

# Introduction to Machine Learning

## PCA

**Bernd Bischl, Christoph Molnar, Daniel Schalk, Fabian Scheipl**

Department of Statistics – LMU Munich



# Introduction

# SUGGESTED LITERATURE

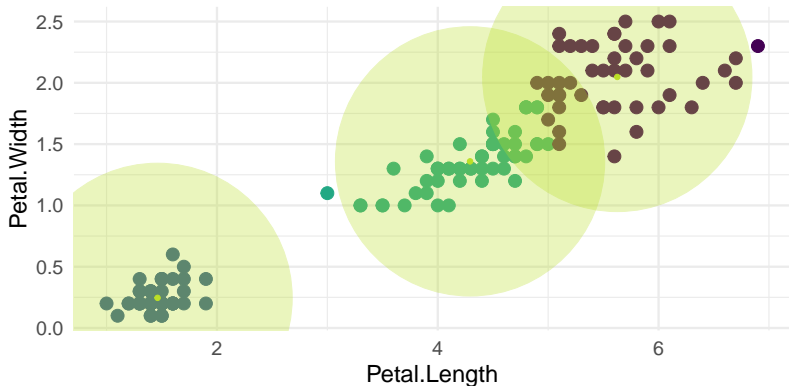
- Hastie, T., Tibshirani, R., Friedman, J. (2009): The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer.
- James, G., Witten, D., Hastie, T., Tibshirani, R. (2013): An Introduction to Statistical Learning with Applications in R. Springer.
- Aggarwal, C. C., & Reddy, C. K. (Eds.). (2013). Data Clustering: Algorithms and Applications. CRC press.

# UNSUPERVISED LEARNING

- Supervised machine learning deals with \*labeled\* data, i.e., we have input data  $x$  and the outcome  $y$  of past events.
- Here, the aim is to learn relationships between  $x$  and  $y$ .
- Unsupervised machine learning deals with data that is \*unlabeled\*, i.e., there is no real output  $y$ .
- Here, the aim is to search for patterns within the inputs  $x$ .

# CLUSTERING TASK

**Goal:** Group data into similar clusters (or estimate fuzzy membership probabilities)



# CLUSTERING: CUSTOMER SEGMENTATION

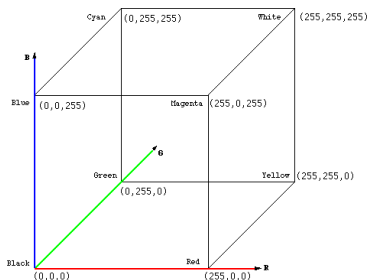
- In marketing, customer segmentation is an important task to understand customer needs and to meet with customer expectations.
- Customer data is partitioned in terms of similarities and the characteristics of each group are summarized.
- Marketing strategies are designed and prioritized according to the group size.

## Example Use Cases:

- Personalized ads (e.g., recommend articles).
- Music/Movie recommendation systems.

# CLUSTERING: IMAGE COMPRESSION

- An image consists of pixels arranged in rows and columns.
- Each pixel contains **RGB** color information, i.e., a mix of the intensity of 3 **primary colors**: **R**ed, **G**reen and **B**lue.
- Each primary color takes intensity values between 0 and 255.



Source: By Ferlixwang CC BY-SA 4.0, from Wikimedia Commons.

# CLUSTERING: IMAGE COMPRESSION

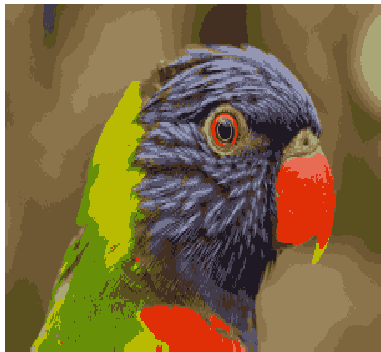
An image can be compressed by reducing its color information, i.e., by replacing similar colors of each pixel with, say,  $k$  distinct colors.

**Example:**

Original Image



Image using 16 Colors

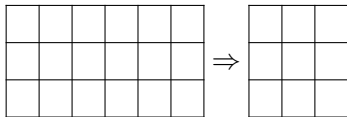




# DIMENSIONALITY REDUCTION TASK

**Goal:** Describe data with fewer features (reduce number of columns).

⇒ there will always be an information loss.



Unsupervised Methods:

- Principle Component Analysis (PCA).
- Factor Analysis (FA).
- Feature filter methods.

Supervised Methods:

- Linear Discriminant Analysis (LDA).
- Feature filter methods.

# Principal Component Analysis

# NORMALIZING DATA

A variable  $X$  can be normalized by subtracting its values with the mean  $\bar{X}$  and dividing by the standard deviation  $s_X$ , e.g.  $\tilde{X} = \frac{X - \bar{X}}{s_X}$ .

## Example:

Consider the following body heights measured in different units:

	Person A	Person B	Person C	mean	sd
body height (cm)	180.00	172.00	175.00	175.67	4.04
body height (m)	1.80	1.72	1.75	1.76	0.04
body height (feet)	5.91	5.64	5.74	5.76	0.13

After normalizing, we always obtain the normalized body height (no matter which unit was used):

	Person A	Person B	Person C	mean	sd
normalized body height	1.07	-0.91	-0.16	0.00	1.00

# NORMALIZING DATA

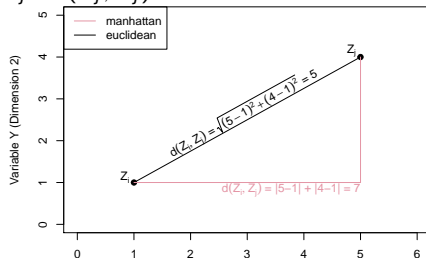
Normalizing all variables in a data set, can have several advantages:

- It puts all variables into \*comparable\* units, i.e., we make sure that all normalized variables have mean 0 and standard deviation of 1.
- It can avoid numerical instabilities in several algorithms, e.g. if a variable has very low / high values.
- It helps in computing meaningful \*distances\* between observations.

# NORMALIZING DATA: DISTANCES

There are many ways to define the distance between two points, e.g.,

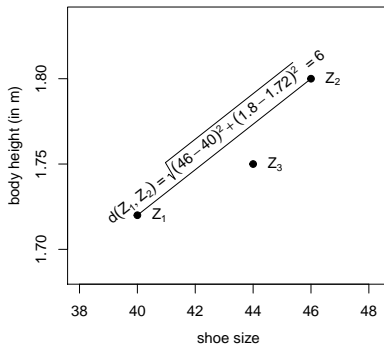
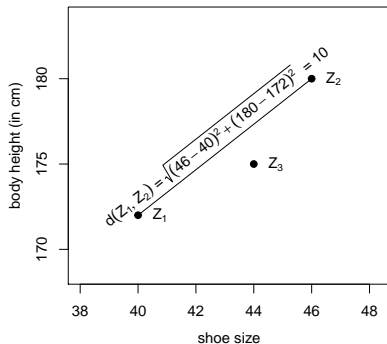
$Z_i = (X_i, Y_i)$  and  $Z_j = (X_j, Y_j)$ :



- manhattan: sum up the absolute distances in each dimension.
- euclidean: remember Pythagoras theorem from school?

# NORMALIZING DATA: DISTANCES

It is often a good idea to *normalize* the data before computing distances, especially when the scale of variables is different, e.g. the euclidean distance between the point  $Z_1$  and  $Z_2$ :



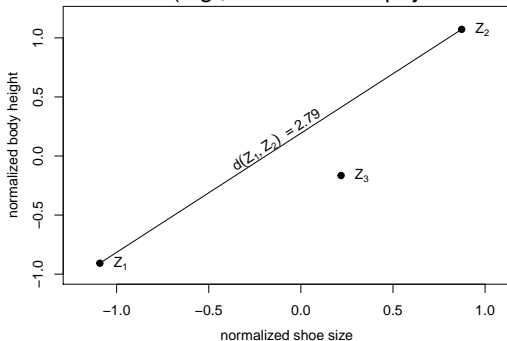
On the right plot, the distance is dominated by “shoe size”.

# NORMALIZING DATA: DISTANCES

The normalized variable  $\tilde{X}_{\text{shoe.size}}$  is computed by `<!-- Normalization of the shoe.size variable means: -->`

$$\tilde{X}_{\text{shoe.size}} = \frac{X_{\text{shoe.size}} - \bar{X}_{\text{shoe.size}}}{s_{X_{\text{shoe.size}}}}.$$

Distances based on normalized data are better comparable and **robust** in terms of linear transformations (e.g., conversion of physical units).



# NORMALIZING: COVARIANCE VS. CORRELATION

The **variance** of a normalized variable is always 1, its mean is always 0. The **covariance** of two normalized variables  $\tilde{X} = \frac{X - \bar{X}}{s_X}$  and  $\tilde{Y} = \frac{Y - \bar{Y}}{s_Y}$  is the same as the **correlation** of the non-normalized variables  $X$  and  $Y$ . One can prove this with the help of

$$s_{\tilde{X}\tilde{Y}} = \frac{1}{n-1} \sum_{i=1}^n (\tilde{x}_i - \bar{\tilde{x}})(\tilde{y}_i - \bar{\tilde{y}}) = \dots = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})}{s_X} \frac{(y_i - \bar{y})}{s_Y} = r_{XY}.$$