

# Data Preprocessing

Sibylle Hess



# BASIC DATA CHARACTERIS- TICS

## Covid: Man offered vaccine after error lists him as 6.2cm tall



## Coronavirus pandemic



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡



---

- categories
- with order
- subjective scale

ID	eyecolor	sex	zip code
1	blue	m	5629NZ
2	green	f	5381
⋮	⋮	⋮	⋮

ID	rating	size	grade
1	★★★★☆	S	7.5
2	★★★☆☆	XL	9
⋮	⋮	⋮	⋮

## Quantitative Features

## Discrete Data:

- countable amount of values
- distances are meaningful

ID	age	nr. children	nr. yes votes
1	22	2	1298
2	24	0	2780
$\vdots$	$\vdots$	$\vdots$	$\vdots$

### Continuous Data:

- measurements
- distances are meaningful

ID	income[€]	height[m]	temp[°C]
1	2260.35	1.75	23
2	3502.60	1.68	26
⋮	⋮	⋮	⋮

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡



# STATISTICS MAKE SENSE FOR ALL TYPES.

\_\_\_\_\_

ID	pic	size
1	1	1
2	0	4
$\vdots$	$\vdots$	$\vdots$
$n$	2	1

9. 11. 2014

\_\_\_\_\_

Create for every value a new binary feature.



ID	pic			
	cat	dog	fish	bird
1	0	1	0	0
2	1	0	0	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	0	0	1	0

**PROBLEMATIC:** increases the dimensionality of the feature space

Usually applied for nominal data

\_\_\_\_\_

statistic		permissible for nominal	ordinal
mode	most frequent value	yes	yes
median	$\begin{cases} x_{(\frac{n+1}{2})} & \text{if } n \text{ is odd} \\ \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2} & \text{otherwise} \end{cases}$	no	yes
mean	$\mu = \frac{1}{n} \sum_{i=1}^n x_i$	no	no
variance	$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$	no	no

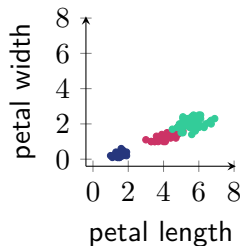
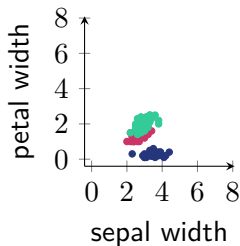
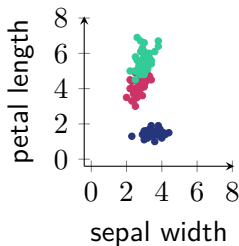
For this, we have a look at the Iris dataset.

\_\_\_\_\_

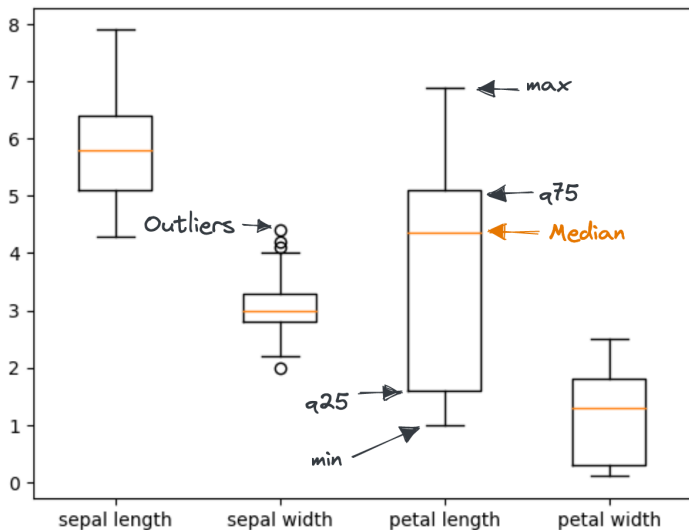


sepal length	sepal width	petal length	petal width	class
5.1	3.5	1.4	0.2	setosa
6.4	3.5	4.5	1.2	versicolor
5.9	3.0	5.0	1.8	virginica
⋮	⋮	⋮	⋮	⋮

\_\_\_\_\_



---





# FEATURE TRANSFORMA- TIONS

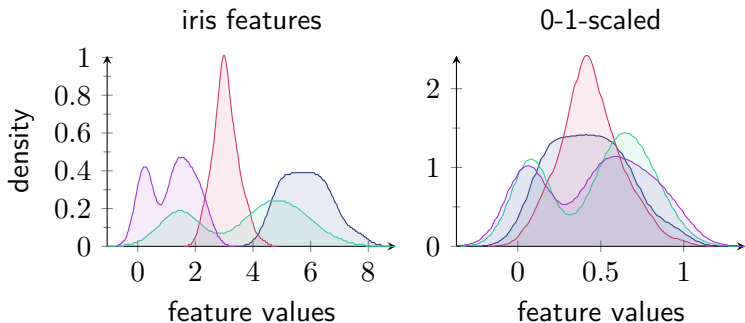
100

- If one feature is measured in kg and the other in g, differences in the feature in g might seem more important
- Most ML methods learn based on similarities/distances between data points. These similarities/distances might be distorted when some of the features have much bigger values than others

\_\_\_\_\_

[illegible]

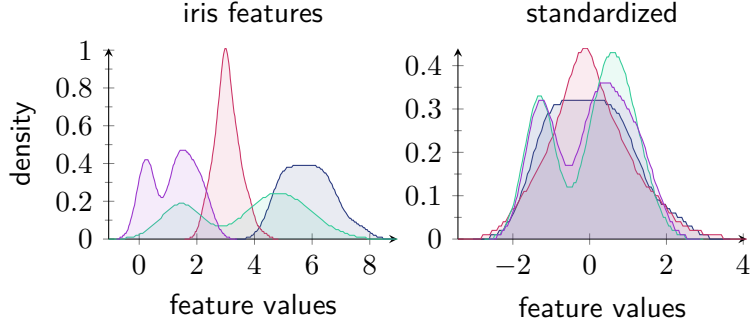
$$\hat{D}_{ik} = \frac{D_{ik} - \min(D_{.k})}{\max(D_{.k}) - \min(D_{.k})}$$

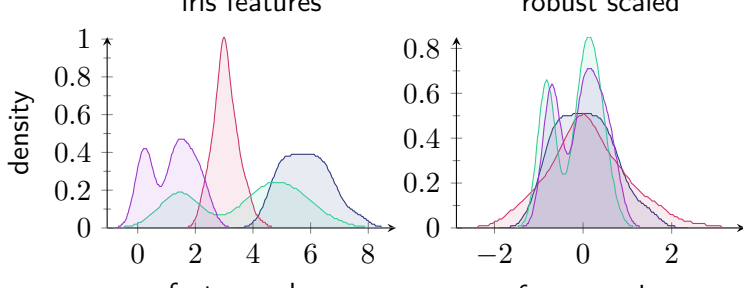


sepal length sepal width petal length petal width

---

$$\hat{D}_{ik} = \frac{D_{ik} - \mu_{F_k}}{\sigma_{F_k}}$$



$$D_{\cdot} = \text{median}(D_{\cdot})$$


1

**MIN-MAX SCALING:** might make sense for ordinal data, but is not robust against outliers

$$\hat{D}_{ik} = \frac{D_{ik} - \min(D_{\cdot k})}{\max(D_{\cdot k}) - \min(D_{\cdot k})}$$

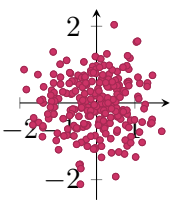
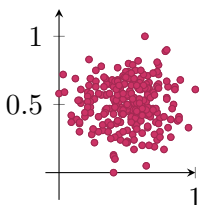
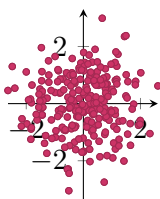
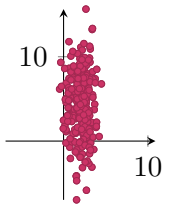
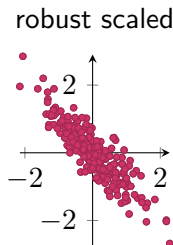
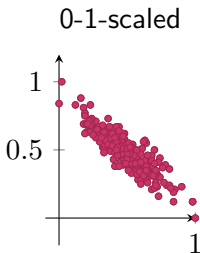
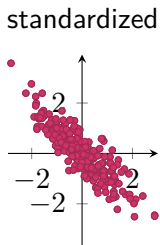
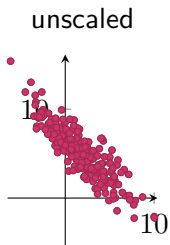
**STANDARDIZATION:** makes mostly sense for normal distributed data, a bit robust to outliers

$$\hat{D}_{ik} = \frac{D_{ik} - \mu_{F_k}}{\sigma_{F_k}}$$

**ROBUST SCALING:** might make sense for ordinal data, and is robust against outliers

$$\hat{D}_{ik} = \frac{D_{ik} - \text{median}(D_{.k})}{q_{75}(D_{.k}) - q_{25}(D_{.k})}$$

1. *Journal of Management Studies*, 1997, 34(1), 1-15.



---





Scaling sounds fair, every feature is treated the same, e.g., every feature gets the same variance.

However, those differences in, e.g., the variance can also tell you about the usefulness of the features.

→ **FEATURE SELECTION**

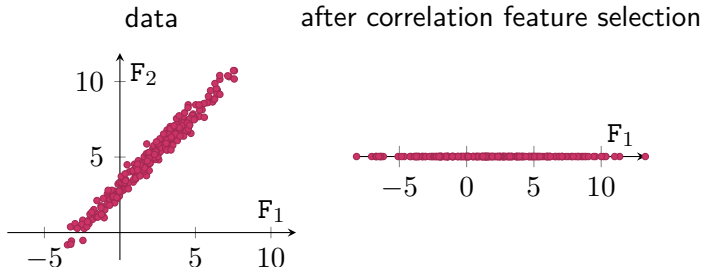
# Feature Selection

Feature selection removes features which are probably not needed.  
This is useful, because:

- every feature adds a dimensionality to your data points,
- in high dimensional data, data points tend to be equidistant to each other, so all observations seem to be alike  
→ **CURSE OF DIMENSIONALITY**,
- having fewer features improves the interpretability of results:  
in many tasks we want to know which features are important to get a good model.



---



# INTERMISSION: COVARIANCE AND CORRELATION

# Covariance as Inner Product Similarity

We compute the sample covariance of two feature vectors by

$$\text{cov}(D_{\cdot k}, D_{\cdot l}) = \frac{1}{n} \sum_{i=1}^n (D_{ik} - \mu_{F_k})(D_{il} - \mu_{F_l}),$$

If we consider the centered data matrix  $C = D - \mathbf{1}\mu_F$ , then the covariance can be written as an inner product:

$$\text{cov}(D_{\cdot k}, D_{\cdot l}) = \frac{1}{n} \sum_{i=1}^n C_{ik} C_{il} = \frac{1}{n} C_{\cdot k}^\top C_{\cdot l}$$

That is, the **COVARIANCE** can be seen as **THE INNER PRODUCT SIMILARITY** of the centered feature values.

# Covariance and Variance

The **VARIANCE** is the covariance of a feature with itself and thus equal to the squared norm:

$$\sigma_{F_k}^2 = \text{cov}(D_{\cdot k}, D_{\cdot k}) = \frac{1}{n} C_{\cdot k}^\top C_{\cdot k} = \frac{1}{n} \|C_{\cdot k}\|^2$$

The covariance between two features is large when the angle between the feature vectors is small and the variance of the features is high

$$\begin{aligned} \text{cov}(D_{\cdot k}, D_{\cdot l}) &= \frac{1}{n} C_{\cdot k}^\top C_{\cdot l} = \frac{1}{n} \cos \angle(C_{\cdot k}, C_{\cdot l}) \|C_{\cdot k}\| \|C_{\cdot l}\| \\ &= \cos \angle(C_{\cdot k}, C_{\cdot l}) \sigma_{F_k} \sigma_{F_l} \end{aligned}$$



# The Covariance Matrix

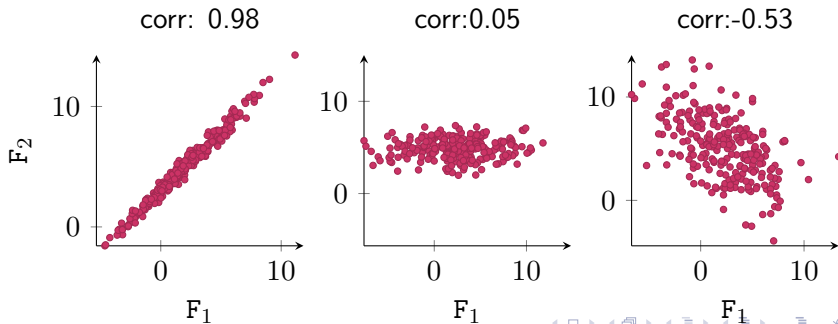
With  $C$  being the centered data matrix, the covariance matrix is given by

$$C^T C = \begin{pmatrix} \sigma_{F_1}^2 & \text{cov}(D_{.1}, D_{.2}) & \dots & \text{cov}(D_{.1}, D_{.d}) \\ \text{cov}(D_{.2}, D_{.1}) & \sigma_{F_1}^2 & \dots & \text{cov}(D_{.2}, D_{.d}) \\ \vdots & \vdots & \vdots & \vdots \\ \text{cov}(D_{.d}, D_{.1}) & \text{cov}(D_{.d}, D_{.2}) & \dots & \sigma_{F_d}^2 \end{pmatrix}$$

---

$$\text{corr}(D_{.k}, D_{.l}) = \frac{\text{cov}(D_{.k}, D_{.l})}{\sigma_{F_k} \sigma_{F_l}} = \cos \angle(C_k, C_l) \in [-1, 1]$$

The correlation between two features is the normalized version of the covariance and it's large when the angle between the centered feature vectors is small.



END OF  
INTERMISSION

# Transformations of the Feature Space

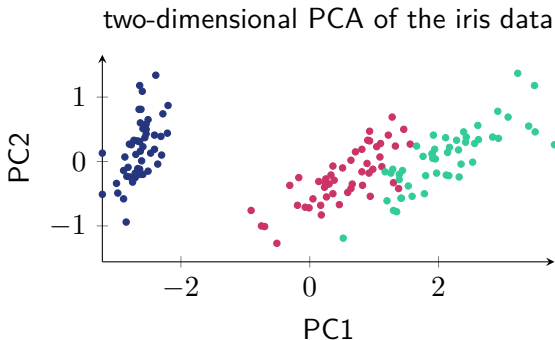
In some cases, it makes sense to transform the features or add transformed features, e.g., transforming the data can transfer nonlinear relations into linear ones.

ID	$F_1$	$F_2$	$\rightarrow$	ID	$F_1^2$	$F_1 F_2$	$F_2^2$	$\log(F_2)$
1	1.1	3.2		1	1.21	3.52	10.24	1.16
2	2.0	2.7		2	4.00	5.40	7.29	1.69
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

The task of dimensionality reduction is to find a  
**TRANSFORMATION INTO A  
LOW-DIMENSIONAL  
FEATURE SPACE** which  
keeps characteristics of the  
data

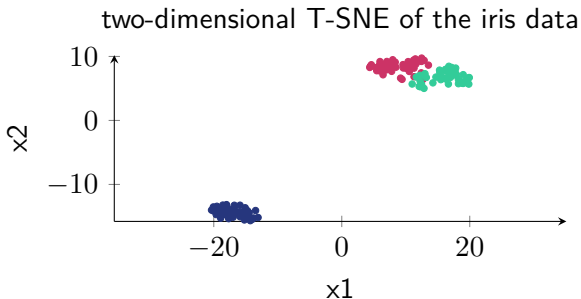
# Principal Component Analysis (PCA)

The goal of PCA is to **MAINTAIN MOST OF THE VARIANCE** in the low-dimensional view.



We will discuss PCA and how it works later in the course.

---



<https://digitall.pub/2016/microad-tano/>