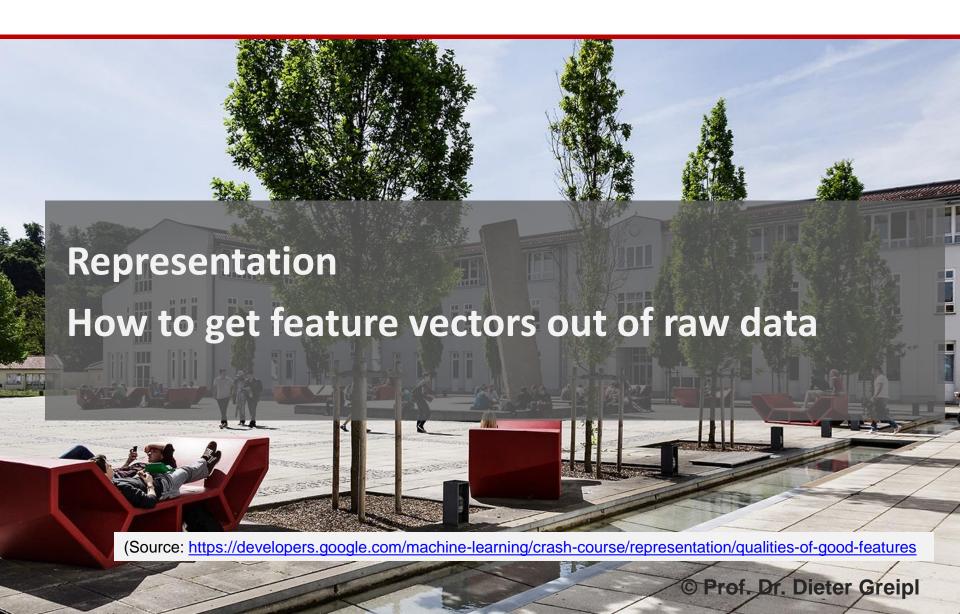
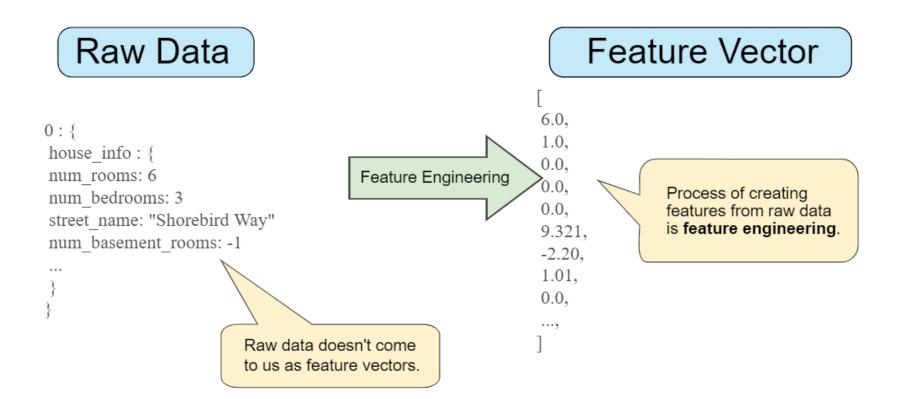
# **Machine Learning - Introduction**





## Feature Engineering: Mapping Raw Data to Features





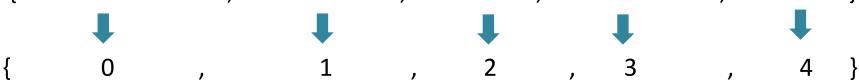
### Mapping Categorical Values



Categorical features have a <u>discrete</u> (finite) set of possible values. For example, there might be a feature called faculty with options that include:

### Simple Solution:

{'Betriebswirtschaft', 'Soziale Arbeit', 'Informatik', 'Maschinenbau', 'E-Technik'}



#### Raw Data

**Feature** 

faculty: "Betriebswirtschaft"

Feat. Engin.

0

#### **Problems**

- How to deal with multiple values?
- Impact of values in a (pseudo)-linear Model?

Prof. Dr. Dieter Greipl

### Mapping Categorical Values: One Hot Encoding



#### Raw Data

faculty: "Betriebswirtschaft"

Feat. Engin.

**Feature** 

[1,0,0,0,0]

faculty: "Informatik"

Feat. Engin.

[0,0,1,0,0]

{'Betriebswirtschaft', 'Soziale Arbeit', 'Informatik', 'Maschinenbau', 'E-Technik'}

Prof. Dr. Dieter Greipl

### Tips for good features



- Avoid rarely used discrete feature values
- Clear and obvious meanings
- Don't mix "magic" values with actual data
- Account for upstream instability

### Cleaning feature values: Know your data



- Cleaning ("Scrubbing"), i.e.
  - handle missing values
  - remove duplicate examples
  - remove false labels
- Outlier Handling (Logarithmic Scaling)
- Binning (if there is obviosuly no linear relationship)

- Scaling

Prof. Dr. Dieter Greipl