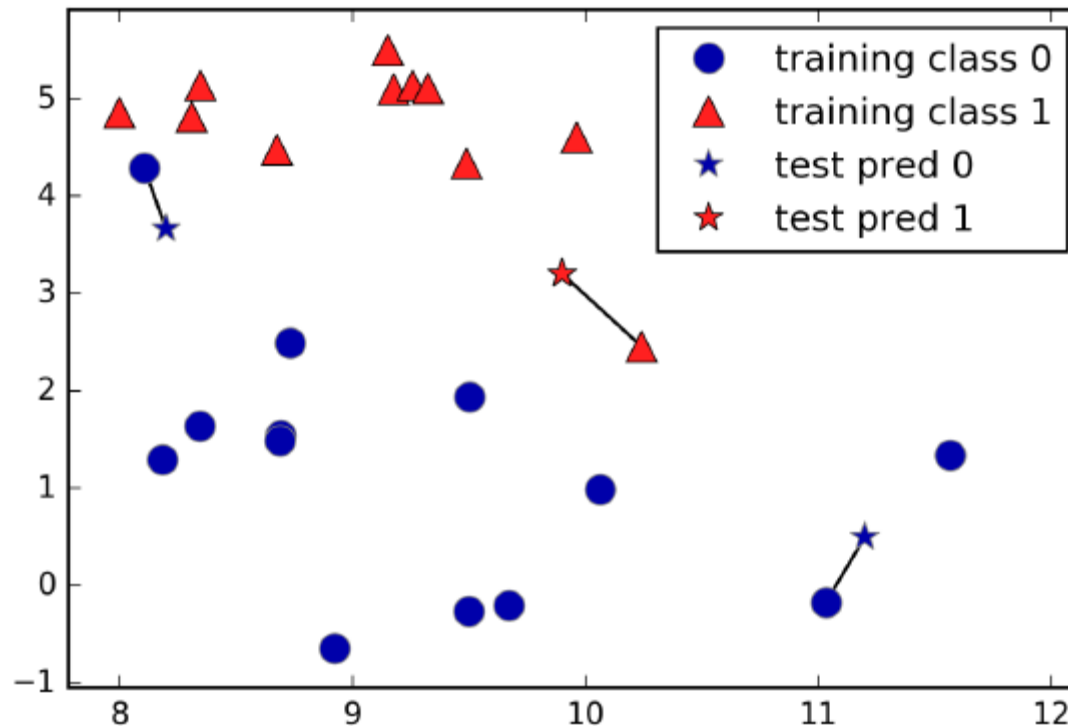
**Dimensionality reduction**

- Reduce number of features

**Question for discussion**
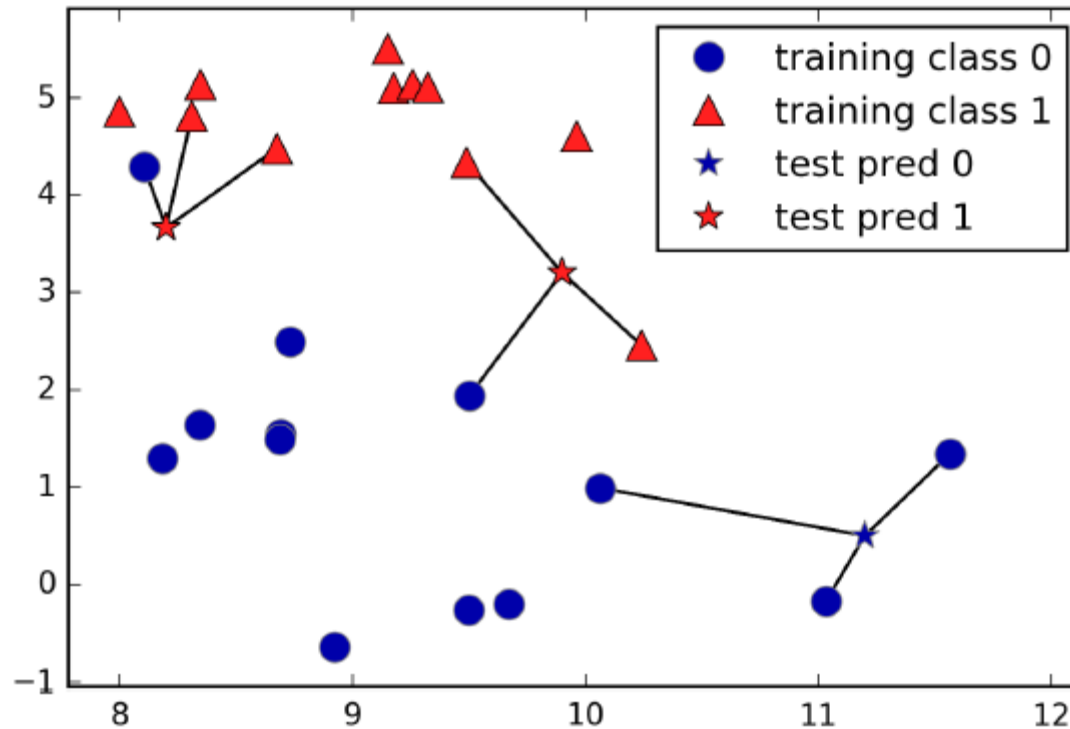
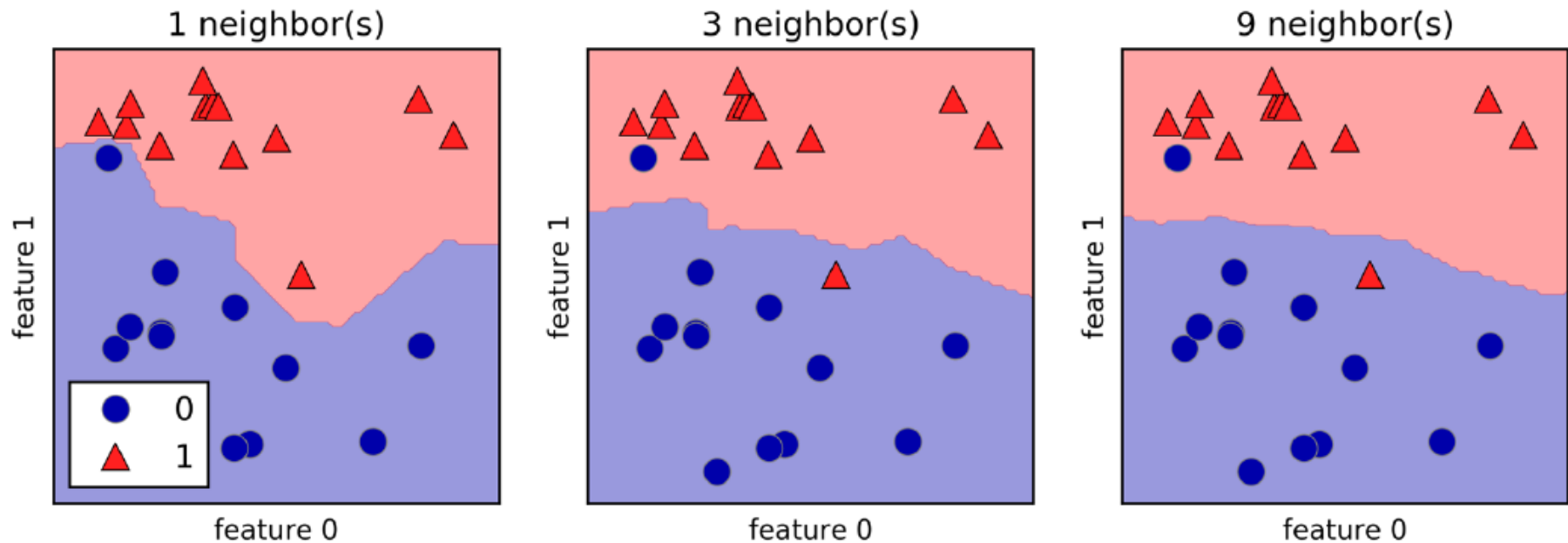- Find examples for each of the 4 categories

# k-Nearest Neighbors – example

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

# k-Nearest Neighbors – Decision Boundary

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

# k-Nearest Neighbors – Accuracy

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

# k-Nearest Neighbors – Summary

**Parameters**

- Number of neighbors k
- Distance metric (Euklidean as default)

**Strengths**

- Easy to understand
- Building model is fast

**Weaknesses**

- Making predictions is very slow on large datasets
- Usually not good with many features (hundreds or more)
- Particularly bad with sparse datasets (many zeros)
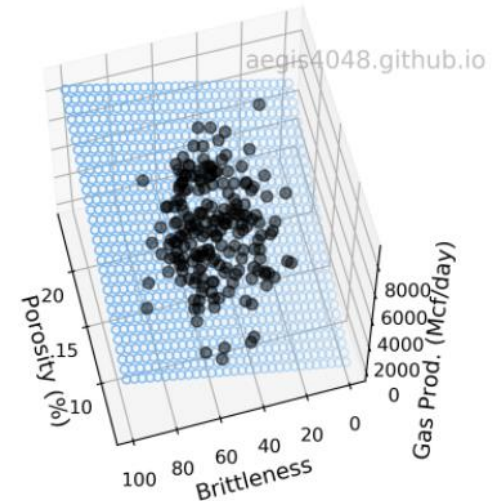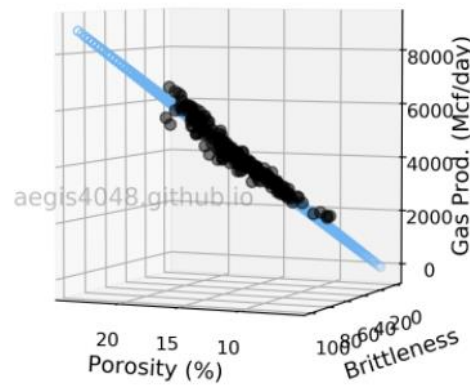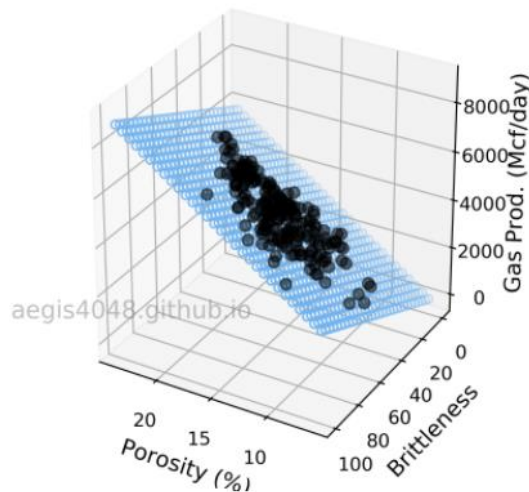- Not robust if features are on different scales

# Linear Regression – one input variable

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

# Linear Regression – two input variables

Source: https://aegis4048.github.io/mutiple_linear_regression_and_visualization_in_python

Given training input data $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$ with labels $y^{(1)}, \dots, y^{(n)} \in \mathbb{R}$.

Try to find $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$ in order to minimize

**Linear regression
(ordinary least squares)**

$$L(w, b) = \sum_{i=1}^{n} \left( y^{(i)} - \left( w \cdot x^{(i)} + b \right) \right)^2$$

**Ridge regression**

$$L_2(w, b) = \sum_{i=1}^{n} \left( y^{(i)} - \left( w \cdot x^{(i)} + b \right) \right)^2 + C\|w\|_2^2$$

**Lasso regression**

$$L_1(w, b) = \sum_{i=1}^{n} \left( y^{(i)} - \left( w \cdot x^{(i)} + b \right) \right)^2 + C\|w\|_1$$
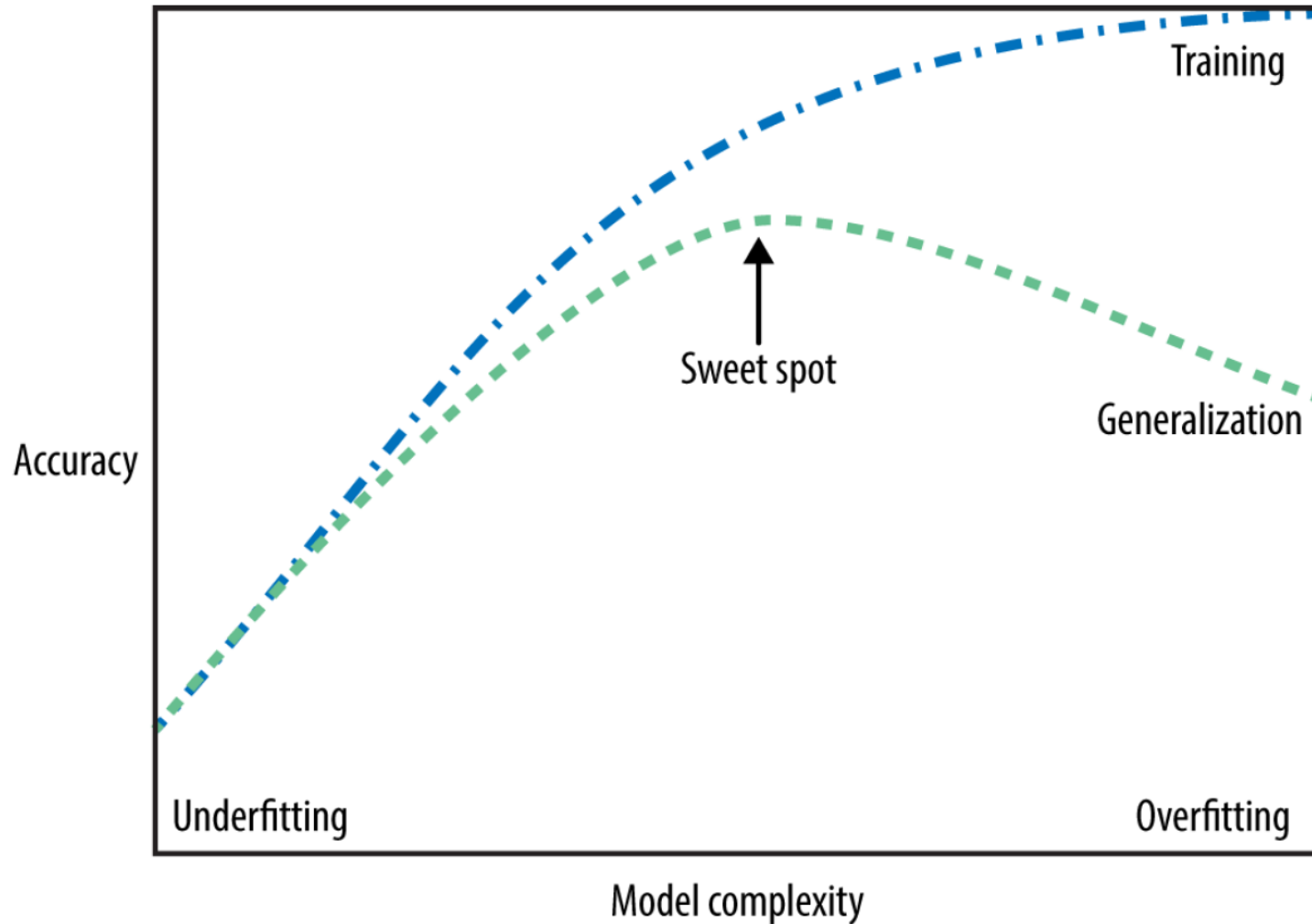
# Linear Models – Summary

**Parameters**

- Regularization parameter alpha and C
- Model type lasso vs. ridge for regression / logistic vs. SVM for classification
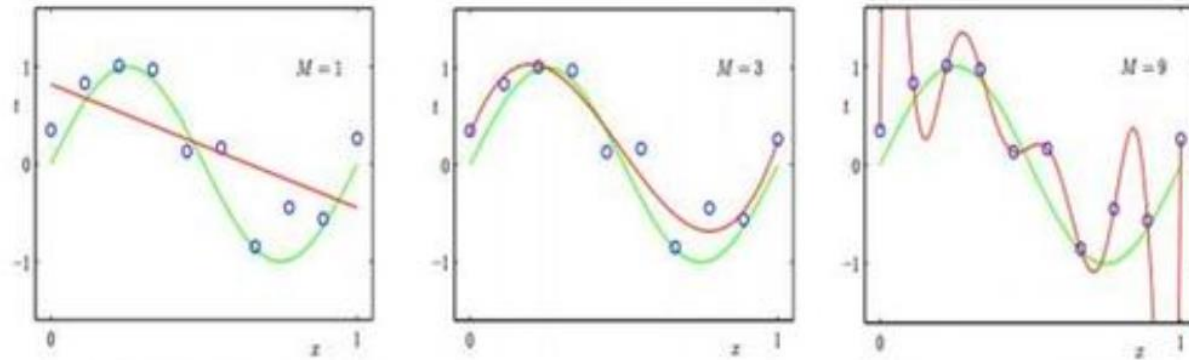
**Strengths**

- Fast to train, fast to predict
- Work well with sparse data
- Relatively easy to understand how predictions are made
- Work well with large number of features
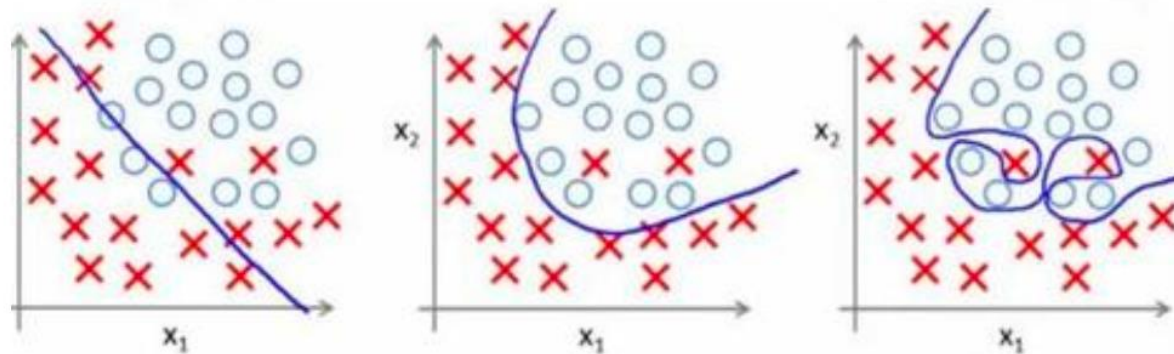
**Weaknesses**

- Coefficients hard to interpret, especially if features are highly correlated
- Sometimes fail with small datasets
- Perform bad with non-linear features and datasets that are not linearly separable

# Machine Learning – Training and Test Error

Source: Müller, A. C., Guido, S. (2016) Introduction to Machine Learning with Python : A Guide for Data Scientists, O'Reilly Media

**Regression**

**Classi-
fication**



Low model complexity
(Underfitting)

High model complexity
(Overfitting)