Knowledge Discovery in Databases II

Lecture 2 – High Dimensional Data Prof. Dr. Peer Kröger, Yifeng Lu Sommer Semester 2019 Credits: Based on material of Eirini Ntoutsi, Matthias Schubert, Arthur Zimek, Peer Kröger, Yifeng Lu



- 1. Intorduction to Feature Spaces
- 2. Challenges of High Dimensional Data
- 3. Supervised Feature Selection
- 4. Feature Reduction and Metric Learning
- 5. Clustering High Dimensional Data

Kapitel 1: Intorduction i

LMU

- 1. Intorduction to Feature Spaces
- 2. Challenges of High Dimensional Data
- 3. Supervised Feature Selection
- 4. Feature Reduction and Metric Learning
- 5. Clustering High Dimensional Data

LMU

Feature Transform

- · Consider the following spaces:
 - $\,\mathbb{U}$ denotes the universe of data objects
 - $\mathbb{F} \subseteq \mathbb{R}^n$ denotes an *n*-dimensional feature space
- A feature transformation is a mapping *f* : U → ℝⁿ of objects from U to the feature space F.

Similarity Model

• A similarity model $S:\mathbb{U} imes\mathbb{U} o\mathbb{R}$ is defined for all objects $p,q\in\mathbb{U}$ as

$$S(p,q) = sim(f(p), f(q))$$

where $sim : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$ is a similarity measure or a dissimilarity (distance) measure in \mathbb{F} .



Comments:

- Often, dissimilarity (distance) is measured instead of similarity
- This is a small but important difference!
 - A similarity measure (sim) assigns high values to similar objects
 - · A dissimilarity measure (dist) assigns low values to similar objects
- The design of *f* and the definition of *sim/dist* are important assumptions about the patterns we want to find later in the data
- As explained before, *f* and *sim/dist* can be derived manually (explicit transformation and coding versus implicit Kernels) or automatically (representation learning)



- · Dissimilarity measures follow the idea of the geometric approach
 - objects are defined by their perceptual representations in a perceptual space
 - perceptual space = psychological space
 - geometric distance between the perceptual representations defines the (dis)similarity of objects
- · Within the scope of Feature-based similarity
 - perceptual space = feature space \mathbb{F} or feature representation space \mathbb{R}^n
 - geometric distance = distance function

Distance Functions

- The distance measure *dist* is a distance function if it is reflexive, non-negative, and symmetric
- A distance function *dist* is a metric if it additionally satisfies the triangle inequality
- · Comments:
 - · Sound mathematical interpretation
 - · Allow domain experts to model their notion of dissimilarity
 - · Metric distances allow to tune efficiency of data mining approaches
 - Long-lasting discussion of whether the distance properties and in particular the metric properties reflect the perceived dissimilarity correctly, see the following contradicting example:



Similarity versus Dissimilarity (again)

- Transformation
 - Let $\mathbb F$ be a feature space and $\textit{dist}:\mathbb F\times\mathbb F\to\mathbb R$ be a distance function
 - Any monotonically decreasing function $f : \mathbb{R} \to \mathbb{R}$ defines a similarity function $s : \mathbb{F} \times \mathbb{F} \to \mathbb{R}$ as follows

$$\forall x, y \in \mathbb{F} : s(x, y) = f(dist(x, y))$$

- Some prominent similarity functions $(x, y \in \mathbb{F})$:
 - exponential:
 s(x,y) = e^{(-dist(x,y))}
 - logarithmic: $s(x,y) = 1 - \log(1 + dist(x,y))$
 - linear: s(x, y) = 1 dist(x, y)





Similarities: Examples (only very few)

• Dot-Product $(x, y \in \mathbb{F} \subseteq \mathbb{R}^d)$

$$x \cdot y^{\mathsf{T}} = \sum_{i=1}^{d} x_i \cdot y_i = \|x\| \cdot \|y\| \cdot \cos \triangleleft(x,y)$$

Cosine
$$(x, y \in \mathbb{F} \subseteq \mathbb{R}^d)$$
$$\frac{x \cdot y^T}{\|x\| \cdot \|y\|}$$

• Pearson Correlation $(x, y \in \mathbb{F} \subseteq \mathbb{R}^d)$

$$\frac{\sum_{i=1}^{d} (x_i - \bar{x}_i) \cdot (y_i - \bar{y}_i)}{\sqrt{\sum_{i=1}^{d} (x_i - \bar{x}_i)^2} \cdot \sqrt{\sum_{i=1}^{d} (y_i - \bar{y}_i)^2}}$$

where \bar{z}_i denotes the mean in attribute *i* over all data points

- Random-Walk Kernel (for graphs x, y)
 - Count common (random) walks in x and y
 - · Walks are sequences of nodes (connected by edges)



Distances: Examples (only very few)



• L_p -norm (aka Minkowski metric) ($x, y \in \mathbb{F} \subseteq \mathbb{R}^d$)

$$L_p(x,y) = \sqrt[p]{\sum_{i=1}^d |x_x - y_i|^p}$$

where

- p < 1: fractional Minkowski distance
- p = 1: Manhattan distance
- p = 2: Euclidean distance
- *p* = ∞: Chebyshev/Maximum distance
- Malahanobis distance
- Hamming distance HammingDist(x,y) = $\sum_{i=1}^{d} \begin{cases} 1 : x_i \neq y_i \\ 0 : else \end{cases}$

Kapitel 2: Challenges i

LMU

1. Intorduction to Feature Spaces

2. Challenges of High Dimensional Data

- 3. Supervised Feature Selection
- 4. Feature Reduction and Metric Learning
- 5. Clustering High Dimensional Data

A Motivating Example

• Let's play the baby shapes game (truly motivating for students ...): Group the items!!!



- · What about grouping based on both shape and color?
- Lesson to learn: there may be different semantic concepts (and their corresponding patterns) hidden in the data (here: shape and color)



The good old days of data mining ...

- Data generation and, to some extend, data storage was costly (hard to imagine but those were the days ...)
- Domain experts carefully considered which features/variables to measure before designing experiments/a feature transform/...
- Consequence: also data sets were well designed and potentially contained only a small number of relevant features

Nowadays, data science is also about integrating everything

- · Generating and storing data is easy and cheap
- People tend to measure everything they can and even more (including even more complex feature transformations)
- The Data Science mantra is often interpreted as "we can analyze data from as many sources as (technically) possible, just record anaything you can"
- Consequence: data sets are high-dimensional containing a large number of features but the relevancy of each feature for the analysis goal is not clear a priori

High-dimensional Data is NOT a Myth

- Example: Image data
 - Low-level image descriptors (color histograms, textures, shape information ...)
 - Regional descriptors: between 16
 and 1,000 features
 - ...
- · Example: Metabolome data
 - Feature = concentration of one metabolite (intermediates/results of metabolism)
 - Bavaria newborn screening (for each baby, the blood concentrations of 43 metabolites are measured in the first 48 hours after birth)
 - between 50 and 2,000 features





More High-dimensional Data

- Example: Microarray data (deprecated)
 - · Features correspond to genes
 - · Up to 20,000 features
 - Dimensionality is much higher than the sample size
- · Example: Text data
 - Term frequency: features
 correspond to words/terms
 - Between 5,000 and 20,000 features (and even more)
 - Often, esp. in social media: abbreviations, colloquial language, special words

What's new at LMU? As usual, the most obvious change from last semester is this term's new crop of first-year students. - Around 8000 of them have arrived in Munich to begin their university careers. For the freshers themselves, of course, virtually everything is new - not just the lecture theaters, the professors and their classmates. Getting to know their new alma mater is their first priority. One of the many newcomers on campus is David Worofka, who is about to embark on a voyage around the bays and inlets of Economics. To ensure that he is well equipped to master the upcoming challenges, David has not only registered for LMU's P2P Mentoring Program but will also take the introductory orientation course (the so-called O Phase) offered by the Faculties of Economics and Business Administration, "For first-year students in particular, the Mentoring Program is a very good idea," he avers. Indeed, university studies are organized along very different lines from the more rigid schedules used in secondary schools and in much of the world of work. "Having a mentor on hand is a great help," he says. David's mentor, Alex Osberghaus, is well aware of how important it is to have someone to turn to for advice and assistance during the early phase of one's first semester: "In the beginning, when everything is unfamiliar, there are lots of questions to be answered," he says. "And mentors who already know the ropes can give their charges valuable tips that can help them to get off to a good start."

> Excerpt from LMU website: http://tinyurl.com/qhq6byz







Overview:

- Distances grow
- · Contrast of distances diminish (concentration problem)
- · Meaning of "neighborhood" concept
- · Growing data space
- · Growing hypothesis space
- · Empty spaces and importance tails
- · Different semantic layers
- ...

So let us have a closer look on these problems ...



The following example uses the Euclidean distance but holds for most distance measures:

- Consider 2D vectors a = (1,2) and b = (4,3)
- The Euclidean distance between *a* and *b* is

$$L_2(a,b) = L_2((1,2),(4,4))$$

= $\sqrt{(1-4)^2 + (2-3)^2}$
= $\sqrt{10}$



which corresponds to the norm of the difference vector c = (3, 1):

$$\|c\|_2 = \sqrt{3^2 + 1^2}$$



With increasing dimensionality, distances grow, too:

- Example: $L_2((1,2),(4,3)) = \sqrt{10}$
- Now double the feature vector length (double the original features): $L_2((1,2,1,2),(4,3,4,3)) = \sqrt{(3^2+1^2+3^2+1^2)} = \sqrt{20}$
- Effect seems not so important, values might be only in a larger scale?
- NOPE:

Contrast of distances is lost in high dimensional data since distances grow more and more alike!

This is know as the Concentration of Distances problem (see next)



Concentration Phenomenon

- As dimensionality grows, distance values grow, too, such that the (numerical) contrast provided by usual measures decreases or even diminishes
- In other words, the distribution of norms in a given distribution of points tends to concentrate
- Example: Euclidean norm of vectors consisting of several variables that are (assumed to be) independent and identically distributed

$$\|y\|_2 = \sqrt{y_1^2 + y_2^2 + \ldots + y_d^2}$$

• In high dimensional spaces this norm behaves unexpectedly ...



Theorem: Concentration of Distances

- Let *y* be a *d*-dimensional vector (*y*₁,..., *y_d*) where all components *y_i*(1 ≤ *i* ≤ *d*) are independent and identically distributed
- Then the mean and the variance of the Euclidean norm are:

$$\mu_{\parallel y \parallel} = \sqrt{a \cdot d - b} + \mathcal{O}(d^{-1})$$
 and $\sigma_{\parallel y \parallel} = b + \mathcal{O}(d^{-1/2})$

where *a* and *b* are parameters depending only on the central moments of order 1, 2, 3, 4.

Interpretation:

- The norm grows proportionally to \sqrt{d} , but the variance remains approx. constant for large *d* (because $\lim_{d\to\infty} d^{-const} = 0$)
- With growing dimensionality, the relative error made by taking $\mu_{\|y\|}$ instead of $\|y\|$ becomes negligible

⁰John A Lee and Michel Verleysen: "Nonlinear Dimensionality Reduction". Springer, 2007.



Implications from the concentration of distances:

- A lot of data mining methods use distances and neighborhoods to define patterns (e.g. *k*NN classifier, density-based clustering, distance-based outlier detection, ...
- Using neighborhoods is based on a key assumption:
 - Objects that are similar to an object *o* are in its neighborhood
 - · Object that are dissimilar to o are not in its neighborhood
- · What if all objects are in the same neighborhood?
 - Consider the above effect on distances: *k*NN distances are almost equal to each other, i.e., the *k* nearest neighbors are random objects

Definition: Unstable Neighborhood

A NN-query is unstable for a given ε if the distance from the query point to most data points is less than (1 + ε) times the distance from the query point to its nearest neighbor

 It can be shown that with growing dimensionality, the probability that a query is unstable converges to 1





Neighborhood Concept Become Meaningless

LMU

 Consider a *d*-dimensional query point *q* and *n d*-dimensional sample points *x*₁,...*x_n* (independent and identically distributed)



• We define:

 $DMIN_d = \min\{L_2(x_i, q) | 1 \le i \le n\}$ (dist to next neighbor) $DMAX_d = \max\{L_2(x_i, q) | 1 \le i \le n\}$ (dist to farthest neighbor)

Theorem

• If
$$\lim_{d\to\infty} \left(\frac{VAR_{L_2(x_i,q)}}{\mu_{L_2(x_i,q)}^2} \right) = 0$$

• Then $orall \epsilon > 0$: $\lim_{d o \infty} \mathbb{P}(\textit{DMAX}_d \leq (1 + \epsilon)\textit{DMIN}_d) = 1$

In other words: if the precondition holds, all points converge to the same distance from the query!

⁰Kevin S. Beyer, Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft: When is "nearest neighbor" meaningful? In ICDT 1999.



Visually: Pairwise distances of a sample of 105 instances drawn from a uniform [0,1] distribution, normalized $(1/\sqrt{d})$.



Neighborhood Concept Become Meaningless

LMU

- Be clear about the precondition of the Theorem!!!
- Consider the feature space of *d* **relevant** features for a given application (i.e., truly similar objects display small distances in most features)
- Now add *d* · *c* additional features being independent of the initial feature space
- With increasing *c* the distance in the independent subspace will dominate the distance in the complete feature space
- So the question is:

How many relevant features must be similar to indicate object similarity?

(or: how many relevant features must be dissimilar to indicate dissimilarity?)

• With increasing dimensionality the likelihood that two objects are similar in every respect gets smaller.

LMU

- OK, the data space grows with increasing dimensionality
- · But what are the problems?
- In low dimensional spaces we have some (intuitive) assumptions on the behavior of volumes (sphere, cube, etc.) and on the distribution of data objects
- However, basic assumptions do not hold in high dimensional spaces:
 - Spaces become sparse or even empty and the probability of one object inside a fixed range tends to become zero
 - Distribution of data has a strange behavior e.g. a normal distribution has only few objects in its center and the tails of distributions become more important

We will have a closer look on these issues ...

- · The more features, the larger the hypothesis space
- · The lower the hypothesis space is,
 - · the easier it is to find the correct hypothesis
 - · the less examples you need to properly test hypothesis



Growing Hypotheses Space

- Consider *f* a unit multivariate normal distribution and normal kernel (KDE)
- The aim is to find an estimate \hat{f} of f at the point 0
- The relative mean square error should be fairly small, e.g. $\frac{\mu_{\hat{f}(0)-f(0)}^2}{f(0)^2} < 0.1$

Dim.	Req. sample size to achieve 0.1 error estimate
1	4
2	19
5	768
8	43.700
10	842.000

Even with only 10 dimensions, we need nearly a million observations to estimate a distribution with an error less than 0.1!!!

⁰B.W. Silverman: "Density Estimation for Statistics and Data Analysis". Chapman and Hall/CRC, 1986.

Empty Spaces and Tails

- Consider a *d*-dimensional space with partitions of constant size 1/*m*
- The number of cells N increases exponentially in d: $N = m^d$
- Suppose *x* points are randomly placed in this space
- In low-dimensional spaces there are few empty partitions and many points per partitions
- In high-dimensional spaces there are far more partitions than points there are many empty partitions







Analogously:

- Consider a simple partitioning scheme, which splits the data in each dimension in 2 halves
- For *d* dimensions we obtain 2^{*d*} partitions
- Consider $n = 10^6$ samples in this space
- For $d \leq 10$ such a partition may make sense
- For d = 100 there are around 10^{30} partitions, so most partitions are empty (given the above 10^6 points)

Empty Spaces and Tails

- Consider a hyper-cube range query with length s in all dimensions, placed arbitrarily in the data space [0,1]^d
- *E* is the event that an arbitrary point lies within the query cube
- The probability for *E* is $\mathcal{P}(E) = s^d$



⇒ with increasing dimensionality, even very large hyper-cube range queries are not likely to contain a point







- The same holds of course for a spherical range query (instead of a cubical range query)
- Consequence: with increasing dimensionality the center of the hyper-cube (or more generally: of the data space) becomes less important and the volume of the data space concentrates in its corners (i.e. randomly distributed points tend to be on the border of the data space ...)
- This seems to be a distortion of space compared to our 3D way of thinking — and that is actually what it is ...

And that also means, that the tails of a distribution become extremely important

- Consider standard density function *f*
- Consider \hat{f} with

$$\hat{f}(x) = \begin{cases} 0 & f(x) < 0.07 \\ f(x) & \text{else} \end{cases}$$

• Rescaling \hat{f} to a density function will make very little difference in 1D, since very few data points occur in regions where *f* is very small







But for high dimensional data:

- More than half of the data has less then 1/100 of the maximum density *f*(0) (for μ = 0)
- Example: 10-dimensional Gaussian distribution X:

$$\frac{f(X)}{f(0)} = e^{(-\frac{1}{2}X^T X)} \approx e^{(-\frac{1}{2}\chi_{10}^2)}$$

since the median of the χ^2_{10} distribution is 9.34, the median of $\frac{f(X)}{f(0)}$ is $e^{\frac{-9.34}{2}} = 0.0094$

- · Thus, most objects occur at the tails of the distribution
- In other words, in contrast to the low dimensional case, regions of relatively very low density can be extremely important parts



But for high dimensional data:

- More than half of the data has less then 1/100 of the maximum density *f*(0) (for μ = 0)
- Example: 10-dimensional Gaussian distribution X:

$$\frac{f(X)}{f(0)} = e^{(-\frac{1}{2}X^T X)} \approx e^{(-\frac{1}{2}\chi^2_{10})}$$

since the median of the χ^2_{10} distribution is 9.34, the median of $\frac{f(X)}{f(0)}$ is $e^{\frac{-9.34}{2}} = 0.0094$

- · Thus, most objects occur at the tails of the distribution
- In other words, in contrast to the low dimensional case, regions of relatively very low density can be extremely important parts

Empty Spaces and Tails



Example: ($\mu = 0, \sigma = 1$)



- 1D: 90% of the mass of the distribution lies between -1.6 and 1.6
- 10D: 99% of the mass of the distribution is at points whose distance from the origin is greater than 1.6
- Thus, it is difficult to estimate the density, except for enormous samples because in very high dimensions virtually the entire sample will be in the tails



- Patterns and models on high-dimensional data are often hard to interpret, e.g. long decision rules
- Efficiency in high-dimensional spaces is often limited because e.g. index structures degenerate and distance computations are much more expensive
- There may be different semantic layers so pattern might only be observable in subspaces or projected spaces (cf. the baby shape game)
- · Cliques of correlated features dominate the object description



- Summarizing: the higher the dimensionality, the worse is the expected outcome of the mining algorithm (i.e., dimensionality is a curse, says Kröger)
- Well, not in general, the Kernel trick shows the opposite: through the extension of the data space with new attributes, the mining algorithm (e.g. a SVM classifier) gets more accurate (i.e., dimensionality is a blessing, says Tresp in his ML course)
- So: Who is right?????? Both What????



- Look at what we assumed for the curse: attributes are independent (and often even uniformly distributed)
- These attributes are likely to be irrelevant for the mining task
- And the blessing: a Kernel (if it works) adds relevant attributes (even more relevant than the original ones)
- Message: high-dimensional data is tricky and the curse can come by as several problems
 - Some are due to irrelevant attributes, so try to get rid of irrelevant attributes and keep the relevant ones
 - Some are instead of relevant attributes, so among the relevant attributes, try to get rid of redundant ones



- 1. Intorduction to Feature Spaces
- 2. Challenges of High Dimensional Data

3. Supervised Feature Selection

- 4. Feature Reduction and Metric Learning
- 5. Clustering High Dimensional Data

Feature Selection

- · A task to remove irrelevant and/or redundant features
 - Irrelevant features:
 - · Not useful for a given task
 - · Probably decrease accuracy
 - Redundant features:
 - · Strongly correlated with another relevant feature
 - · Does not drop the accuracy, but may drop efficiency, explainability, etc.
- Deleting irrelevant and redundant features can improve the quality as well as the efficiency of the methods and the found patterns.
- New feature space: Delete all useless features from the original feature space.

Keep in mind...

Feature selection \neq Dimensionality reduction Feature selection \neq Feature extraction





Feature y is irrelevant, because if we omit x, we have only one cluster, which is uninteresting.



Features x and y are redundant, because x provides (appr.) the same information as feature y with regard to discriminating the two clusters

Prof. Dr. Peer Kröger: KDD2 (SoSe 2019) — Lecture 2 – High Dimensional Data — 3. Feature Selection

⁰ Source: Feature Selection for Unsupervised Learning, Dy and Brodley, Journal of Machine Learning Research 5 (2004)





Feature y separates well the two classes. Feature x is irrelevant. Its addition "destroys" the class separation.

Individually irrelevant together relevant



Redundancy



Features x_1 and x_2 are redundant.

O Source: http://www.kdnuggets.com/2014/03/machine-learning-7-pictures.html

Problem Definition

- Input: Vector space $F = d_1 \times \cdots \times d_n$, dimensions $D = \{d_1, \ldots, d_n\}$.
- **Output:** a minimal subspace *M* over dimensions *D*′ ⊆ *D* which is optimal for a given data mining task.
 - Minimality increases the efficiency, reduces the effects of the curse of dimensionality and increases interpretability.

Challenges:

- Optimality depends on the given task.
- There are 2^d possible solution spaces (exponential complexity)
- This search space is similar to the frequent itemset mining problem, but:
 - There is often no monotonicity in the quality of subspace (which is important for efficient searching)
 - Features might only be useful in combination with other certain features.
- \Rightarrow For many popular criteria, feature selection is an exponential problem.
- \Rightarrow Most algorithms employ search heuristics.



- 1. Feature subset generation
 - Single dimensions
 - · Combinations of dimensions (subspaces)
- 2. Feature subset evaluation
 - Importance scores like information gain, χ^2
 - Performance of a learning algorithm
- \Rightarrow How to select/evaluate features? How to traverse the search space?



- 1. Filter methods
 - Explores the general characteristics of the data, independent of the learning algorithm.
- 2. Wrapper methods
 - The learning algorithm is used for the evaluation of the subspace.
- 3. Embedded methods
 - The feature selection is part of the learning algorithm.



- · Filter methods
 - Basic idea: assign an "importance" score to each feature to filter out useless ones
 - Examples: information gain, χ^2 -statistic, TF-IDF for text...
 - Disconnected from the learning algorithm.
 - Pros:
 - Fast and generic
 - o Simple to apply
 - Cons:
 - Doesn't take into account interactions between features
 - Individually irrelevant features, might be relevant together
 - Too generic?



- · Wrapper methods
 - A learning algorithm is employed and its performance is used to determine the quality of selected features.
 - Pros:
 - o take feature dependencies into account
 - $\circ\;$ interaction between feature subset search and model selection
 - Cons:
 - o higher risk of overfitting than filter techniques
 - very computationally intensive, especially if building the classifier has a high computational cost.



- · Embedded methods
 - Such methods integrate the feature selection in model building
 - Example: decision tree induction algorithm: at each decision node, a feature has to be selected.
 - Pros:
 - $\circ~$ less computationally intensive than wrapper methods.
 - Cons:
 - specific to a learning method



- Forward selection
 - Start with an empty feature space and add relevant features
- Backward selection
 - Start with all features and remove irrelevant features
- · Branch-and-bound
 - Find the optimal subspace under the monotonicity assumption
- Randomized
 - Randomized search for a k dimensional subspace

• ...