




TU HCI-KDD

Andreas Holzinger
185.A83 Machine Learning for Health Informatics
2017S, VU, 2.0 h, 3.0 ECTS
Lecture 07 - Module 05 – Week 18 - 02.05.2017

Dimensionality Reduction and Subspace Clustering: Example for the Doctor-in-the-Loop

a.holzinger@hci-kdd.org
<http://hci-kdd.org/machine-learning-for-health-informatics-course>




Holzinger Group, hci-kdd.org 1 Machine Learning Health 07

TU HCI-KDD

Red thread through this lecture

- 01 Classification vs Clustering
- 02 Feature spaces, feature engineering
 - Feature selection, feature extraction
- 03 The curse of dimensionality
- 04 Dimensionality reduction
 - PCA, ICA, FA, MDS, LDA – Isomap, LLE, Autoencoder
- 05 Subspace clustering and analysis
- 06 Projection Pursuit: “What is interesting?”



Holzinger Group, hci-kdd.org 4 Machine Learning Health 07

TU HCI-KDD

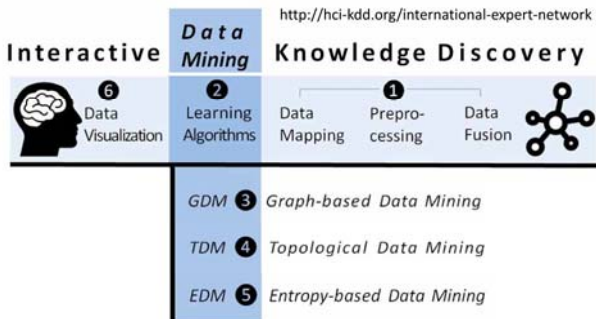
01 Classification vs. Clustering

Holzinger Group, hci-kdd.org 7 Machine Learning Health 07

TU ML needs a concerted effort fostering integrated research HCI-KDD

<http://hci-kdd.org/international-expert-network>

Interactive Data Mining Knowledge Discovery




6 Data Visualization
2 Learning Algorithms
3 GDM Graph-based Data Mining
4 TDM Topological Data Mining
5 EDM Entropy-based Data Mining

Privacy, Data Protection, Safety and Security

Holzinger, A. 2014. Trends in Interactive Knowledge Discovery for Personalized Medicine: Cognitive Science meets Machine Learning. IEEE Intelligent Informatics Bulletin, 15, (1), 6-14.

Holzinger Group, hci-kdd.org 2 Machine Learning Health 07

TU HCI-KDD



Holzinger Group, hci-kdd.org 5 Machine Learning Health 07

TU Key Challenges HCI-KDD

- Uncertainty, Validation, Curse of Dimensionality
- Large spaces gets sparse
- Distance Measures get useless
- Patterns occur in different subspaces
- Most pressing question “What is interesting?”

Holzinger Group, hci-kdd.org 8 Machine Learning Health 07

TU ML-Jungle Top Level View HCI-KDD



Maths Cognition Visualization Data structure Challenges
Perception Preprocessing
Decision Interaction Integration

Always with a focus/application in health informatics

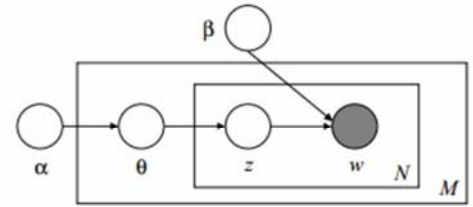
CONCEPTS THEORIES PARADIGMS MODELS METHODS TOOLS

Curse of Dim Bayesian p(x) unsupervised Gaussian P. Regularization Python
NFL-Theorem DR Complexity supervised Graphical M. Validation Julia
Overfitting KL-Divergence Semi-supv. NN DL Aggregation Etc.
Non-Parametric Info Theory online SVM Nature Inspired Azure
Exp. & Eval. iML Linear Models Privacy ML
RL PL AL D. Trees

Holzinger Group, hci-kdd.org 3 Machine Learning Health 07

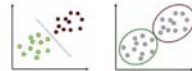
TU LDA = ? HCI-KDD

- Latent Dirichlet Allocation
- LDA = linear discriminant analysis (Attention!)



Holzinger Group, hci-kdd.org 6 Machine Learning Health 07

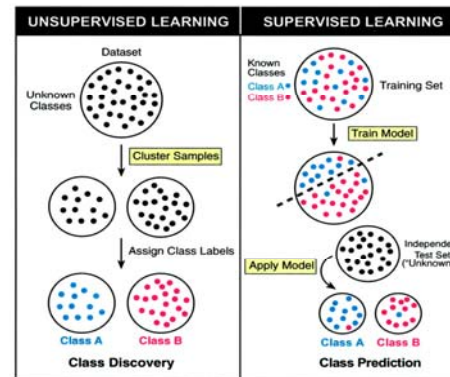
TU Classification (A) vs. Clustering (C) – Intro Quiz HCI-KDD

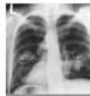
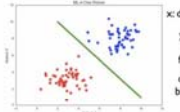


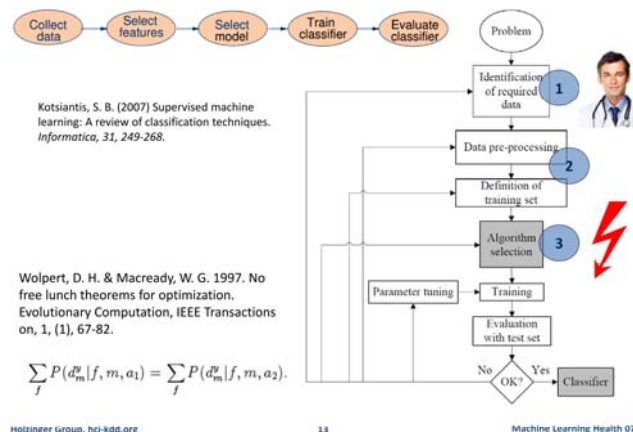
- The data is not labeled (A/C)?
- Identify structure/patterns (A/C)?
- Predicting an item set, identifying to which set of categories a new observation belongs (A/C)?
- Assigning a set of objects into groups (A/C)?
- Having many labelled data points (A/C)?
- Using the concept of supervised learning (A/C)?
- Grouping data items close to each other (A/C)?
- Used to explore data sets (A/C)?

Holzinger Group, hci-kdd.org 9 Machine Learning Health 07

- **Classification (Supervised learning, Pattern Recogn., Prediction)**
 - Supervision = the training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations;
 - New data is **classified based on the training set**
 - Important for clinical decision making
 - Example: Benign/Malign Classification of Tumors
- **Clustering (Unsupervised learning, class discovery,)**
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of **establishing the existence of clusters** in the data;



-  C_1 : Cancer present
 C_2 : Cancer absent
- x – set of pixel intensities
-  x : data points
 y : labels
features
decision boundary
- Typical questions include:
 - Is this protein functioning as an enzyme?
 - Does this gene sequence contain a splice site?
 - Is this melanoma malign?
 - Given object x – predict the class label y
 - If $y \in \{0,1\} \rightarrow$ binary classification problem
 - If $y \in \{1, \dots, n\}$ and is $n \in \mathbb{N} \rightarrow$ multiclass problem
 - If $y \in \mathbb{R} \rightarrow$ regression problem

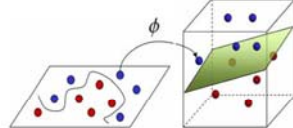


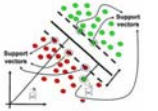
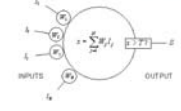
- Naïve Bayes (NB) – see Bayes' theorem with independent assumptions (hence “naïve”)
- Decision Trees (e.g. C4.5)
- NN – if x_1 is most similar to $x_2 \Rightarrow y_1 = y_2$

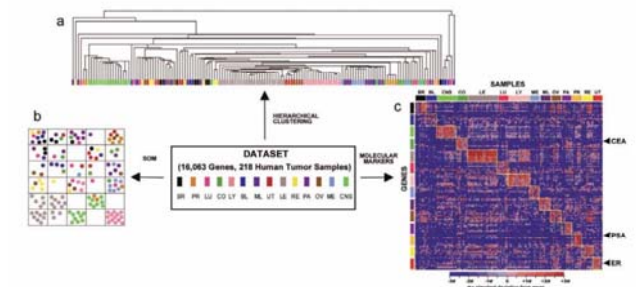
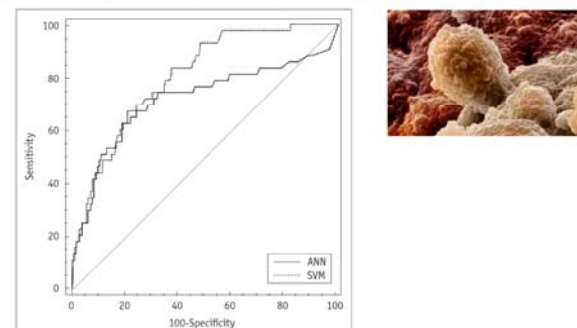
$$x_j = \operatorname{argmin}_{x \in D} \|x - x_i\|^2 \Rightarrow y_i = y_j$$
- SVM – a plane/hyperplane separates two classes of data – very versatile for classification and clustering – also via the Kernel trick in high-dimensions

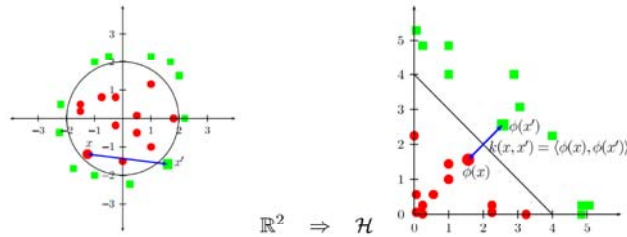
```

1: Input:  $(x_1, y_1), \dots, (x_n, y_n), C, \epsilon$ 
2:  $S_i \leftarrow \emptyset$  for all  $i = 1, \dots, n$ 
3: repeat
4:   for  $i = 1, \dots, n$  do
5:      $H(y) \leftarrow \Delta(y, y_i) + w^T \phi(x_i, y) - w^T \phi(x_i, y_i)$ 
6:     compute  $\hat{y} = \operatorname{argmax}_{y \in C} H(y)$ 
7:     compute  $\xi_i = \max\{0, \max_{y \in C} H(y)\}$ 
8:     if  $H(\hat{y}) > \xi_i + \epsilon$  then
9:        $S_i \leftarrow S_i \cup \{y\}$ 
10:     $w \leftarrow$  optimize primal over  $S = \bigcup_i S_i$ 
11:   end if
12: end for
13: until no  $S_i$  has changed during iteration
    
```

- Uses a nonlinear mapping to transform the original data (input space) into a higher dimension (feature space)
- 
- = classification method for both linear and nonlinear data;
 - Within the new dimension, it searches for the linear optimal separating hyperplane (i.e., “decision boundary”);
 - By nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated with a hyperplane;
 - The SVM finds this hyperplane by using **support vectors** (these are the “essential” training tuples) and **margins** (defined by the support vectors);

- **SVM**
 - Deterministic algorithm
 - Nice generalization properties
 - Hard to learn – learned in batch mode using quadratic programming techniques
 - Using kernels can learn very complex functions
 - **ANN**
 - Nondeterministic algorithm
 - Generalizes well but doesn't have strong mathematical foundation
 - Can easily be learned in incremental fashion
 - To learn complex functions—use multilayer perceptron (nontrivial)
-  





Borgwardt, K., Gretton, A., Rasch, J., Kriegel, H.-P., Schölkopf, B. & Smola, A. 2006. Integrating structured biological data by kernel max. mean discrepancy. *Bioinformatics*, 22, 14, e49-e57.

- Partite a data set into k clusters so that intra-cluster variance is a minimum
 - V ... variance (objective function)
 - S_i ... cluster
 - Y_i ... mean
 - D ... set of all points x_j
 - k ... number of clusters

$$V(D) = \sum_{i=1}^k \sum_{x_j \in S_i} (x_j - \mu_i)^2$$

- What is the computational time of k-means?
- NP-hard in Euclidean space, however, if k and d can be fixed then it can be solved within:

$\mathcal{O}(npkt)$
compute kn distances
in p dimensions

number of iterations
Can be small if there's
indeed a cluster
structure in the data

Jain, A. K. 2010. Data clustering: 50 years beyond K-means. *Pattern recognition letters*, 31, (8), 651-666.

- C4.5**
 - for generation of decision trees used for **classification**, (statistical classifier, Quinlan (1993));
- k-means**
 - simple iterative method for partition of a dataset in a user-specified n of **clusters**, k (Lloyd (1957));
- Apriori**
 - for finding frequent item sets using candidate generation and **clustering** (Agrawal & Srikant (1994));
- EM**
 - Expectation-Maximization algorithm for finding maximum likelihood estimates of parameters in models (Dempster et al. (1977));
- PageRank**
 - a search ranking algorithm using hyperlinks on the Web (Brin & Page (1998));
- Adaptive Boost**
 - one of the most important ensemble methods (Freund & Shapire (1995));
- k-Nearest Neighbor**
 - a method for **classifying** objects based on closest training sets in the feature space (Fix & Hodges (1951));
- Naive Bayes**
 - can be trained efficiently in a supervised learning setting for classification (Domingos & Pazzani (1997));
- CART**
 - Classification** And Regression Trees as predictive model mapping observations about items to conclusions about the goal (Breiman et al 1984);
- SVM** support vector machines offer one of the most robust and accurate methods among all well-known algorithms (Vapnik (1995));

Algorithm 1: Example for a classical weight balanced k -means algorithm

Input: $d, k, n \in \mathbb{N}$, $X := \{x_1, \dots, x_n\} \subset \mathbb{R}^d$, $S := \{s_1, \dots, s_k\} \subset \mathbb{R}^d$

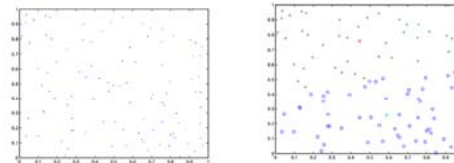
Output: Clustering $C = (C_1, \dots, C_k)$ of X and the arithmetic means c_1, \dots, c_k as sites

- Partition X into a clustering $C = (C_1, \dots, C_k)$ by assigning $x_j \in X$ to a cluster C_i that is closest to site $s_i \in S$.
- Update each site s_i as the center of gravity of cluster C_i ; if $|C_i| = 0$, choose $s_i = x_i$ for a random $i \leq n$ with $x_i \neq s_j$ for all $j \leq k$. If the sites change, go to (1..).

Merely an increase in awareness of physicians on risk factors for ARA in children can be sufficient to change their attitudes towards antibiotics prescription. Our results can also be useful when preparing recommendations for antibiotics prescription and to guide the standardized health data record.



Yildirim, P., Majnarić, L., Ekmekci, O. I. & Holzinger, A. 2013. On the Prediction of Clusters for Adverse Reactions and Allergies on Antibiotics for Children to Improve Biomedical Decision Making. In: Lecture Notes in Computer Science LNCS 8127. 431-445

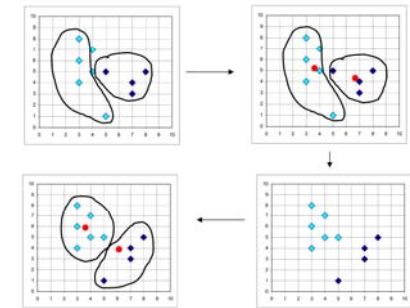
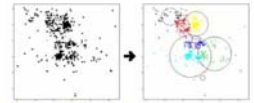


- Centroid:** mean of the points in the cluster.

$$\mu = \frac{1}{|C|} \sum_{x \in C} x$$

- Medoid:** point in the cluster that is closest to the centroid.
 $m = \arg \min_{x \in C} d(x, \mu)$

- Group similar objects into clusters together, e.g.
 - For image segmentation
 - Grouping genes similarly affected by a disease
 - Clustering patients with similar diseases
 - Cluster biological samples for category discovery
 - Finding subtypes of diseases
 - Visualizing protein families
- Inference: given x_i , predict y_i by learning f
- No training data set – learn model and apply it



02 Feature Engineering

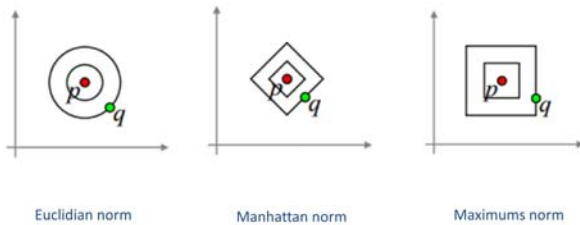


"Applied ML is basically feature engineering.
Andrew Yan-Tak Ng"

- Feature:= specific measurable property of a phenomenon being observed.
- Feature engineering:= using domain knowledge to create features useful for ML. (**"Applied ML is basically feature engineering. Andrew Ng"**).
- Feature learning:= transformation of raw data input to a representation, which can be effectively exploited in ML.

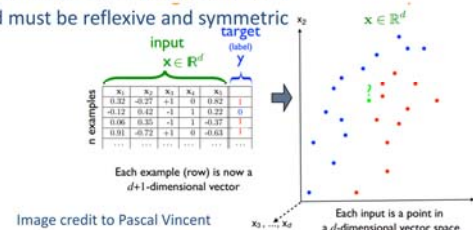
Let do a Quiz again: Similarities of feature vectors

Look at the examples below, which distance measures would you select?



03 Curse of Dimensionality

- Intuitively: a domain with a distance function
- Formally: Feature Space $\mathcal{F} = (\mathcal{D}, d)$
 - \mathcal{D} = ordered set of features
 - $d: \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}_0^+$... a total distance function; true for
 - $\forall p, q \in \mathcal{D}, p \neq q: d(p, q) > 0$ (strict)
 - and must be reflexive and symmetric



Feature Selection: Overview

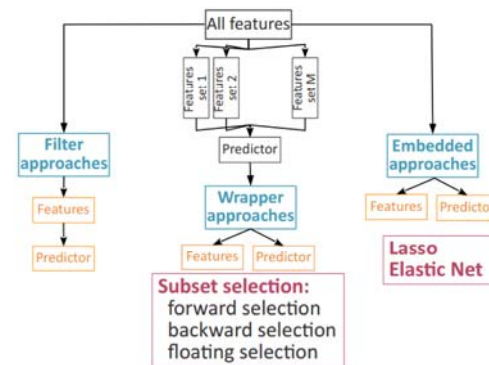
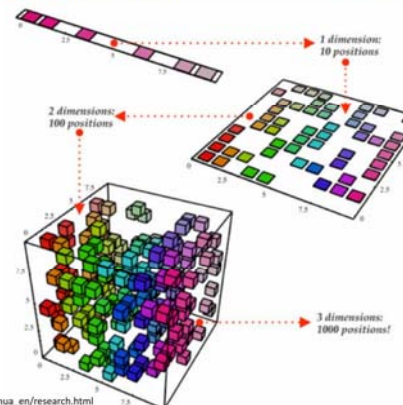


Image credit to Chloe Azencott

Remember: The curse of dimensionality



Bengio, S. & Bengio, Y. 2000. Taking on the curse of dimensionality in joint distributions using neural networks. IEEE Transactions on Neural Networks, 11, (3), 550-557.

http://www.imo.umontreal.ca/~bengio/yoshua_en/research.html

A **Metric Space** is a pair (X, d) where X is a set and $d: X \times X \rightarrow \mathbb{R}^+$, called the metric, s.t.

- For all $x, y, z \in X$, $d(x, y) \leq d(x, z) + d(z, y)$.
- For all $x, y \in X$, $d(x, y) = d(y, x)$.
- $d(x, y) = 0$ if and only if $x = y$.

Remark 1. One example is \mathbb{R}^d with the Euclidean metric. Spheres S^n endowed with the spherical metric provide another example.

$$d: \mathcal{X} \rightarrow \mathbb{R}$$

$$d(x, x) = 0$$

$$d(x^1, x^2) = d(x^2, x^1) \text{ symmetry}$$

$$d(x^1, x^2) \leq d(x^1, x^3) + d(x^3, x^2) \text{ triangle inequality}$$

Feature Selection vs. Feature Extraction

- Feature selection is just selecting a subset of the existing features without any transformation
- Feature extraction is *transforming* existing features into a lower dimensional space

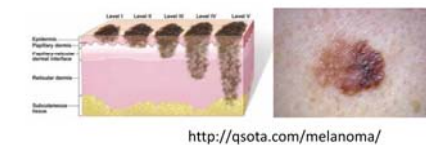
$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \xrightarrow{\text{feature selection}} \begin{bmatrix} x_{i_1} \\ x_{i_2} \\ \vdots \\ x_{i_M} \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \xrightarrow{\text{feature extraction}} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} = f \left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \right)$$

Blum, A. L. & Langley, P. 1997. Selection of relevant features and examples in machine learning. Artificial intelligence, 97, (1), 245-271.

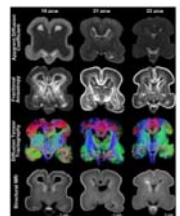
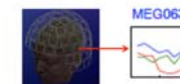
Examples for High-Dimensional Biomedical Data

- Medical Image Data (16 - 1000+ features)



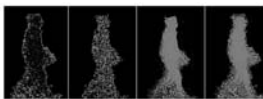
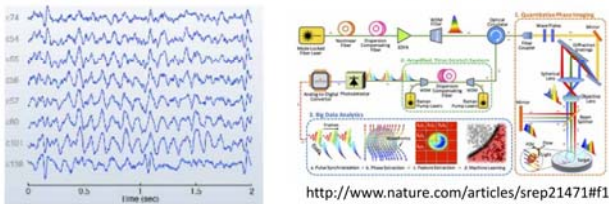
MEG Brain Imaging

120 locations x 500 time points
x 20 objects



Nature 508, 199–206
doi:10.1038/nature13185

Biomedical Signal Data (10 - 1000+ features)



Holzinger Group, hci-kdd.org



Machine Learning Health 07

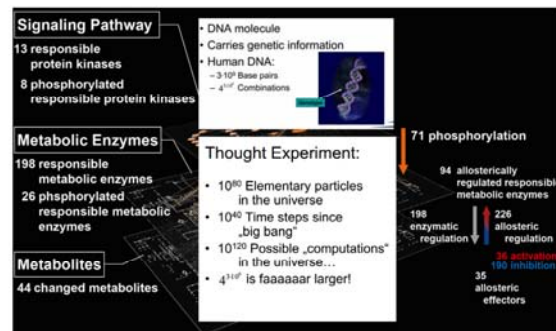
- Text > 10^9 documents \times 10^6 words/n-grams features correspond to words or terms, between 5k to 20k features

- Text (Natural Language) is definitely very important for health:

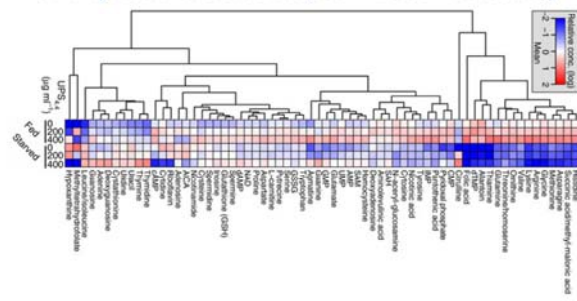
- Handwritten Notes, Drawings
- Patient consent forms
- Patient reports
- Radiology reports
- Voice dictations, annotations
- Literature !!!



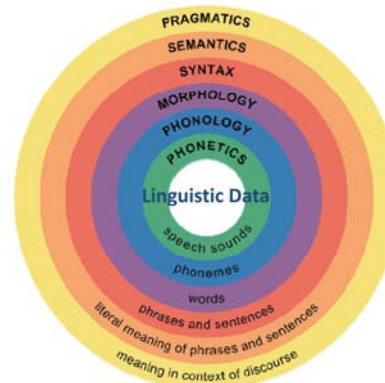
https://www.researchgate.net/publication/255723699_An_Answer_to_Who_Needs_a_Stylus_o_n_Handwriting_Recognition_on_Mobile_Devices
Holzinger Group, hci-kdd.org



Metabolome data (feature is the concentration of a specific metabolite; 50 – 2000+ features)



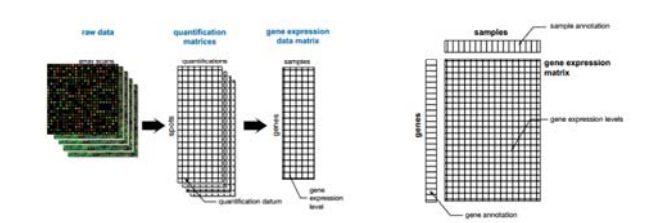
Holzinger Group, hci-kdd.org



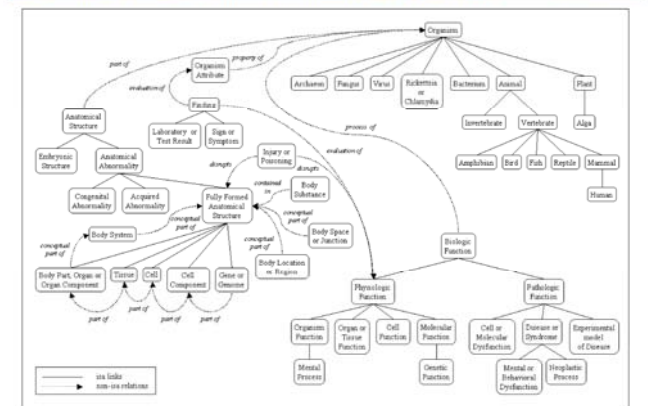
Holzinger Group, hci-kdd.org

- Hyperspace is large – all points are far apart
- Computationally challenging (in time and space)
- Complexity grows with n of features
- Complex models less robust – more variance
- Statistically challenging – hard to learn
- Hard to interpret and hard to visualize
- Problem with redundant features and noise
- Question: Which algorithms will provide worse results with increasing irrelevant features?
- Answer: Distance-based algorithms generally trust all features of equal importance

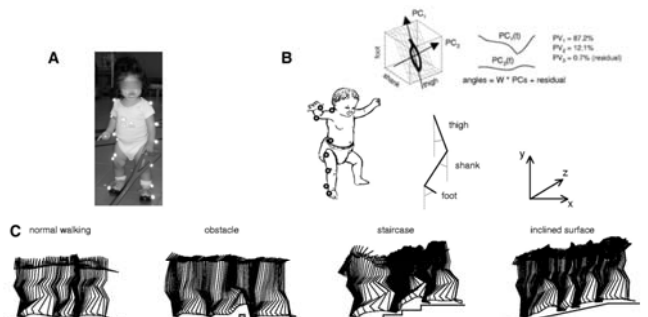
Microarray Data (features correspond to genes, up to 30k features)



Holzinger Group, hci-kdd.org

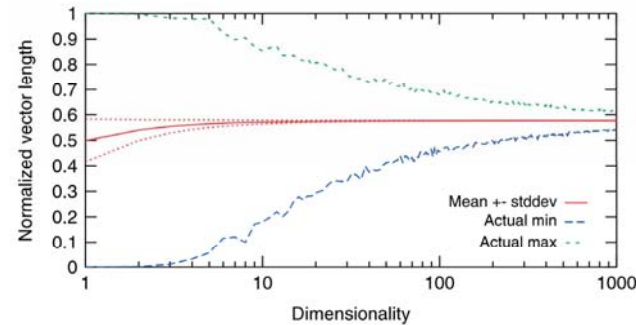


Holzinger Group, hci-kdd.org



- Aspect 1: Optimization Problem
- Aspect 2: Concentration Effect
- Aspect 3: Irrelevant Attributes
- Aspect 4: Correlated Attributes

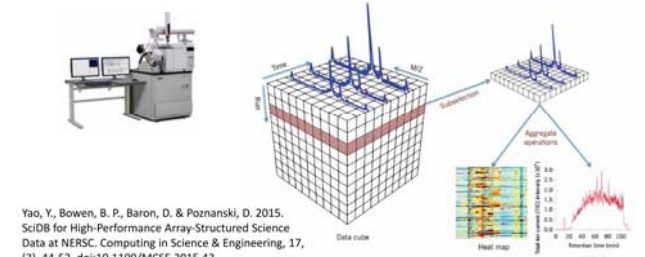
Kriegel, H. P., Kröger, P. & Zimek, A. 2012. Subspace clustering. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2, (4), 351-364, doi:10.1002/widm.1057.



Zimek, A., Schubert, E. & Kriegel, H. P. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5, (5), 363-387, doi:10.1002/sam.11161.



Fourteen amino acids and 29 fatty acids are analysed from a single blood spot using MS/MS. The concentrations are given in $\mu\text{mol/L}$.



Yao, Y., Bowen, B. P., Baron, D. & Poznanski, D. 2015. SciDB for High-Performance Array-Structured Science Data at NERSC. *Computing in Science & Engineering*, 17 (3), 44-52, doi:10.1109/MCSE.2015.43.

04 Dimensionality Reduction

- Data visualization only possible in \mathbb{R}^2 (\mathbb{R}^3 cave)
- Human interpretability only in $\mathbb{R}^2/\mathbb{R}^3$ (visualization can help sometimes with parallel coordinates)
- Simpler (=less variance) models are more robust
- Computational complexity (time and space)
- Eliminate non-relevant attributes that can make it more difficult for algorithms to learn
- Bad results through (many) irrelevant attributes?
- *Note again: Distance-based algorithms generally trust that all features are equally important.*

- Linear methods (unsupervised):
 - PCA
 - FA
 - MDS
- Supervised methods:
 - LDA
- Non-linear methods (unsupervised):
 - Isomap (Isometric feature mapping)
 - LLE (locally linear embedding)
 - Autoencoders

TU Why should we reduce the dimensionality?

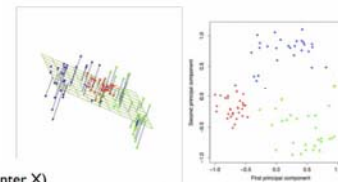
- Given n data points in d dimensions
 - Conversion to m data points in $r < d$ dimensions
 - Challenge: **minimal loss of information *)**
- *) this is always a grand challenge, e.g. in k-Anonymization – see later in this
 - Very dangerous is the “modeling-of-artifacts”

- *) this is always a grand challenge, e.g. in k-Anonymization – see later in this

- Very dangerous is the “modeling-of-artifacts”

TU Approaches **HCI-KDD**

Example 1: PCA



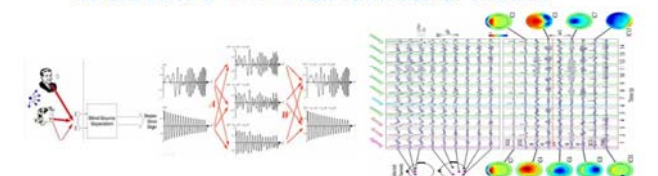
- Subtract mean from data (center X)
- (Typically) scale each dimension by its variance
 - Helps to pay less attention to magnitude of dimensions
- Compute covariance matrix $S = \frac{1}{N} X^T X$
- Compute k largest eigenvectors of S
- These eigenvectors are the k principal components

Hastie, T., Tibshirani, R. & Friedman, J. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Second Edition. New York, Springer. doi:10.1007/978-0-387-84858-7.

Example 2 ICA (Motivation: Blind Source Separation)

- Suppose that there are k unknown independent sources

$$\mathbf{s}(t) = [s_1(t), \dots, s_k(t)]^T \text{ with } E\mathbf{s}(t) = \mathbf{0}$$
- A data vector $\mathbf{x}(t)$ is observed at each time point t , such that $\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t)$
 where \mathbf{A} is a $n \times k$ full rank scalar matrix



Holzinger, A., Scherer, R., Seeber, M., Wagner, J. & Müller-Putz, G. 2012. Computational Sensemaking on Examples of Knowledge Discovery from Neuroscience Data: Towards Enhancing Stroke Rehabilitation. In: Böhm, C., Khuri, S., Lhotská, L. & Renda, M. (eds.) *Information Technology in Bio- and Medical Informatics, Lecture Notes in Computer Science, LNCS 7451*. Heidelberg, New York: Springer, pp. 166-168

- Factor Analysis describes the variability of observations in terms of unobserved latent variables (these are called "factors") and noise
 - The factors explain the correlation between the var
- Variance can be explained by Gaussian noise (and can be calculated)
- Advantage: generative approach and models BOTH the noise of the observations and their correlation!
- You can make assumptions on the distributions of noise and factors

A Global Geometric Framework for Nonlinear Dimensionality Reduction

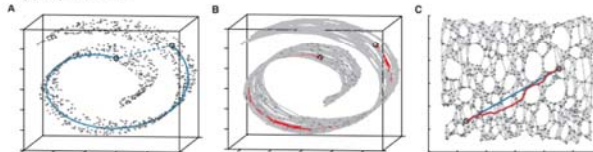
Joshua B. Tenenbaum,^{1,2} Vin de Silva,³ John C. Langford⁴

Scientists working with large volumes of high-dimensional data, such as global climate patterns, stellar spectra, or human gene distributions, regularly confront the problem of dimensionality reduction: finding meaningful low-dimensional structures hidden in their data. In this paper, we describe a new approach to this problem: we view the data as points on a manifold, and we use the global geometry of the manifold to find a low-dimensional representation. This approach is based on the idea of "isomap" (isometric manifold projection), and it is capable of discovering the nonlinear degrees of freedom that underlie complex natural phenomena. In contrast to previous algorithms for nonlinear dimensionality reduction, isomap is a global optimization problem, and, for an important class of data manifolds, is guaranteed to converge asymptotically to the true structure.

Goal: Find projection onto nonlinear manifold

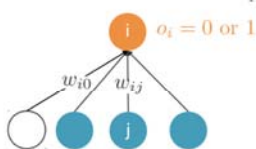
- Construct neighborhood graph G :
For all x_i, x_j
If $\text{distance}(x_i, x_j) < \epsilon$
Then add edge (x_i, x_j) to G
- Compute shortest distances along graph $\delta_{ij}(x_i, x_j)$ (e.g., by Floyd's algorithm)
- Apply multidimensional scaling to $\delta_{ij}(x_i, x_j)$

<http://isomap.stanford.edu/>



- Based on Information processing in dynamical systems: Foundations of harmony theory by Smolensky (1986): Stochastic neural networks where the unit activation i = probabilistic

$$Pr(o_i = 1) = \frac{1}{1 + e^{-w_{i0} + \sum_j w_{ij} o_j}}$$



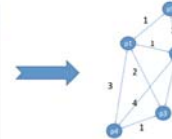
Right: A restricted Boltzmann machine with binary hidden units and softmax visible units

Salakhutdinov, R., Mnih, A. & Hinton, G. (2007) Restricted Boltzmann machines for collaborative filtering. ICML, 791-798.

Find a set of points whose pairwise distances match a given distance matrix

- Given $n \times n$ matrix of pairwise distances between data points
- Compute $n \times k$ matrix X with coordinates of distances with some linear algebra magic
- Perform PCA on this matrix X

	p1	p2	p3	p4	p5
p1	0	1	2	3	1
p2	1	0	2	4	1
p3	2	2	0	1	3
p4	3	4	1	0	1
p5	1	1	3	1	0



x_i Point in d dimensions
 y_i Corresponding point in $r < d$ dimensions
 δ_{ij} Distance between x_i and x_j
 d_{ij} Distance between y_i and y_j

- Define (e.g.) $E(y) = \sum_{i,j} (d_{ij} - \delta_{ij})^2$
- Find y_i 's that minimize E by gradient descent
- Invariant to translations, rotations and scalings



Compact representation of input

$$\min_{f,g} \sum_x \Delta(f \circ g, x)$$

- History: Dim-reduction with NN: Learning representations by back-propagating errors
- Goal: output matches input

Rumelhart, D. A., Hinton, G. E. & Williams, R. J. 1986. Learning representations by back-propagating errors. Nature, 323, 533-536.

Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y. & Manzagol, P.-A. 2010. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. The Journal of Machine Learning Research, 11, 3371-3408.

- Goal: Having $m < p$ features
- Feature selection via
 - A) Filter approaches
 - B) Wrapper approaches
 - C) Embedded approaches (Lasso, Electric net, see Tibshirani, Hastie ...)
- Feature extraction
 - A) Linear: e.g. PCA
 - B) Non-linear: Autoencoders (map the input to the output via a smaller layer)

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here, two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using comparative analyses to compare known genomes, concluded that today's *Arabidopsis*, on its own, is sustained with just 250 genes, and that the earliest life forms required a mere 125 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 80 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, these *predictions* are a step toward understanding the minimum number of genes needed for life.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down: Computer analysis yields an estimate of the minimum number of genes needed for life.

SCIENCE • VOL. 272 • 24 MAY 1996

- Sigmoidal neurons and backpropagation: Rumelhart*, D. A., Hinton, G. E. & Williams, R. J. 1986. Learning representations by back-propagating errors. Nature, 323, 533-536.

$$\Delta(y, x) = \|y - x\|_2^2$$

- Linear autoencoders: Baldi, P. & Hornik, K. 1989. Neural networks and principal component analysis: Learning from examples without local minima. Neural networks, 2, (1), 53-58.

$$\min_{A,B} \sum_x \|ABx - x\|_2^2$$

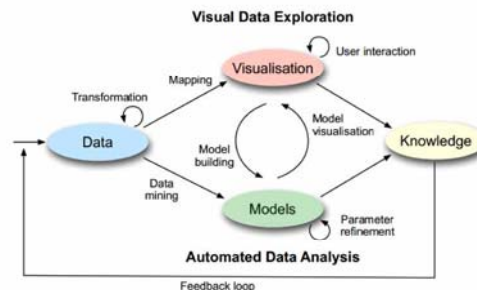
*) David Rumelhart (1942-2011) was Cognitive Scientist working on math. Psychology

05 Subspace Clustering* & Subspace Analysis

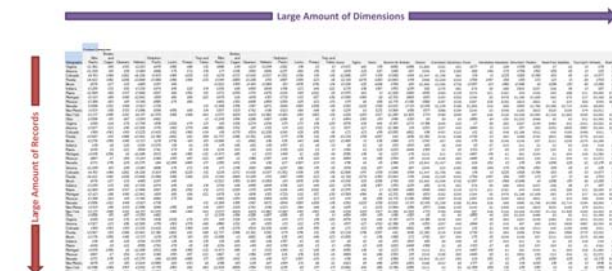
- * Two major issues
 - (1) the algorithmic approach to clustering and
 - (2) the definition and assessment of similarity versus dissimilarity.

- K clusters
- N data points
- D dimensions (original space)
- d dimensions (latent subspace)
- SC = clustering data whilst reducing the d of each cluster to a cluster-dependent subspace

Agrawal, R., Gehrke, J., Gunopulos, D. & Raghavan, P. 1998. Automatic subspace clustering of high dimensional data for data mining applications. SIGMOD Rec., 27, (2), 94-105, doi:10.1145/276305.276314.



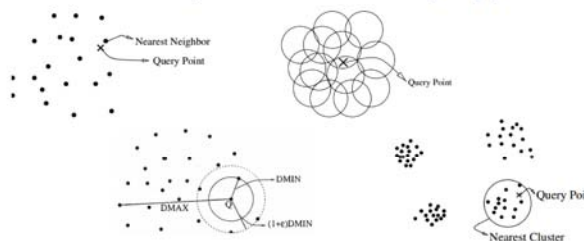
Keim, D., Kohlhammer, J., Ellis, G. & Mansmann, F. (eds.) 2010. Mastering the Information Age: Solving Problems with Visual Analytics, Goslar: Eurographics.
<http://www.vismaster.eu/wp-content/uploads/2010/11/VisMaster-book-lowres.pdf>



- Irrelevant Dimensions
- Correlated and Redundant Dimensions
- Conflicting Dimensions
- Challenging Interpretation of data and analysis results

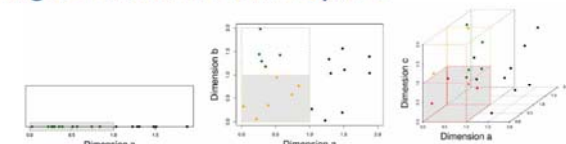
Beyer, K., Goldstein, J., Ramakrishnan, R. & Shaft, U. 1999. When is "nearest neighbor" meaningful? In: Beeri, C. & Buneman, P. (eds.) Database Theory ICDT 99, LNCS 1540. Berlin: Springer, pp. 217-235.

- NN problem: Given n data points and a query point in an m –dimensional metric space
- find the data point closest to the query point.

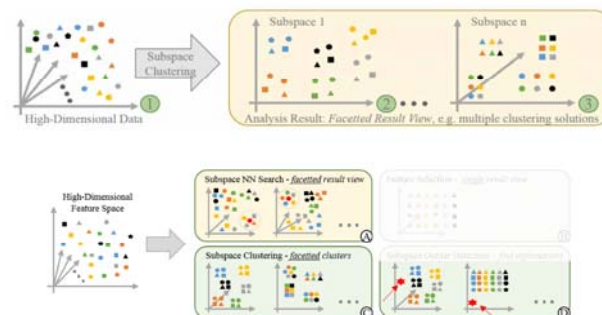
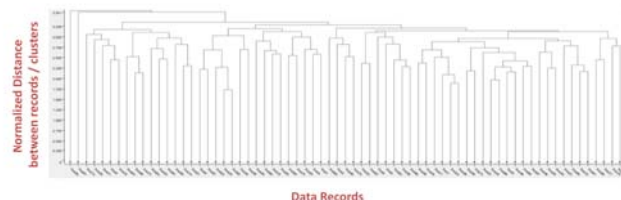


Beyer, K., Goldstein, J., Ramakrishnan, R. & Shaft, U. 1999. When is "nearest neighbor" meaningful? In: Beeri, C. & Buneman, P. (eds.) Database Theory ICDT 99, LNCS 1540. Berlin: Springer, pp. 217-235.

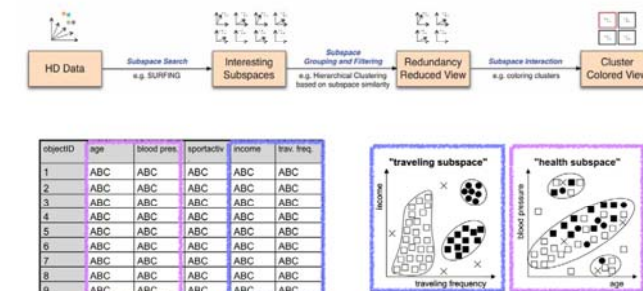
- Concentration Effect
 - Discriminability of similarity gets lost
 - Impact on usefulness of a similarity measure
- High-Dimensional Data is Sparse



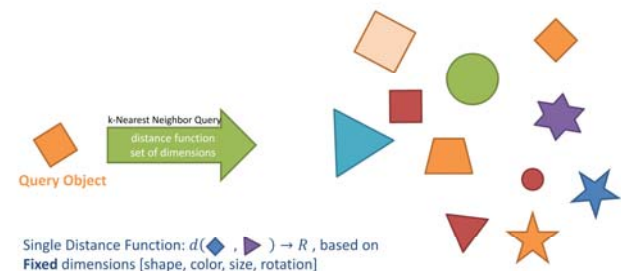
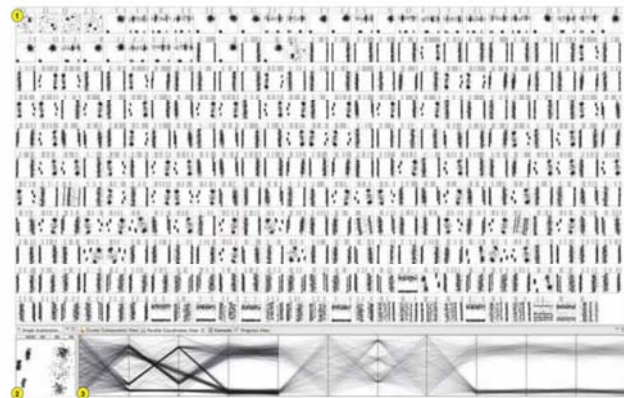
Optimization Problem and Combinatorial Issues
 Feature selection and dimension reduction
 $2^d - 1$ possible subsets of dimensions (\rightarrow subspaces)



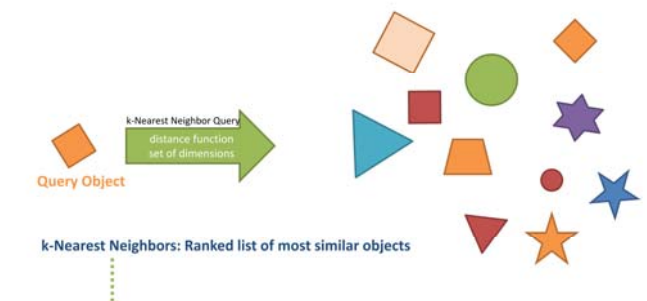
- Patterns may be found in subspaces (dimension combinations)
- Patterns may be complementary or redundant to each other



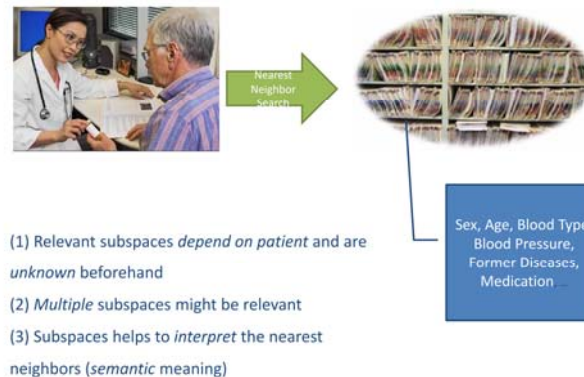
Tattu, A., Maass, F., Faerber, I., Bertini, E., Schreck, T., Seidl, T. & Keim, D. Subspace search and visualization to make sense of alternative clusterings in high-dimensional data. IEEE Symposium on Visual Analytics Science and Technology (VAST), 2012 Seattle. IEEE, 63-72, doi:10.1109/VAST.2012.6400488.



Hund, M., Behrisch, M., Färber, I., Sedlmair, M., Schreck, T., Seidl, T. & Keim, D. 2015. Subspace Nearest Neighbor Search-Problem Statement, Approaches, and Discussion. *Similarity Search and Applications*. Springer, pp. 307-313.

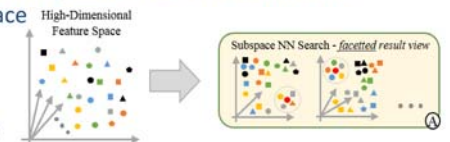


- Attention: Similarity measures lose their discriminative ability
- Noise, irrelevant, redundant, and conflicting dimensions appear



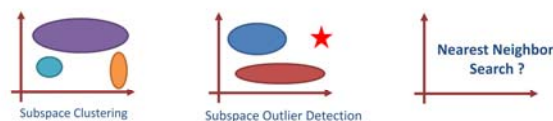
- (1) Relevant subspaces *depend on patient* and are *unknown beforehand*
- (2) *Multiple* subspaces might be relevant
- (3) Subspaces helps to *interpret* the nearest neighbors (*semantic meaning*)

1. Detect all previously unknown subspaces that are relevant for a NN-search
2. Determine the respective set of NN within each relevant subspace



Characteristics:

- Search for different NN's in different subspaces
- Consider local similarity (instead of global)
- Subspaces are query dependent
- Subspaces are not an abstract concept but helps to semantically interpret the nearest neighbors



Subspace clustering aims at finding clusters in different axis-parallel or arbitrarily-oriented subspaces [1]

Subspace Outlier Detection search for subspaces in which an arbitrary, or a user-defined object is considered as outlier [2].

[1] Kriegel, H. P., Kroger, P. & Zimek, A. 2009. Clustering High-Dimensional Data: A Survey on Subspace Clustering, Pattern-Based Clustering, and Correlation Clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 3, (1), 1-58, doi:10.1145/1497577.1497578.

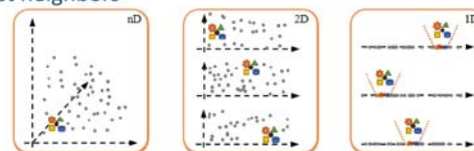
[2] Zimek, A., Schubert, E. & Kriegel, H. P. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5, (5), 363-387.

Relevance of Nearest Neighbors

A set of objects a, b, c are NN of the query q in a subspace s , iff a, b , and c are similar to q in all dimensions of s .

Relevance of a Subspace

A subspace is considered **relevant**, iff it contains relevant nearest neighbors



Dimensionality

Hund, M., Behrisch, M., Färber, I., Sedlmair, M., Schreck, T., Seidl, T. & Keim, D. 2015. Subspace Nearest Neighbor Search-Problem Statement, Approaches, and Discussion. *Similarity Search and Applications*. Springer, pp. 307-313.

- Interpretability:** reflects the semantic meaning
 - In which way are NN's similar to the query?
 - In all dimensions of the subspace
- Fulfills the downward-closure property**
 - Make use of *Apriori-like algorithms* for subspace search
- No global distance function necessary**
 - Heterogeneous subspaces can be described
 - Compute the nearest neighbors in every dimension separately (with an appropriate distance function)
 - Compute subspace by intersection



Supplementary Material

- <http://files.dbvis.de/sisap2015>

Dataset

- USDA National Nutrition Database
- <http://ndb.nal.usda.gov/>

Experiment

- Full Space (Eucl. distance, 50 dim.)
- Subspaces (our model)

Full Space	Subspace 1	Subspace 2
butter, whipped	butter, whipped	butter, whipped
butter, without salt	butter, oil, salt/sodium	butter, without salt
butter, oil, salt/sodium	butter, oil, salt/sodium	salt, strong, mayo
ketchup, full fat	lard	margarine
margarine	salt, strong, mayo	chicken, broilers
peas/beans	oil, saffron	pork, bacon
sauces	oil, stout	saunders, herbaceous
cream	oil, olive	saunders, hard
cheese, cream	oil, saffron	saunders, rich/brine
pie crust	vegetable oil, palm kernel	saunders, main backbone
cheese, ricotta/milk	oil, clove	cheating gun
ketchup, tomato	oil, saffron	pudding, vanilla
snop	margarine	glides
cheese, burdure	chickening	swastikas, tobacco
peppers	chicken, broilers	hemp, cone
strong tobacco	oil, cone, peanut, and olive	hemp, rope

TU Discussion and Open Research Questions



(1) Determine Nearest Neighbors per Dimension



(2) Efficient Search Strategy



(4) Subspace Quality Criterion (Depends on Analysis Task)

(5) Evaluation Methods and Development of Benchmark Datasets



(6) Multi-input Subspace Nearest Neighbor Search

(7) Visualization and User Interaction

Example Clust Nails Tatu et al (2012)

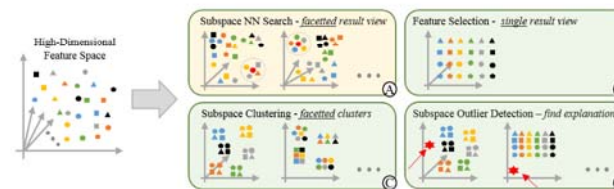
- Which dimensions occur more often in clusters?
- Which occur often together?
- Which values do records in a specific cluster have?



Tatu, A., Albuquerque, G., Eisemann, M., Schneidewind, J., Theisel, H., Magnor, M. & Keim, D. Combining automated analysis and visualization techniques for effective exploration of high-dimensional data. *Visual Analytics Science and Technology*, 2009. VAST 2009. IEEE Symposium on, 2009. IEEE, 59-66.

Tatu, A., Maass, F., Faerber, I., Bertini, E., Schreck, T., Seidl, T. & Keim, D. Subspace search and visualization to make sense of alternative clusterings in high-dimensional data. IEEE Symposium on Visual Analytics Science and Technology (VAST), 2012 Seattle. IEEE, 63-72, doi:10.1109/VAST.2012.6400488.

TU Subspace Clustering HCI-KDD



Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majnari, L. & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Results Based on Subspace Clustering. In: Guo, Y., Friston, K., Aldo, F., Hill, S. & Peng, H. (eds.) *Brain Informatics and Health, Lecture Notes in Artificial Intelligence* LNAI 9250. Cham: Springer International Publishing, pp. 358-368, doi:10.1007/978-3-319-23344-4_35.

TU Wien Further Subspace Cluster Visualization Techniques HCI-KDD

- VISA by Assent et al. (2007)
- CoDa by Günnemann et al (2010)
- Morpheus by Müller et al. (2008)
- Visual Analytics Framework by Tatu et al. (2012), see before



TU Witten Projected Clustering / Subspace Clustering / Alternative

- Variety of different algorithms, e.g. PROCLUS [1], CLIQUE [2], RESCUE [3]
- Example CLIQUE:



- Challenges
- Exponential # of possible subspaces
- Result highly depend on parameters
- Highly redundant results (clusters + subspaces)

Visual Analytics for Subspace Steering HCL-KDD

- Existing techniques: **exploration** of subspace clusters
- Visualizations to **make sense** of clusters and its subspaces

Is the parameter setting appropriate for the data?

What happens if algorithms cannot scale with the #dimensions?

- We need methods to **steer algorithms** while computing relevant subspaces
 - Pruning of intermediate results
 - Adjust parameters to domain knowledge
 - ...

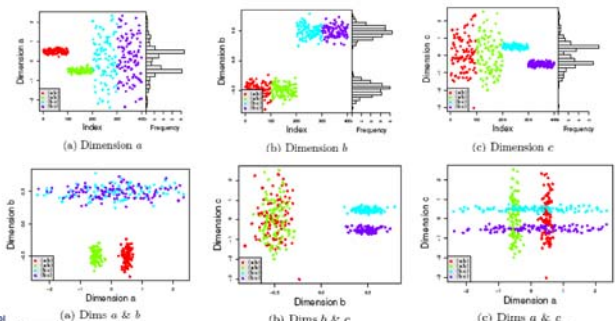
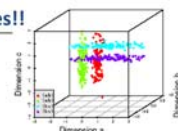


Fig. 3 A screenshot of our visual analysis tool SubViz. It enables the user to interactively explore a large number of subspaces. A general overview of the similarity between the subspaces is given by an MDS plot (A). Small multiples (B) allows to preview projections of different distance functions and a quick change of the MDS plot. On the very top (C) the user is provided with some distribution properties of the subspaces such as the #dimensions. A heatmap (D) provides more details of relationships between the pair-wise distances. An aggregation table (E) shows the values of the aggregated cluster members and the table lens (F) provides details on demand.

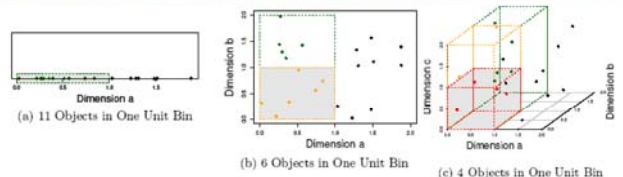
Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majnarić, L. & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Results Based on Subspace Clustering. In: Guo, Y., Friston, K., Aldo, F., Hill, S. & Peng, H. (eds.) Brain Informatics and Health, Lecture Notes in Artificial Intelligence LNAI 9250. Cham: Springer International Publishing, pp. 358-368, doi:10.1007/978-3-319-23344-4_35.

Interesting Clusters may ONLY exist in subspaces!!

Parsons, L., Haque, E. & Liu, H. 2004. Subspace clustering for high dimensional data: a review. SIGKDD Explorations 6, (1), 90-105.

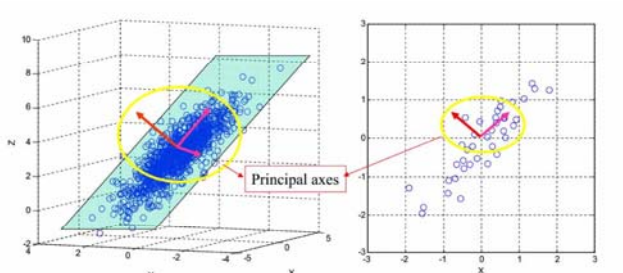


06 “What is interesting?” Projection Pursuit



Data in only one dimension is relatively packed
Adding a dimension “stretch” the points across that dimension, making them further apart
Adding more dimensions will make the points further apart—high dimensional data is extremely sparse
Distance measure becomes meaningless—due to equidistance

Similar concept : Principal Component Analysis (PCA)



Huber (1985): “What is interesting?”

- **Projection pursuit** : Find a subset of coordinates of the data which display “interesting” features. Often the selection of the subset of coordinates is manual, but there are automated algorithms which can find these subsets automatically also. Finally one has to inspect each projection and decide if its “interesting”.

Huber P.J.: Projection pursuit. *Ann. Statist.* 13, 2 (1985), 435-525.

- **Dataset** - consists of a matrix of data values, rows represent individual instances and columns represent dimensions.
- **Instance** - refers to a vector of d measurements.
- **Cluster** - group of instances in a dataset that are more similar to each other than to other instances. Often, similarity is measured using a distance metric over some or all of the dimensions in the dataset.
- **Subspace** - is a subset of the d dimensions of a given dataset.
- **Subspace Clustering** – seek to find clusters in a dataset by selecting the most *relevant* dimensions for each cluster separately.
- **Feature Selection** - process of determining and selecting the dimensions (features) that are most relevant to the data mining task.

Black-Box approach

Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples

Nicolas Papernot and Patrick McDaniel
The Pennsylvania State University
University Park, PA
{ngp5056, mcdaniel}@cse.psu.edu

Ian Goodfellow
OpenAI
San Francisco, CA
ian@openai.com

ABSTRACT

Many machine learning models are vulnerable to adversarial examples: inputs that are specially crafted to cause a machine learning model to produce an incorrect output. Adversarial examples that affect one model often affect another model, even if the two models have different architectures or were trained on different training sets, so long as both models were trained to perform the same task. An attacker may therefore train their own substitute model, craft adversarial examples against the substitute, and transfer them to a victim model, with very little information about the victim. Recent work has further developed a technique that uses the victim model as an oracle to label a synthetic training set for the substitute, so the attacker need not even collect a training set to mount the attack. We extend these recent techniques using reservoir sampling to greatly enhance the efficiency of the training procedure for the substitute model. We introduce new transferability attacks between previously unexplored (substitute, victim) pairs of machine learning model classes, most notably SVMs and decision trees. We demonstrate our attacks on two commercial machine learning classification systems from Amazon (96.19% misclassification rate) and Google (88.94%) using only 800 queries of the victim model, thereby showing that existing machine

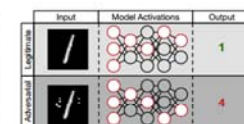


Figure 1: An adversarial sample (bottom row) is produced by slightly altering a legitimate sample (top row) in a way that forces the model to make a wrong prediction whereas a human would still correctly classify the sample [15].

Adversarial sample transferability is the property that some adversarial samples produced to mislead a specific model f can mislead other models f' even if their architecture greatly differ [23, 15, 24]. A practical impact of this property is that it leads to oracle-based black box attacks. In

Least Gaussian projections of the data (interesting?)

how to define non-Gaussianity?

covariance and mean given: Gaussian distribution maximizes the entropy

Objective: minimize $H(t)$ for $t = w^T x$
 t is normalized to zero mean and unit variance

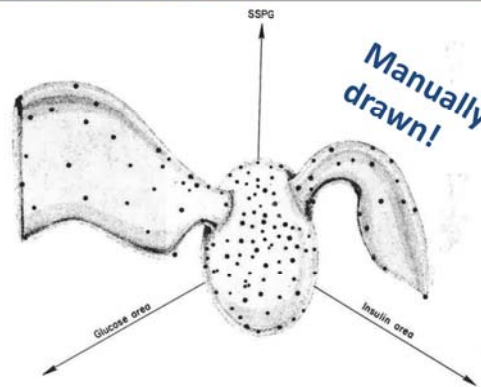
This is difficult to optimize

- finding unimodal super-Gaussians
- finding multimodal distributions

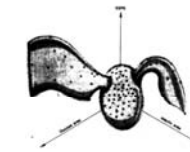
Other criteria are given for ICA: kurtosis and different contrast functions which measure non-Gaussianity

- 145 diabetes patients
- 6 dimensional data set:
 - 1) age,
 - 2) relative weight,
 - 3) fasting plasma glucose,
 - 4) area under the plasma glucose curve for the three hour glucose tolerance test (OGTT),
 - 5) area under the plasma insulin curve for the OGTT,
 - 6) steady state plasma glucose response.
- Method: Projection Pursuit (PP)
- Result: $\mathbb{R}^6 \rightarrow \mathbb{R}^3$

Reaven, G. & Miller, R. (1979) An attempt to define the nature of chemical diabetes using a multidimensional analysis. *Diabetologia*, 16, 1, 17-24.



Reaven, G. & Miller, R. (1979) An attempt to define the nature of chemical diabetes using a multidimensional analysis. *Diabetologia*, 16, 1, 17-24.



Given a point cloud data set X and a covering U
 \Rightarrow simplicial complex

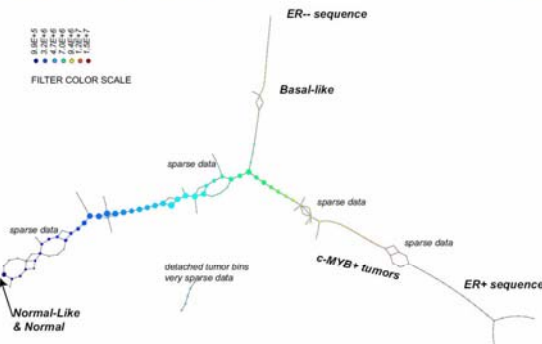
$$f: X \rightarrow \mathbb{R}$$

$$f: X \rightarrow \mathbb{Z}$$

$$U = \{U_\alpha\}_{\alpha \in A}$$

$$f_\varepsilon(x) = C_\varepsilon \sum_y \exp\left(\frac{-d(x, y)^2}{\varepsilon}\right)$$

Singh, G., Mémoli, F. & Carlsson, G. (2007). Topological methods for the analysis of high dimensional data sets and 3D object recognition. *Eurographics Symposium on Point-Based Graphics*, Euro Graphics Society, 91-100.



Nicolau, M., Levine, A. J. & Carlsson, G. (2011) Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival. *Proceedings of the National Academy of Sciences*, 108, 17, 7265-7270.

- Time (e.g. entropy) and Space (e.g. topology)
- Knowledge Discovery from “unstructured” ;-) (Forrester: >80%) data and applications of structured components as methods to index and organize data -> Content Analytics
- Open data, Big data, sometimes: small data
- Integration in “real-world” (e.g. Hospital), mobile
- How can we measure the benefits of visual analysis as compared to traditional methods?
- Can (and how can) we develop powerful visual analytics tools for the non-expert end user?

- Why would we wish at all to reduce the dimensionality of a data set?
- Why is feature selection so important? What is the difference between feature selection and feature extraction?
- What types of feature selection do you know?
- Can Neural Networks also be used to select features?
- Why do we need a human expert in the loop in subspace clustering?
- What is the advantage of the Projection Pursuit method?
- Why is algorithm selection so critical?

- What are the problems in high-dimensional spaces?
- When is the human-in-the-loop beneficial?
- What is a Autoencoder and when would you use it?
- When would you use PCA?
- What did the authors of the Miller-Reaves study do?
- Why is the question “what is interesting?” a hard question?



Thank you!