

# Machine Learning

## Introduction & distance function

S. Nõmm

<sup>1</sup>Department of Software Science, Tallinn University of Technology

30.01.2018

## Course organization

- Lectures on Tuesdays 16:00 - 17:30 ICT-A1
- Practices on Thursdays 16:00 - 17:30 ICT-401 (Computer room).
- Programming language: MATLAB only!!!
- Each home assignment gives 10% of the final grade. NB! all the assignments should be accepted with positive grade!!!
- Two closed book tests (mandatory!) each gives 10% of the final grade.
- Final exam (Written report on assign topic (programming is required) + short presentation) will give 50% of the final grade.
- Assignment submission via [ained.ttu.ee](http://ained.ttu.ee) ITI8565 Machine Learning 2018 only!
- For correspondence please use [sven.nommm@ttu.ee](mailto:sven.nommm@ttu.ee) Please do not call me by phone!!!
- Consultations by appointment only. Free time in the computer class on Fridays ICT-405 12:00 - 14:00 your teacher is not present this time!!!

## Organisation clarified

**No Plagiarism in any of your tests and final project!!!**. You should cite all the references, including libraries you use to complete your computational assignments. The student should be able to explain the meaning of all the computations performed, interpret and present the results. Grading: Detailed information about the grading may be found in <https://ois.ttu.ee>.

- Final project gives 50% of the final grade. Guidelines will be available from the course page one month before the examination date. During examination each student will be given 5 min. to explain the problem, describe chosen methodology, and present the results. This part will be followed by more detailed examination of the implementation, student may be asked to change implementation on the fly and/or run it on the new data set. Precise project preparation guidelines will appear one month before the deadline.
- Note, your own implementation graded higher than usage of third party libraries.

# References

Present lectures are based on:

- Machine Learning: A Probabilistic Perspective Textbook by Kevin R. Murphy. MIT Press, Aug 24, 2012 - Computers - 1067 pages (available in TUT Library).
- The Elements of Statistical Learning Book by Jerome H. Friedman, Robert Tibshirani, and Trevor Hastie, Springer (available through the TUT library as e-book, UNIID required).
- Pattern Recognition and Machine Learning Book by Christopher Bishop. Springer. (available in TUT Library)
- Data Mining The Text Book by C. Aggarwal. 2015 Springer.

In course of the lectures, whenever necessary additional references will be shared.

# Prerequisites

- Knowledge of Calculus, Foundations: function, derivative etc.
- Knowledge of Linear algebra, Foundations: matrixes, matrix operations, polynomials etc.
- Medium level of programming skills.
- Foundations of probability theory and statistics.
- You are **STRONGLY** advised to write your own notes!!!
- Lecture slides could not substitute attendance, there are quite many things are discussed during lectures which are not reflected in the slides.
- It is assumed that students are: attending the lectures, read books and do home works regularly!!!

# Course plan

- Introduction: Metric (distance) function. Note on data preparation.
- Unsupervised learning: Clustering and validation. Representative methods, agglomerative methods density based methods, probabilistic methods. impact of metric function.
- Supervised learning: Classification and regression.  $k$ -nearest neighbours, decision trees, support vector machines, neural networks.
- Ensemble learning, bagging, boosting.
- Markov models.
- Introduction to deep learning.

# What is Learning?

According to:

<http://www.merriam-webster.com/dictionary/learning>

**learning** noun

: the activity or process of gaining knowledge or skill by studying, practicing, being taught, or experiencing something : the activity of someone who learns

- the act or experience of one that learns
- knowledge or skill acquired by instruction or study
- modification of a behavioral tendency by experience (as exposure to conditioning)

# What is *Machine Learning*?

Arthur Samuel, 1959 has defined Machine learning as: Field of study that gives computers the ability to learn without being explicitly programmed.

Tom Mitchell, 1997 has defined Machine learning as:

A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

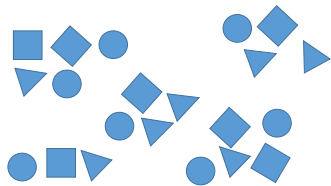


# Main steps

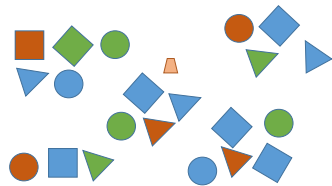
- Data acquisition.
- Data preparation.
- Feature selection.
- Model training.
- Tuning (sometimes).
- Model validation.

From this point it is expected that when solving an exercise student chooses and executes all necessary steps and able to explain their decision.

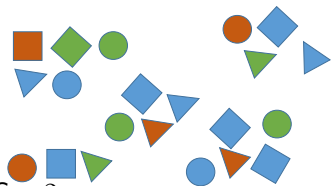
# Similarity and distance



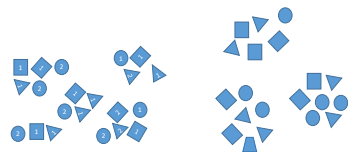
Set 1



Set 3



Set 2



set 4

# Distance function

Distance function is one of most fundamental notions in Machine learning and Data mining. Formally defined in pure mathematics as *metric* function. It provides measure of similarity or distance between two elements.

## Definition

A function  $S : X \times X \rightarrow \mathbb{R}$  is called metric if for any elements  $x$ ,  $y$  and  $z$  of  $X$  the following conditions are satisfied.

- 1 Non-negativity or separation axiom

$$S(x, y) \geq 0$$

- 2 Identity of indiscernible, or coincidence axiom

$$S(x, y) = 0 \Leftrightarrow x = y$$

- 3 Symmetry

$$S(x, y) = S(y, x)$$

- 4 Subadditivity or triangle inequality

$$S(x, z) \leq S(x, y) + S(y, z)$$

# Distance function: Examples 1 (Most common distance functions)

- Euclidean distance

$$S(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- Manhattan distance also referred as city block distance or taxicab distance

$$S(x, y) = \sum_{i=1}^n |x_i - y_i|$$

- Chebyshev distance

$$S(x, y) = \max_i (|x_i - y_i|)$$

# Distance function: Examples 2

## Euclidean distance

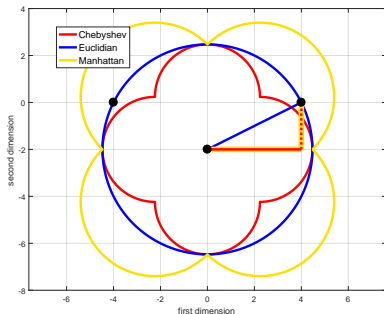
$$S(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

## Manhattan distance

$$S(x, y) = \sum_{i=1}^n |x_i - y_i|$$

## Chebyshev distance

$$S(x, y) = \max_i (|x_i - y_i|)$$



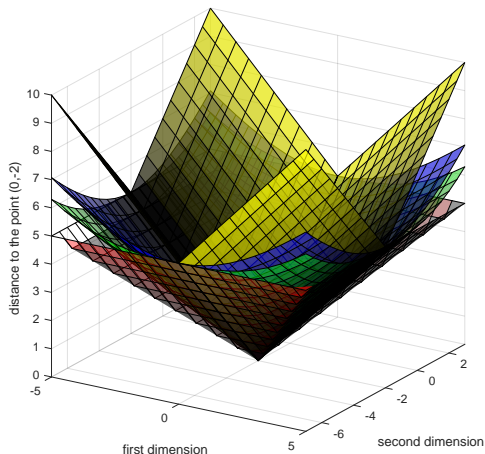
## Distance function: Examples 3 Minkowsky distance

$$S(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}$$

- $p < 1$  triangle inequality is violated, therefore for the values of  $p$  smaller than one, equation above is not a distance function.
- $p = 1$  case of Manhattan distance.
- $p = 2$  case of Euclidian distance.
- $p \rightarrow \infty$  case of Chebyshev distance.

## Distance function: Examples 4

3D representation of the Minkovski distances for different values of parameter  $p$ .  $p = 1$  - yellow surface, Manhattan;  $p = 2$  - blue surface, Euclidean,;  $p = 3$  - green surface;  $p \rightarrow \infty$  - red surface, Chebyshev.

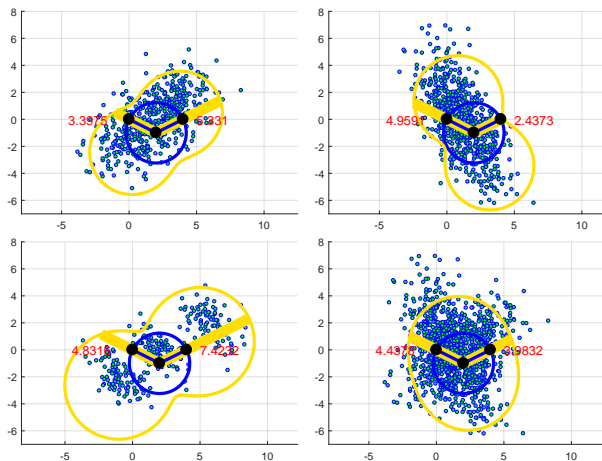


# Distance function: Examples 5

## Mahalanobis distance

$$S(x, y) = \sqrt{(x - y)^T C^{-1} (x - y)}$$

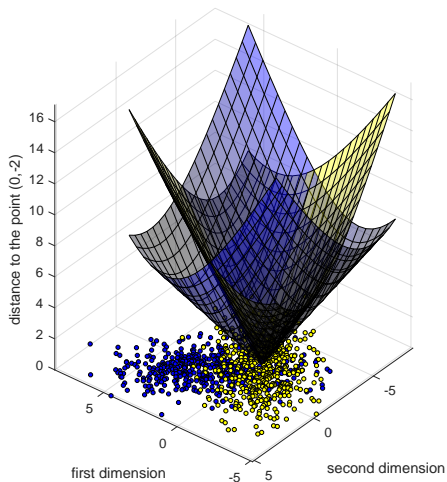
where  $C$  is the covariance matrix. Takes into account impact of data distribution.





## Distance function: Examples 6

- Impact of the rotation of underlying data set.



## Distance function: Examples 7

- Canberra distance

$$S(x, y) = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i| + |y_i|}$$

weighted version of Manhattan distance.

- Cosine distance Cosine similarity is the measure of the angle between two vectors

$$S_c(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

Usually used in high dimensional positive spaces, ranges from  $-1$  to  $1$ . Cosine distance is defined as follows

$$S_C(x, y) = 1 - S_c(x, y)$$

## Distance function: Examples 8

- Levenshtein or SED distance. SED - minimal number of single-character edits required to change one string into another. Edit operations are as follows:
  - ▶ insertions
  - ▶ deletions
  - ▶ substitutions
- $SED(\text{delta}, \text{delata})=1$  delete "a" or  $SED(\text{kitten}, \text{sitting})=3$  : substitute "k" with "s", substitute "e" with "i", insert "g".
- Hamming distance Similar to Levenshtein but with substitution operation only. Frequently used with categorical and binary data.
- Specialized similarity measures Distance and similarity functions applicable to the graphs, temporal data etc. These topics are left outside of the framework of the present course.

## Impact of High Dimensionality (Curse of Dimensionality)

*Curse of dimensionality* - term introduced by Richard Bellman. Referred to the phenomenon of efficiency loss by distance based data-mining methods. Let us consider the following example.

- Consider the unit cube in  $d$  - dimensional space, with one corner at the origin.
- What is the Manhattan distance from the arbitrary chosen point inside the cube to the origin?

$$S(\bar{0}, \bar{Y}) = \sum_{i=1}^d (Y_i - 0)$$

Note that  $Y_i$  is random variable in  $[0, 1]$

- The result is random variable with a mean  $\mu = d/2$  and standard deviation  $\sigma = \sqrt{d/12}$
- The ratio of the variation in the distances to the mean value is referred as *contrast*

$$G(d) = \frac{S_{max} - S_{min}}{\mu} = \sqrt{\frac{12}{d}}$$

# Data Preparation

Note: some steps listed below are frequently considered in pair of corresponding problem.

- **Feature Extraction.** Feature Extraction is the process of selecting the attributes and features relevant to the goal of analysis.
- **Data Cleaning.** Handling missing entries, handling incorrect and inconsistent entries, scaling and normalization.
- **Data Reduction and Transformation.** Sampling, feature subset selection dimensionality reduction. Conversion between different data types e.g. Numeric to Categorical data or Categorical to Numeric data.

# Conversion

- **Discretization:** Numeric to Categorical Data: Divide range of numeric attribute into finite number of intervals. To each data point assign categorical value of the interval containing its numerical attribute.
  - ▶ Equi-width ranges: the ranges have the same length.
  - ▶ Equi-log ranges:  $\log(b) - \log(a)$  have the same length for all the intervals. Here  $a$  is a beginning and  $b$  is the ending of the intervals.
  - ▶ Equi-depth ranges: Each range contains the same number of intervals.
- **Biniarization:** Categorical to Numeric Data
- Text to Numeric Data
- Time Series to Discrete Sequence Data
- Time series to Numeric Data
- Discrete Sequence to Numeric Data
- Spatial to Numeric Data
- Graphs to Numeric Data

# Data Cleaning

- Missing Entries
- Incorrect and Inconsistent Entries
- Scaling and Normalization: Different features represent different scales and not always comparable.
  - ▶ Normalization Let  $j^{th}$  attribute has mean  $\mu_j$  and standard deviation  $\sigma_j$  then  $j^{th}$  attribute value  $x_i^j$  of the record  $\bar{X}_i$  may be normalized as follows

$$z_i^j = \frac{x_i^j - \mu_j}{\sigma_j} \quad (1)$$

- ▶ Min - max scaling:

$$y_i^j = \frac{x_i^j - \min(x^j)}{\max(x^j) - \min(x^j)} \quad (2)$$

# Data Reduction and Transformation

- Sampling
  - ▶ Sampling for Static data
    - ★ Biased sampling
    - ★ Stratified sampling
  - ▶ Reservoir Sampling for Data Streams
- Feature Subset Selection
- Dimensionality Reduction
  - ▶ Principal Component Analysis
  - ▶ Singular Value Decomposition
  - ▶ Latent Semantic Analysis



## Example: Principal Component Analysis

Problem: Significant number of correlations may exist between different attributes. Usually used after the mean centering (subtracting the mean of the data set from each data point). The goal of PCA is to rotate the data into a coordinate system where the greatest amount of variance is captured in a smaller number of dimensions.

Let  $\mathcal{D}$  be  $n \times d$  data matrix and  $\mathcal{C}$   $d \times d$  covariance matrix. The covariance matrix  $\mathcal{C}$  is positive semi-definite. Since  $\mathcal{C}$  is positive semi-definite

$$\mathcal{C} = P\Lambda P^T \quad (3)$$

$P$  contains orthonormal eigenvectors of  $\mathcal{C}$  and diagonal matrix  $\Lambda$  - corresponding nonnegative eigenvalues.

# Principal Component Analysis

- Both eigenvectors and eigenvalues have a geometric interpretation.
- It may be shown that  $\binom{d}{2}$  covariances of transformed features are zero.
- Matrix  $\Lambda$  is the covariance matrix after axis rotations.
- Eigenvectors with large eigenvalues preserve greater variance and referred as principal components
- Transformed data matrix is computed as follows

$$\mathcal{D}' = \mathcal{D}P \quad (4)$$

# Exercises for self practice

Please note this is not a mandatory Home Assignment, nevertheless some or all of the exercises may be included into Home Assignments.

## Exercises

- Implement in MATLAB your own Minkowsky distance function.
- Implement in MALTAB your own Mahalanobis distance function.
- Implement in MATLAB your own Canberra distance functions.
- Implement in MATLAB the code to mimic Figures on the slides.

## Refresh in your memory the following notions:

Matrix, Matrix operations, inverse matrix, matrix determinant, positive define, probability, mean, standard deviation, covariance matrix, distribution.