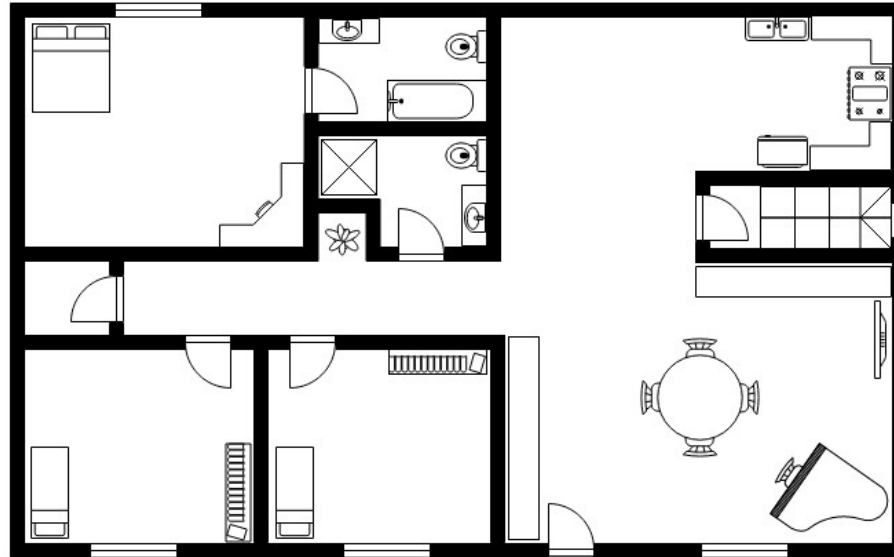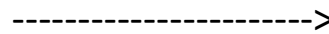# Machine Learning

- Learning from Data
- Models
- Features and Targets
- Dimensions of Machine Learning

## Estimating Appartment Prices



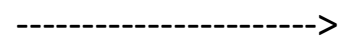Features?                    ----------------------->          Target: Price

Regression: Predicting a real number value

Kind of Iris Plant



Features?                        ---------------------->              Target: Kind
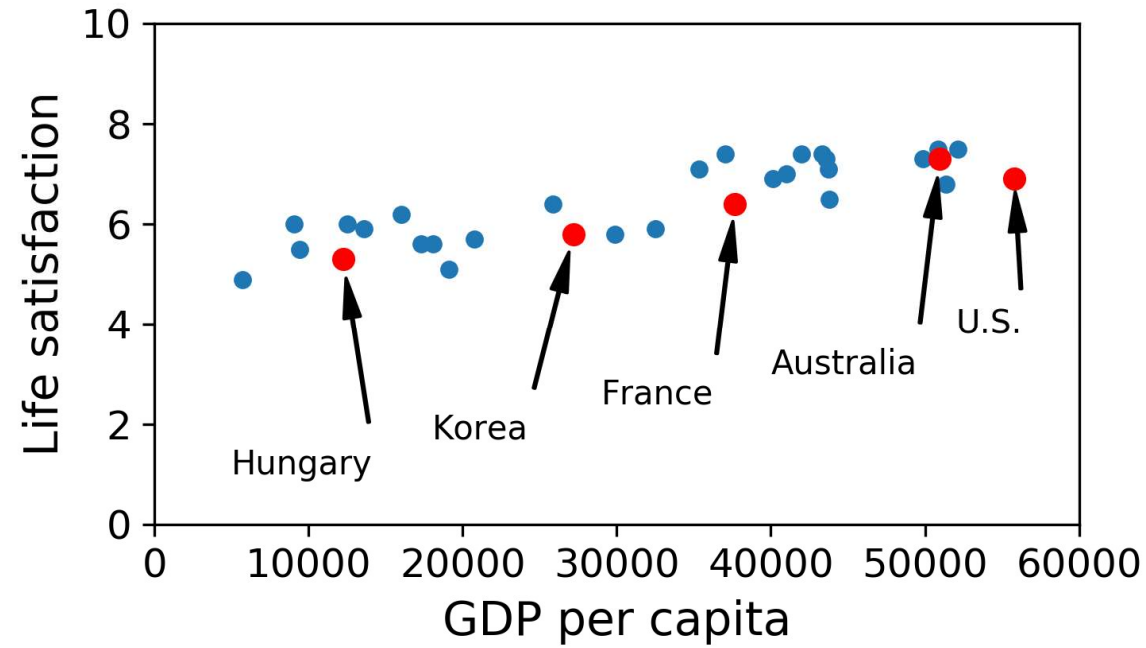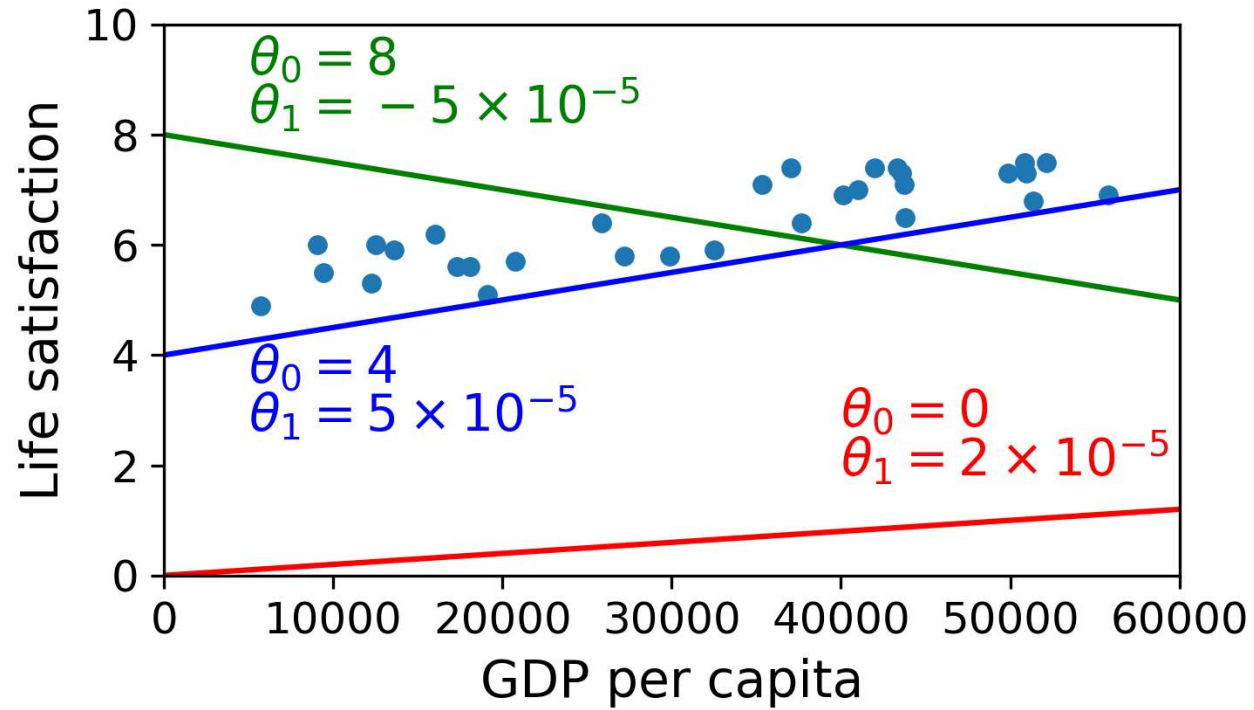
Classification: Predicting a label
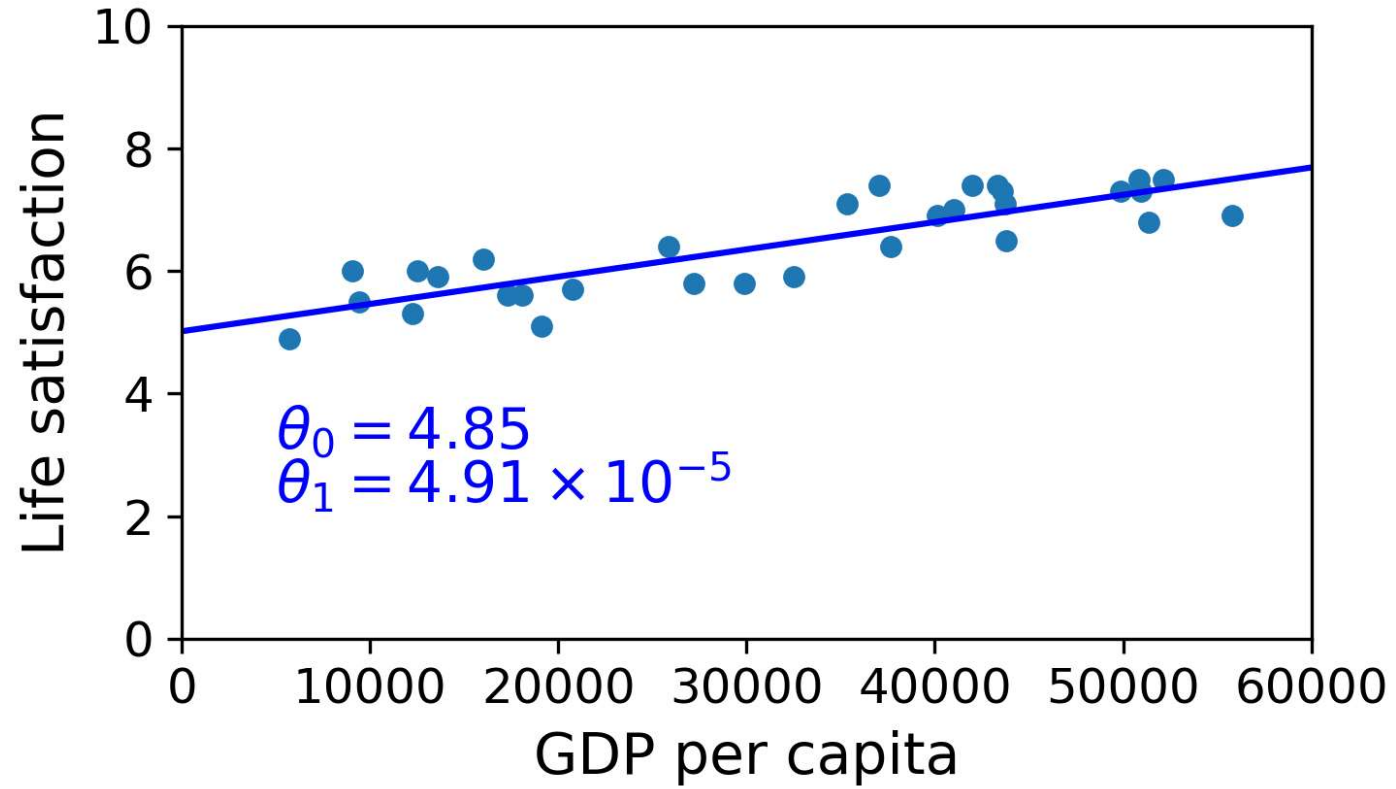
Is there a trend?



Linear model: $sat = \theta_0 + \theta_1 * gdp\_per\_capita$

Finding good models?
How is "good" defined

$$\theta_0 = 4.85$$
$$\theta_1 = 4.91 \times 10^{-5}$$

Donnerstag, 6. Mai 2021        10:30
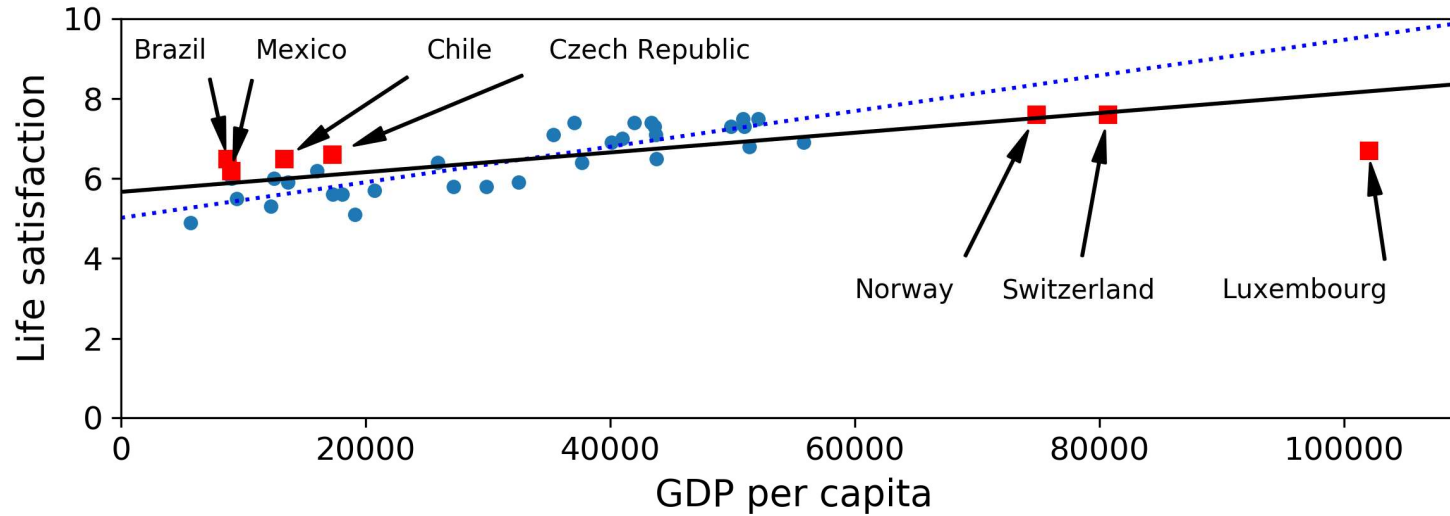


Take GDB per capita and predict sat

- Not enough data
- Data is not representative
- Bad data quality
- Irrelevant features
- Overfitting
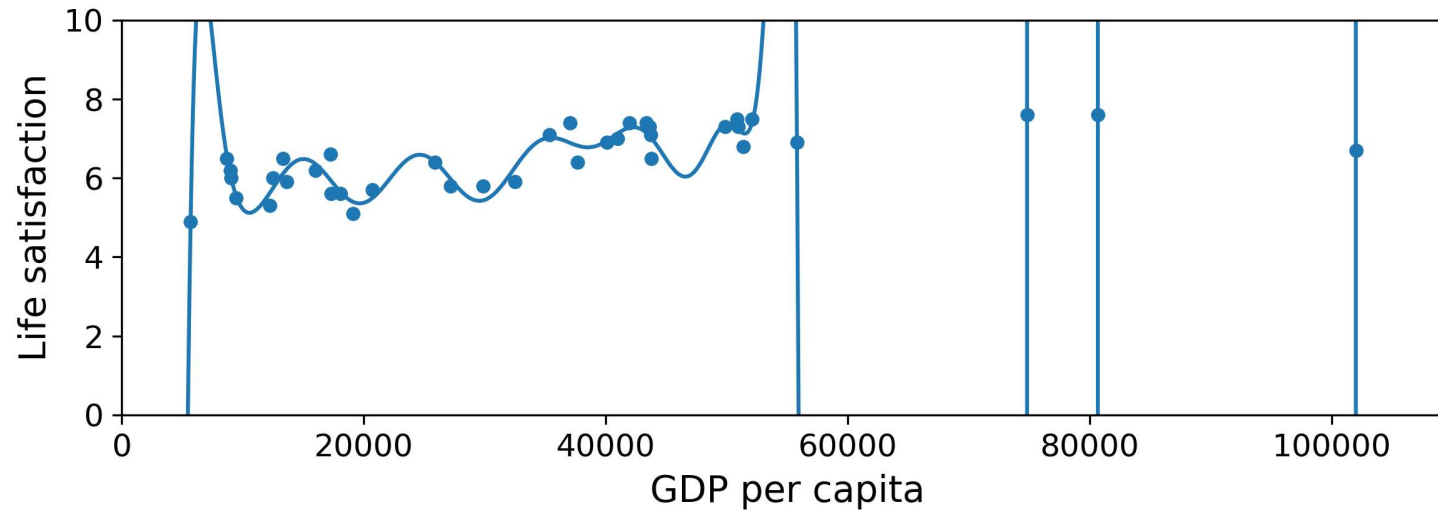- Underfitting

Adding Brasil, Mexico, …, Norway, … changes model
Linear model does not really fit



Countries with high GDP per capita but lower sat
Countries with low GDP per capita but higher sat

Komplex model
Polynomial with high degree

Goal:
- Find relationship between features and targets
- Discover patterns
- Find correlations

Common terms:
- Machine Learning
- Data Mining
- Statistical Learning

Reasons why learning from data
- Constructing an abstraction for the relationship between features and targets (model)
- Complex problems where algorithmic formulation is hard
- Environments where data changes continuously
- Problems that require a long list of complicated rules

Analysis of spam emails, identifiction of offensive words and phrases



Long list of rules - not easily maintainable
Adaption to changes are hard to do

Identification auf words and phrases by learning from data using suitable algorithms.



Data

Analyse Problem → Train → Evaluate → Deploy

Evaluate → Analyse Errors → Analyse Problem

Adaption to changes much easier: just take new data and learn again

## Features

- Attributes that characterize a particular instance/object
- Also called predictors

## Targets

- Scalar value in case of regression
- Nominal value in case of classification

Features / Targets

### Regression

| Bedrooms | Bathrooms | Latitude | Longitude | Price |
|----------|-----------|----------|-----------|-------|
| 1 | 1 | 40.71 | -73.94 | 3055 |
| … | … | … | … | … |

### Classification

| Income | Job | IsMarried | Age | Loan |
|--------|-----|-----------|-----|------|
| … | … | … | … | Yes |
| … | … | … | … | No |

| x1 | = | x11 | x12 | x13 | ... | x1n |
|----|---|-----|-----|-----|-----|-----|
| x2 | = | x21 | x22 | x23 | ... | x2n |
| ... | | ... | | | | |
| xm | = | xm1 | xm2 | xm3 | ... | xmn |

| y1 |
|----|
| y2 |
| ... |
| ym |

X                    y

| xi | Feature vector |
|----|----------------|
| xij | Feature value |
| yi | Target value |
| X | Feature matrix |
| y | Target vector |

- All entries must be numbers, i.e. all feature and target values must be transformed into numbers
- Each feature vector is one instance/object

Boxplot

each point is
one person
(object, instance)

Scatterplot

Learning should
- be as accurate as possible
- generalize to new data as good as possible

Conflicting goals!

Just remember target value wrt. feature combination
- Good accuracy
- Bad generalization
- Combinations might not exist
- Combinations might not be unique

| bedrooms | bathrooms | latitude | longitude | price |
|---|---|---|---|---|
| 0 | 1.0000 | 40.7073 | -73.9664 | 2650 |
| 0 | 1.0000 | 40.7073 | -73.9664 | 2850 |
| 0 | 1.0000 | 40.7073 | -73.9664 | 2950 |
| 0 | 1.0000 | 40.7073 | -73.9664 | 2850 |

Build coarse grained groups, e.g. price per bedrooms and bathroom combinations
- Bad accuracy
- Good generalization

| Bedrooms | Bathroom | Price |
|---|---|---|
| 1 | 1 | 3000 |
| 2 | 1 | 3700 |
| 1 | 2 | 3010 |
| ... | ... | ... |

Learn from neighbors
(k nearest neighbor)
- Powerful
- But must store complete training data set
- No real learning
- Finding nearest neighbors requires search, might be slow

Build a linear model:
Price = w0 +
    w1*bedrooms +
    w2*bathrooms +
    w3*latitude +
    w4*longitude
Might be not suitable: higher number of bathrooms doesn't mean higher price

Supervision
- Supervised
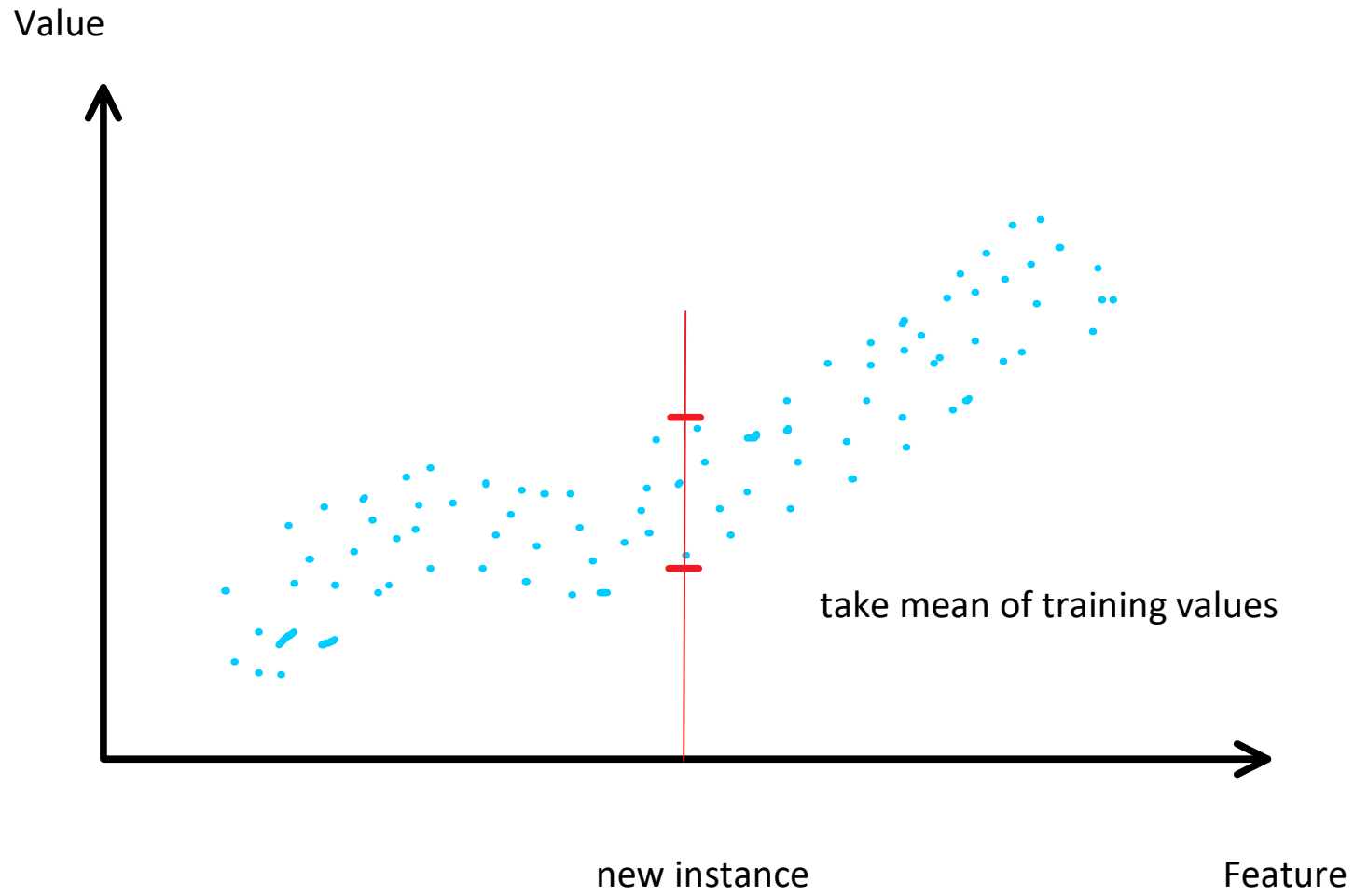- Unsupervised
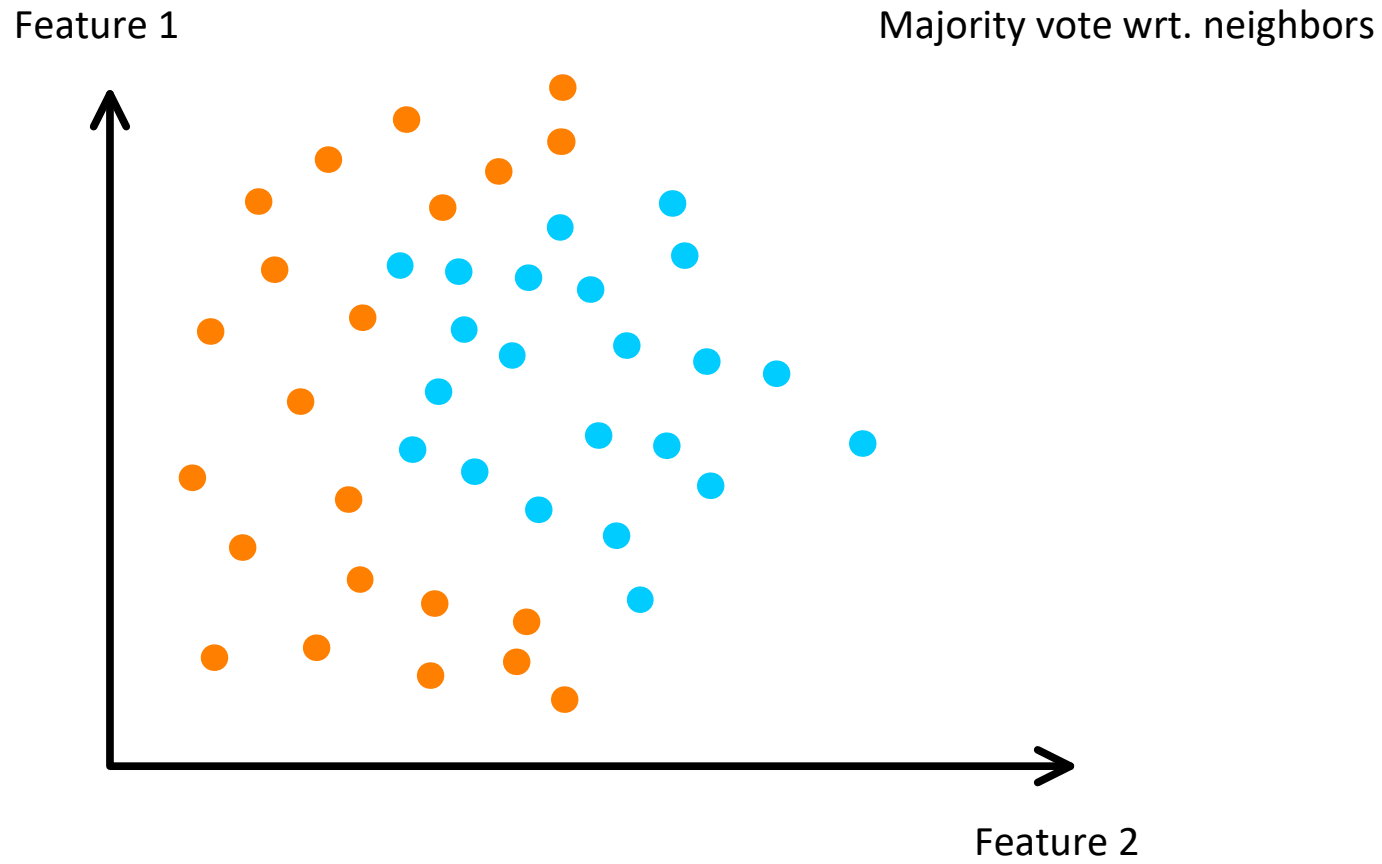- Semisupervised
- Reinforcement Learning

Data Usage
- Batch
- Online

Method
- Instance-Based
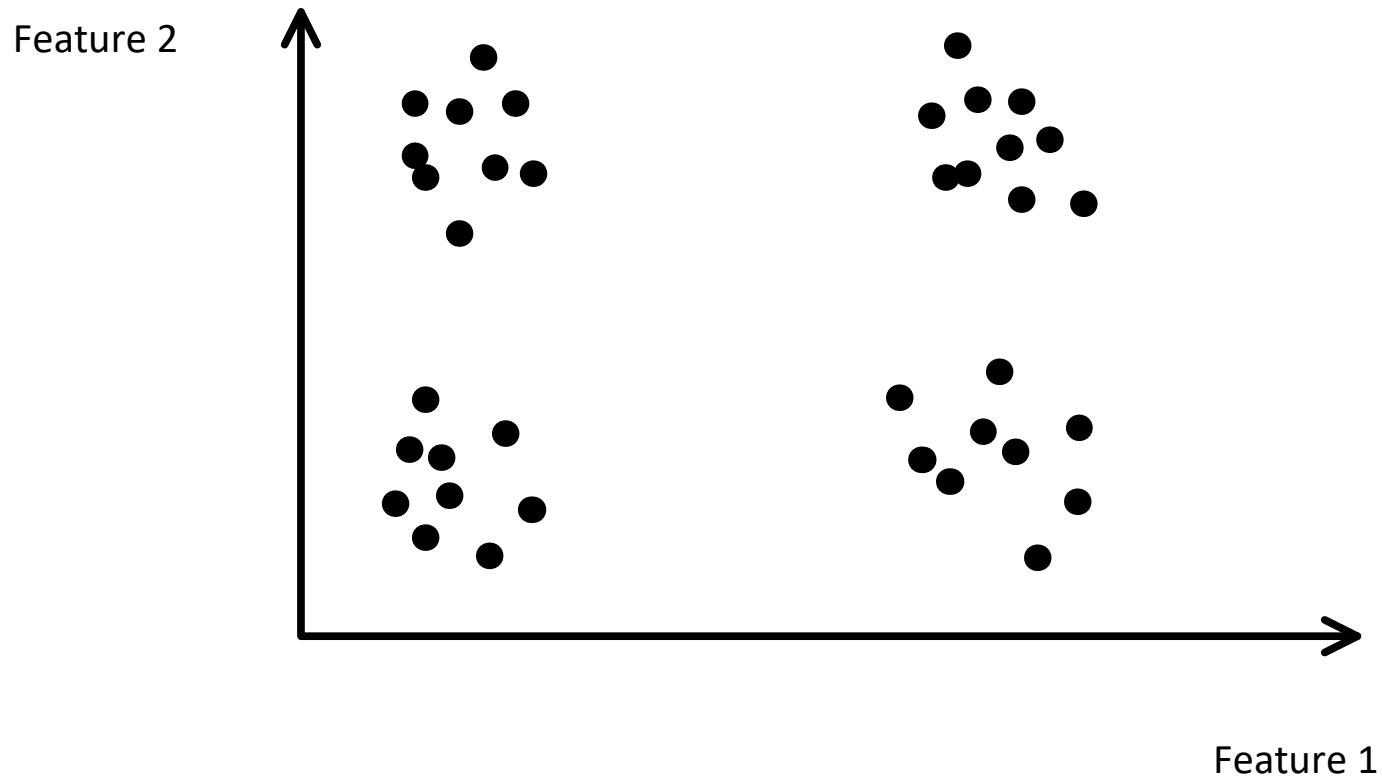- Model-Based

Features of instances + target values

Feature 1

Majority vote wrt. neighbors



Feature 2

Class tags:
- Class1: orange
- Class2: blue

Features  of instances and no class tags
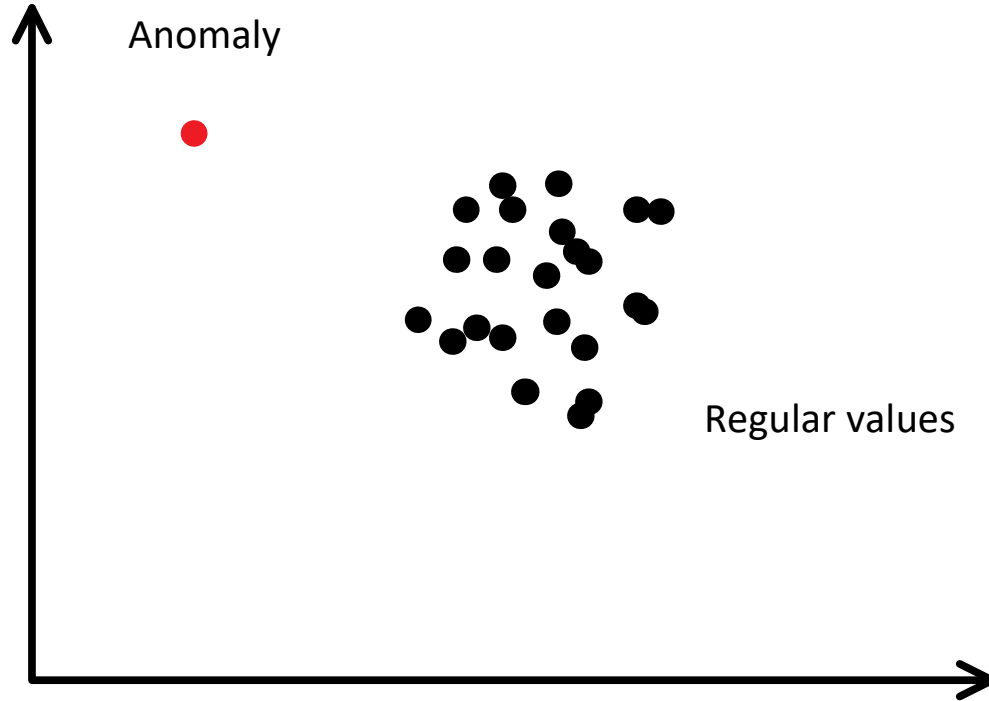
Feature 2

Feature 1

Find sets of data with strong coupling
and large distance between sets

Features of instances

Compute cluster of regular instances

Calculate distance



Feature 2

Anomaly

Regular values

Feature 1

Market Basket Analysis
Learning of rules

{ Milk, Bread } -> { Butter }

{ PC, Monitor, Graphics Card } -> { PC Game }

Friday evening
{ Beer } -> { Diapers }

Example: Identification of persons on your private fotos
- Upload of fotos
- System detects same persons on different fotos (unsupervised)
- You tag with name (supervised)

- Agent observes environment
- Chooses action according to strategy
- Executes action
- Gets reward or punishment
- Learns from feedback - adopts strategy
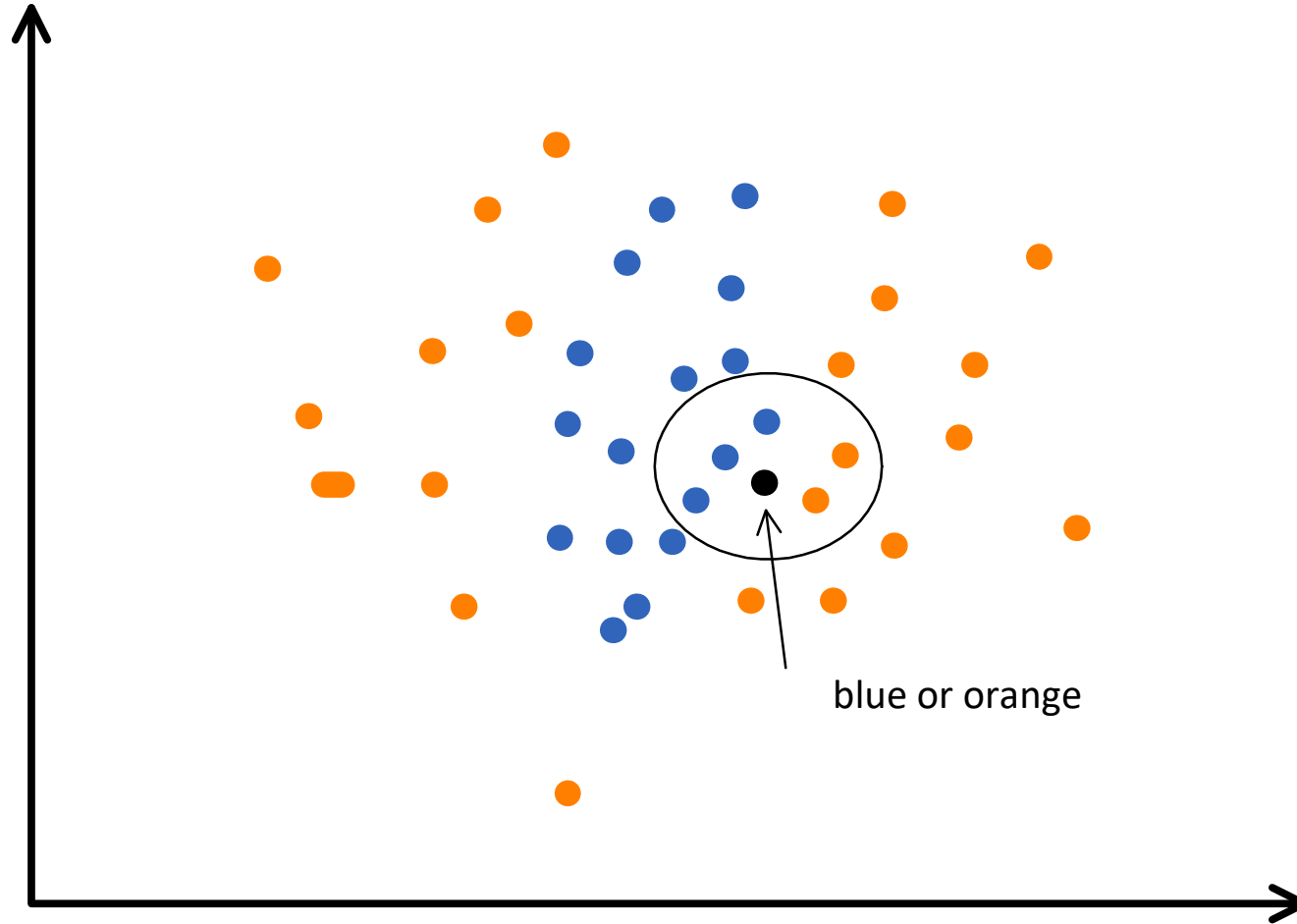- Repeat until optimal solution found (or at least a satisfying solution)

Batch Learning
- Take all data at once and learn from it
- Complete new take if data changes
- Requires long execution time and ressources in case of large training data set
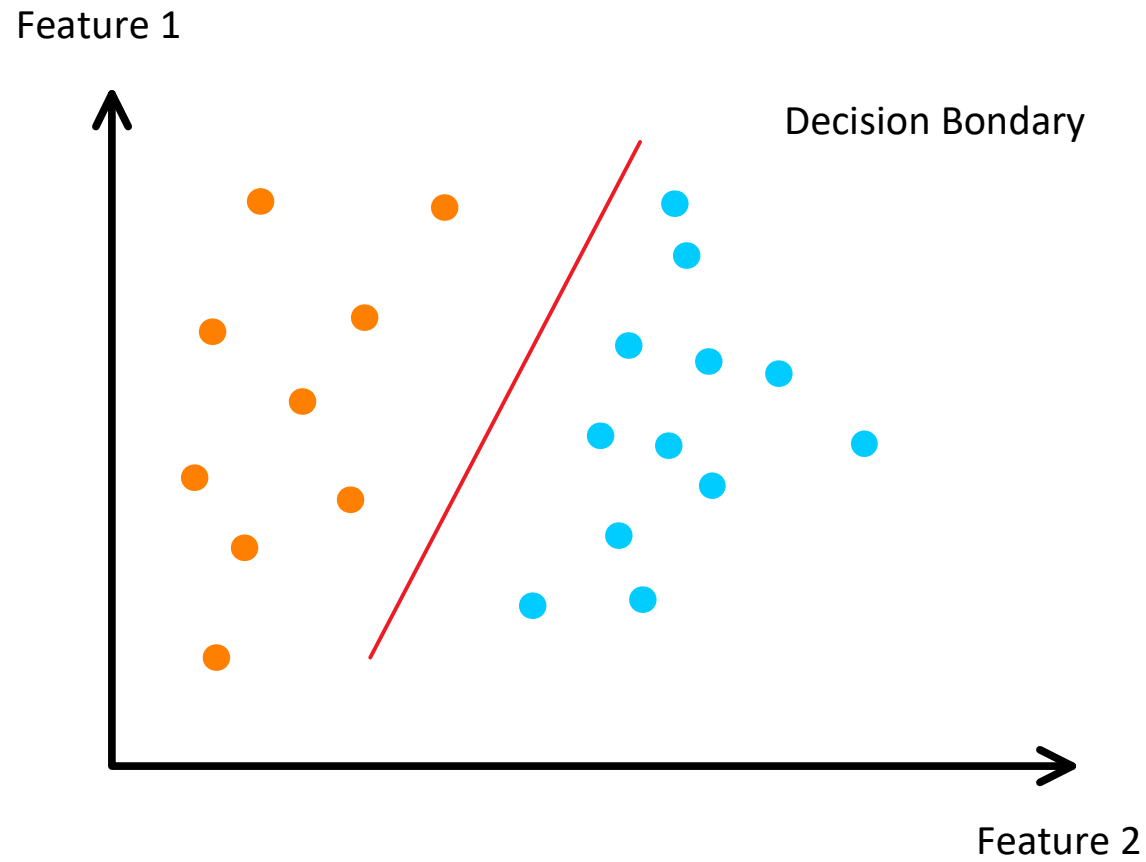
Online Learning
- Incremental usage of data
- Single data points or minibatches
- Faster adoption to new data
- More volatile models

Similarity to training data



blue or orange

Keep all training data
Use distance metrics

Donnerstag, 6. Mai 2021      08:20



No storage of trainings data
Just keeping model parameters, e.g. slope and intercept of separating line

Parr, Howard: The Mechanics of Machine Learning, https://mlbook.explained.ai/

James, Witten, Hastie, Tibshirani: An Introduction to Statistical Learning,
Springer, 2013, http://www.statlearning.com/

Géron: Hands-On Machine Learning with Scikit-Learn and TensorFlow,
O'Reilly, 2017, http://shop.oreilly.com/product/0636920052289.do