



Andreas Holzinger
VO 709.049 Medical Informatics
07.12.2016 11:15-12:45



Lecture 07 Dimensionality Reduction and Subspace Clustering with the Doctor-in-the-Loop

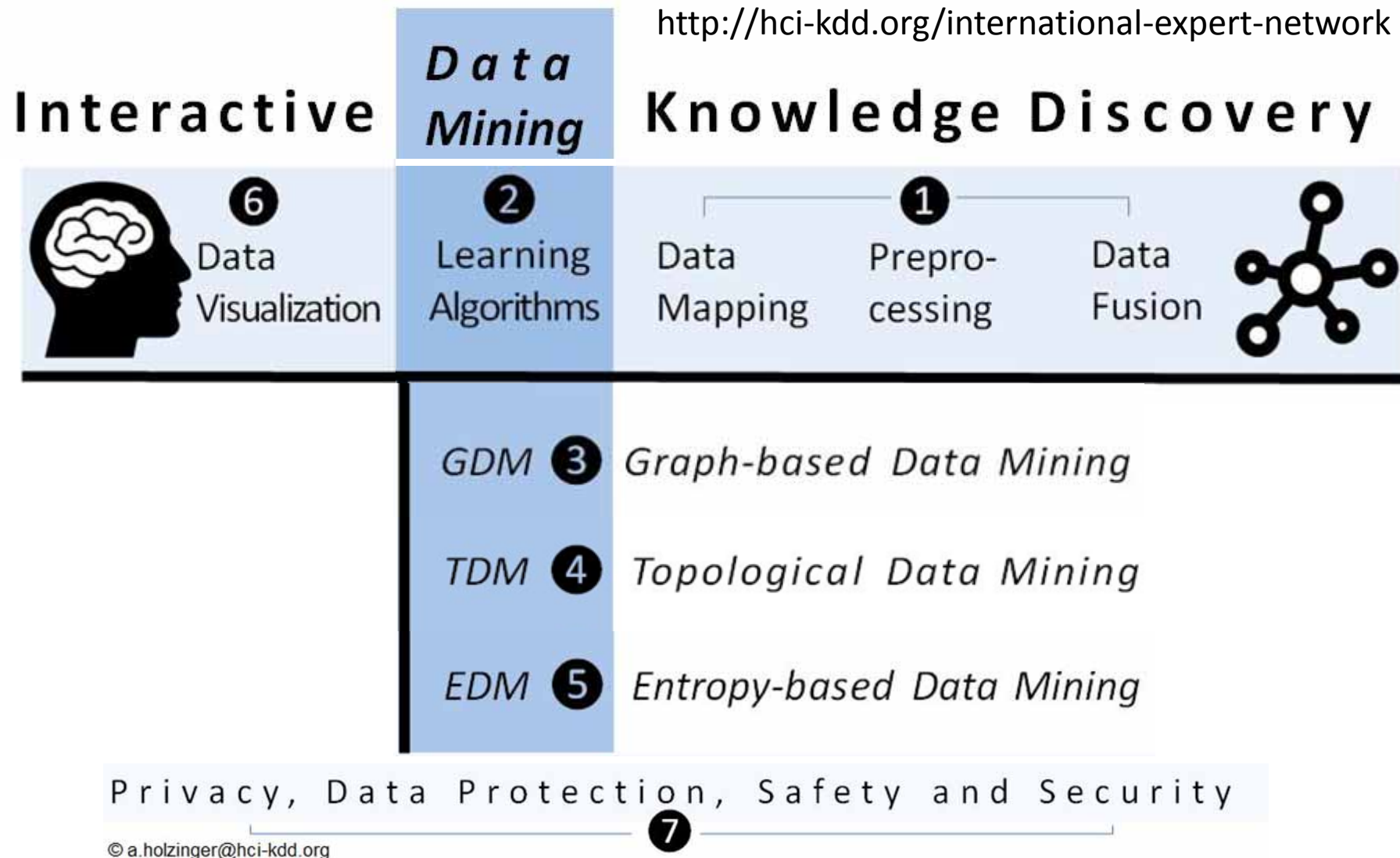
a.holzinger@tugraz.at

Tutor: markus.plass@student.tugraz.at

<http://hci-kdd.org/biomedical-informatics-big-data>



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



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

- Classification
- Clustering
- Curse of dimensionality
- Dimensionality reduction
- Interestingness
- Feature extraction
- Feature selection
- Mapping
- Subspace analysis
- Subspace clustering

- **Artificial neural network (ANN)** = a computational adaptive model (inspired by biological neural networks), consisting of interconnected groups of artificial neurons; processes information using a connectionist approach.
- **Association rule learning** = a set of techniques for discovering interesting relationships, i.e., “association rules,” among variables in large databases used for data mining;
- **Classification** = a set of techniques to identify the categories in which new data points belong, based on a training set containing data points that have already been categorized; these techniques are often described as supervised learning because of the existence of a training set; they stand in contrast to cluster analysis, a type of unsupervised learning; used e.g. for data mining;
- **Cluster analysis** = statistical method for classifying objects that splits a diverse group into smaller groups of similar objects, whose characteristics of similarity are not known in advance; a type of unsupervised learning because training data are not used - in contrast to classification; used for data mining.
- **Data mining** = a set of techniques to extract patterns from large data by combining methods from statistics and machine learning with database management (e.g. association rule learning, cluster analysis, classification, regression, etc.);
- **Knowledge Discovery (KD)** = process of identifying valid, novel, useful and understandable patterns out of large volumes of data

- **Deep Learning** = class of machine learning algorithms using layers of non-linear processing units for feature extraction (remember: features are key for learning and understanding) - learning representations from data;
- **Knowledge Extraction** = is the creation of knowledge from structured (relational databases, XML) and unstructured (text, documents, images) sources;
- **Multimedia** = several data of different modalities are processed at the same time, i.e. encompassing audio data (sound, speech), image data (b/w and colour images), video data (time-aligned sequences of images), electronic ink (sequences of time aligned 2D and 3D coordinates of a stylus, pen, data gloves etc.)
- **Principal Component Analysis (PCA)** = statistical technique for finding patterns in high-dimensional data;
- Sparse Data =
- **Supervised learning** = inferring a function from supervised training data on the basis of training data which consist of a set of training examples, the input objects (typically vectors) and a desired output value (also called the supervisory signal).
- **Supervised learning algorithm** = analyzes the training data and produces an inferred function, called a classifier (if the output is discrete) or a regression function (if the output is continuous); the algorithm generalizes from the training data to unseen situations.
- **Support vector machine (SVM)** = concept for a set of related supervised learning methods to analyze data and recognize patterns, used for classification and regression analysis.
- **Unsupervised learning** = establishes clusters in data, where the class labels of training data is unknown.

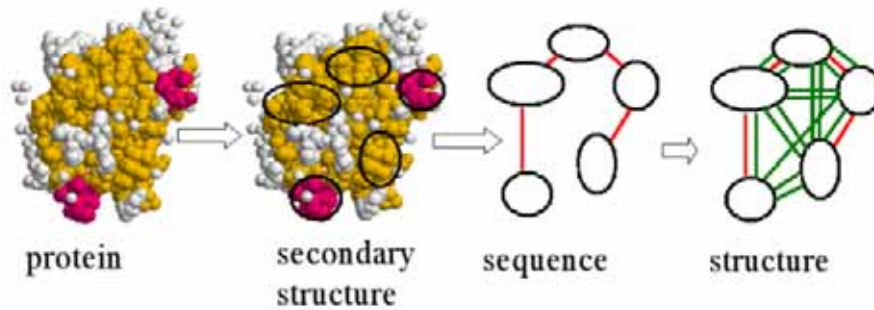
- ANN = Artificial Neural Network
- ANOVA = Analysis of Variance
- AUC - area under the curve
- CDT = Clinical Decision Tree
- DM = Data Mining
- KDD = Knowledge Discovery from Data(bases)
- LLE = Locally Linear Embedding
- MDS = Multi Dimensional Scaling
- MELD = model for end-stage liver disease
- MM = Multimedia
- NLP = Natural Language Processing
- PCA = Principal Components Analysis
- ROC = Receiver Operating Characteristic
- SVM = Support Vector Machine

- ... know the differences between **classification** and **clustering** and why it is important for health;
- ... are aware that **features** are key to learning and understanding;
- ... understand the **curse of dimensionality**;
- ... have an idea of **dimensionality reduction**;
- ... recognize the value of **subspace clustering** and analysis with the **doctor-in-the-loop**;
- Understand why the question “**What is interesting?**” is not easy to answer;

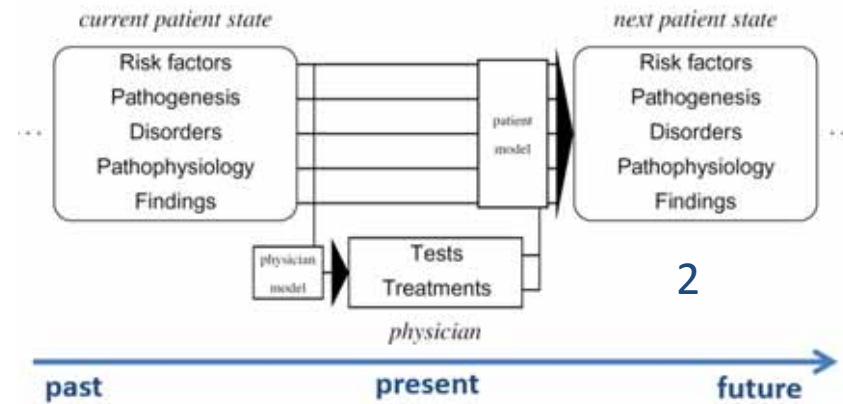
- Uncertainty
- Validation
- Curse of Dimensionality
- Large spaces gets sparse
- Distance Measures get useless
- Patterns occur in different subspaces
- Pressing question: “What is interesting?”

- **00 Reflection – follow-up from last lecture**
- **01 Classification vs. Clustering**
- **02 Feature Engineering**
- **03 Curse of Dimensionality**
- **04 Dimensionality Reduction**
- **05 Subspace Clustering and Analysis**
- **06 Projection Pursuit**
- **07 Conclusion and Future Challenges**

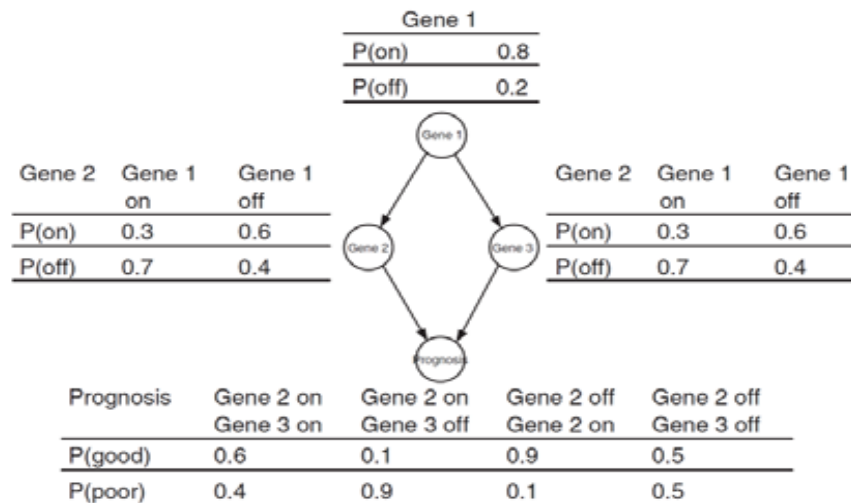
00 Reflection



1



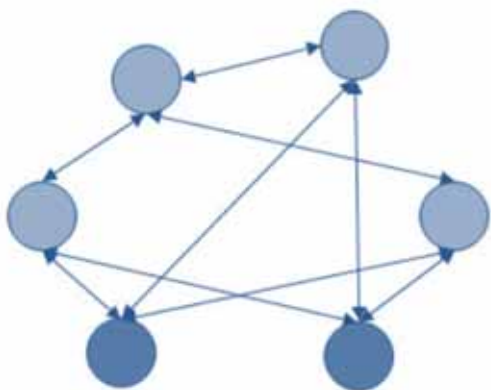
2



3

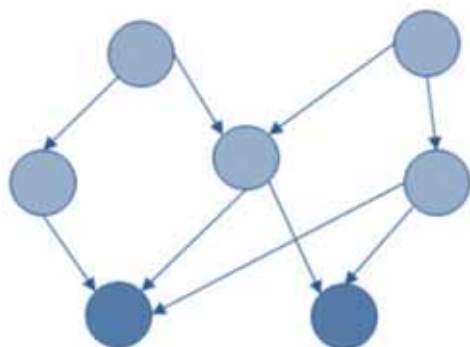
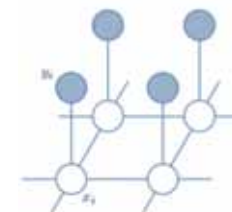


4



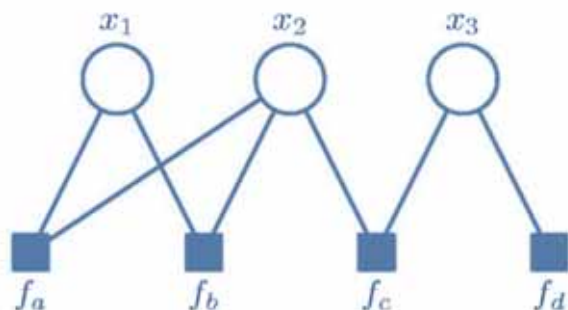
Undirected: Markov random fields, useful e.g. for computer vision (Details: Murphy 19)

$$P(\mathbf{X}) = \frac{1}{Z} \exp \left(\sum_{ij} W_{ij} x_i x_j + \sum_i x_i b_i \right)$$



Directed: Bayes Nets, useful for designing models (Details: Murphy 10)

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \text{pa}_k)$$

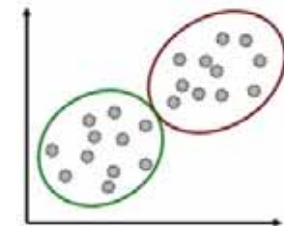
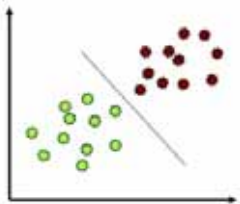


Factored: useful for inference/learning

$$p(\mathbf{x}) = \prod_s f_s(\mathbf{x}_s)$$

Tutorial on Factor Graphs <http://deepdive.stanford.edu/inference>

01 Classification vs. Clustering





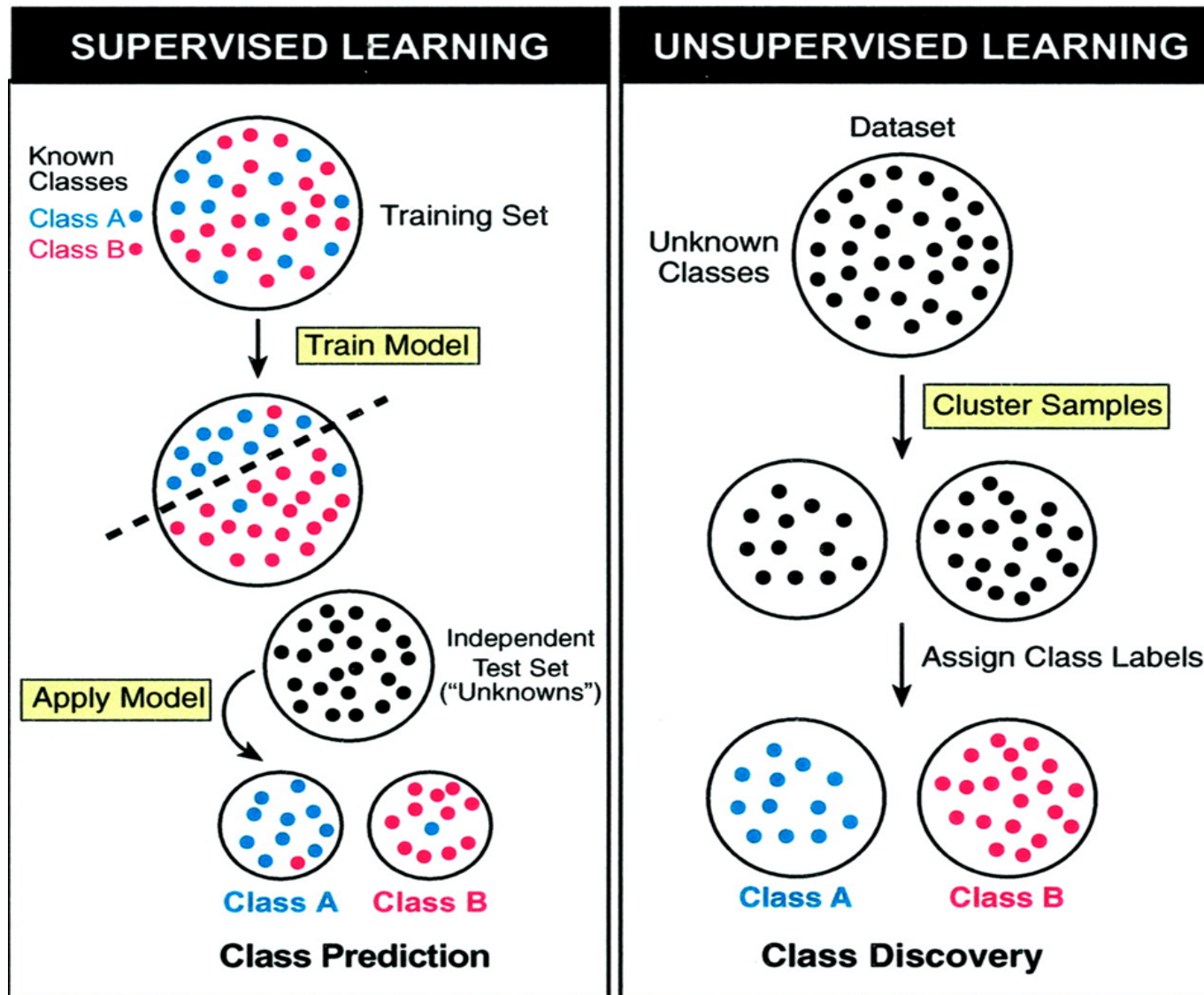
- 1) The data is not labeled (clA/Clu)?
- 2) Identify structure/patterns (clA/Clu)?
- 3) Predicting an item set, identify to which set of categories a new observation belongs (clA/Clu)?
- 4) Assigning a set of objects into groups (clA/Clu)?
- 5) Having many labelled data points (clA/Clu)
- 6) Using the concept of supervised learning (clA/Clu)?
- 7) Grouping data items close to each other (clA/Clu)?
- 8) Used to explore data sets (clA/Clu)?

■ Classification

- Supervised learning, Pattern Recognition, 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

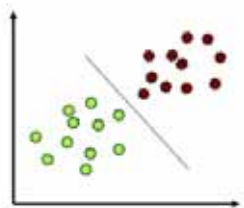
■ Clustering

- 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;
- Example: K-Means Algorithm for disease clustering



Ramaswamy, S. & Golub, T. R. (2002) DNA Microarrays in Clinical Oncology. *Journal of Clinical Oncology*, 20, 7, 1932-1941.

Why do we need Classification in Health Informatics?

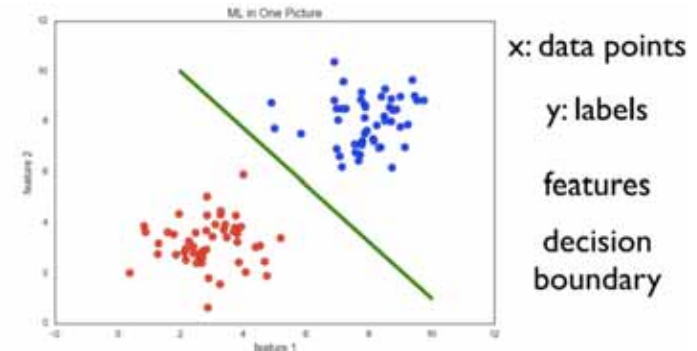




x -- set of pixel intensities

C_1 : Cancer present

C_2 : Cancer absent



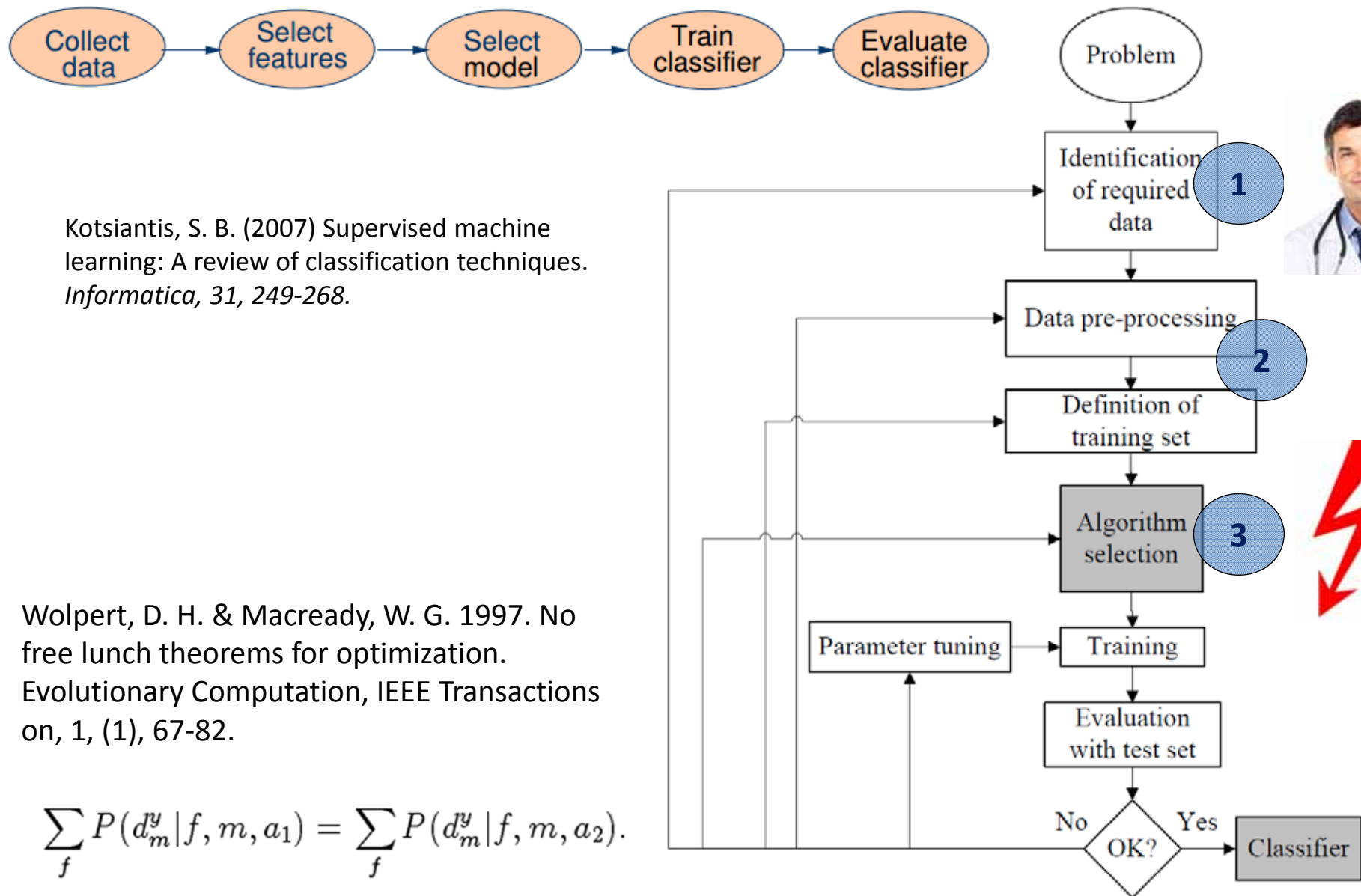
- Typical questions include:
 - Is this protein functioning as an enzyme?
 - Does this gene sequence contain a splice site?
 - Is this melanoma malign?
- Given object x – predict the class label y
 - If $y \in \{0,1\} \rightarrow$ binary classification problem
 - If $y \in \{1, \dots, n\}$ and is $n \in \mathbb{N} \rightarrow$ multiclass problem
 - If $y \in \mathbb{R} \rightarrow$ regression problem



Kotsiantis, S. B. (2007) Supervised machine learning: A review of classification techniques. *Informatica*, 31, 249-268.

Wolpert, D. H. & Macready, W. G. 1997. No free lunch theorems for optimization. *Evolutionary Computation*, IEEE Transactions on, 1, (1), 67-82.

$$\sum_f P(d_m^y | f, m, a_1) = \sum_f P(d_m^y | f, m, a_2).$$



- 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 $x_2 \Rightarrow y_1 = y_2$

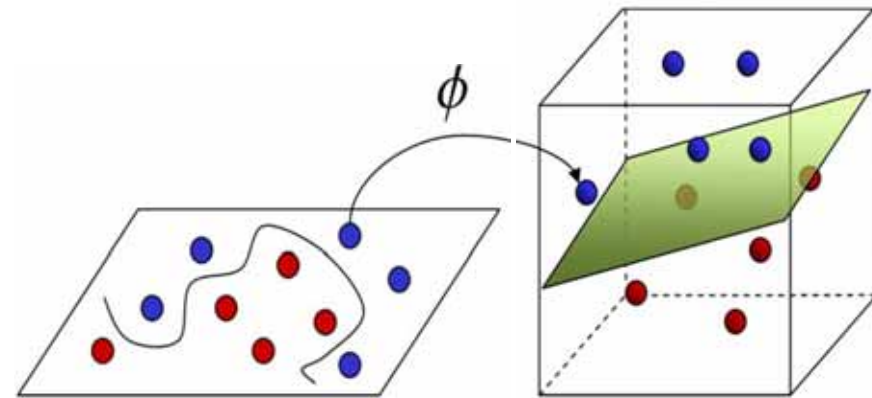
$$x_j = \operatorname{argmin}_{x \in D} ||x - x_i||^2 \Rightarrow y_i = y_j$$

- 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:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), C, \epsilon$ 
2:  $S_i \leftarrow \emptyset$  for all  $i = 1, \dots, n$ 
3: repeat
4:   for  $i = 1, \dots, n$  do
5:      $H(y) \equiv \Delta(y_i, y) + \mathbf{w}^T \Psi(\mathbf{x}_i, y) - \mathbf{w}^T \Psi(\mathbf{x}_i, y_i)$ 
6:     compute  $\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} H(y)$ 
7:     compute  $\xi_i = \max\{0, \max_{y \in S_i} H(y)\}$ 
8:     if  $H(\hat{y}) > \xi_i + \epsilon$  then
9:        $S_i \leftarrow S_i \cup \{\hat{y}\}$ 
10:     $\mathbf{w} \leftarrow \text{optimize primal over } S = \bigcup_i S_i$ 
11:    end if
12:  end for
13: until no  $S_i$  has changed during iteration
```

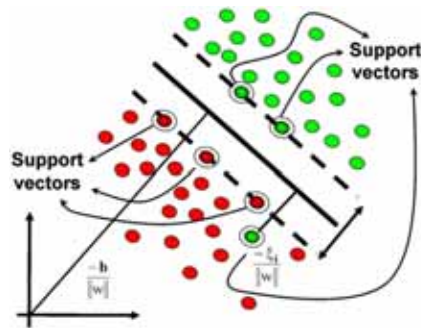
Finley, T. & Joachims, T. Supervised clustering with support vector machines. Proceedings of the 22nd international conference on Machine learning, 2005. ACM, 217-224.

- Uses a nonlinear mapping to transform the original data (input space) into a higher dimension (feature space)



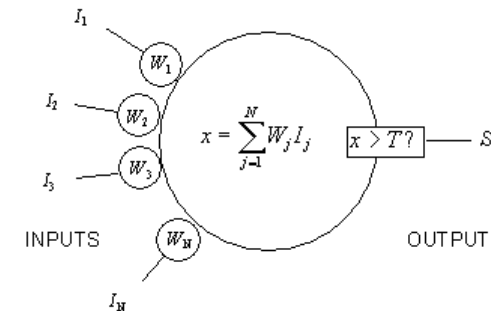
- = classification method for both linear and nonlinear data;
- Within the new dimension, it searches for the linear optimal separating **hyperplane** (i.e., “decision boundary”);
- By nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated with a hyperplane;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

■ SVM

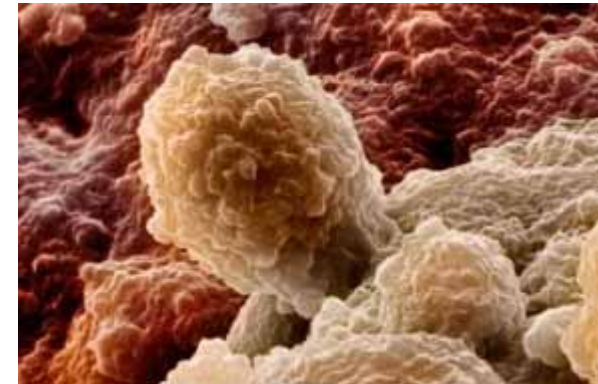
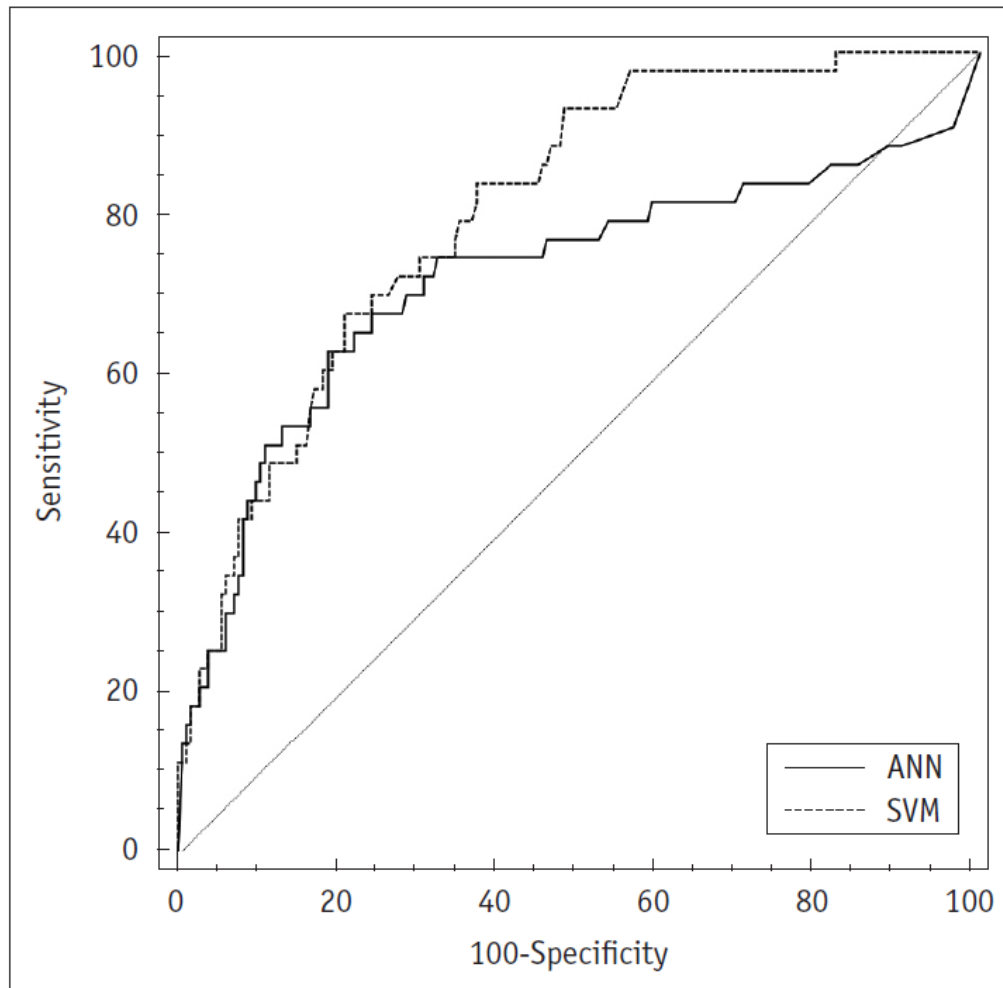


- Deterministic algorithm
- Nice generalization properties
- Hard to learn – learned in batch mode using quadratic programming techniques
- Using kernels can learn very complex functions

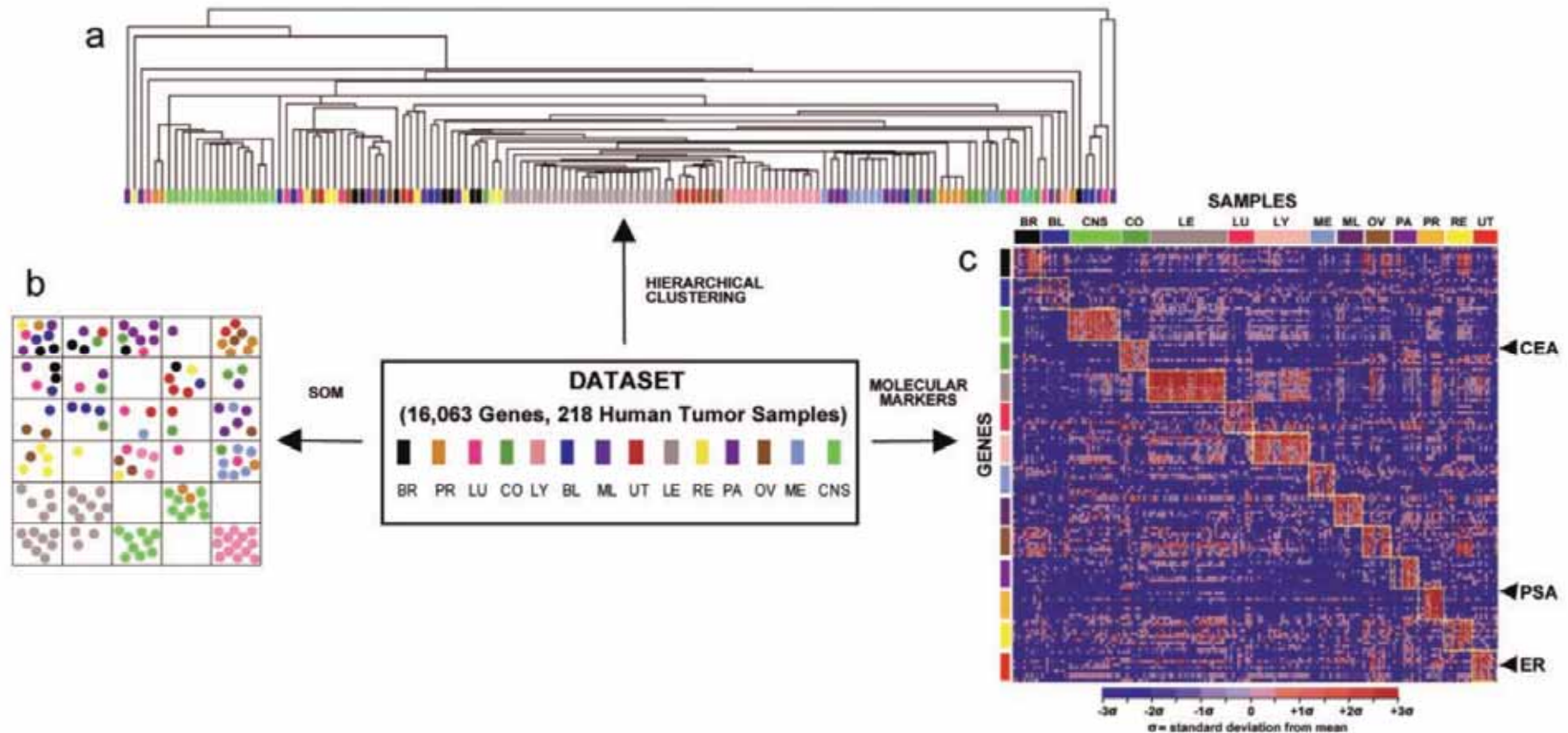
■ ANN



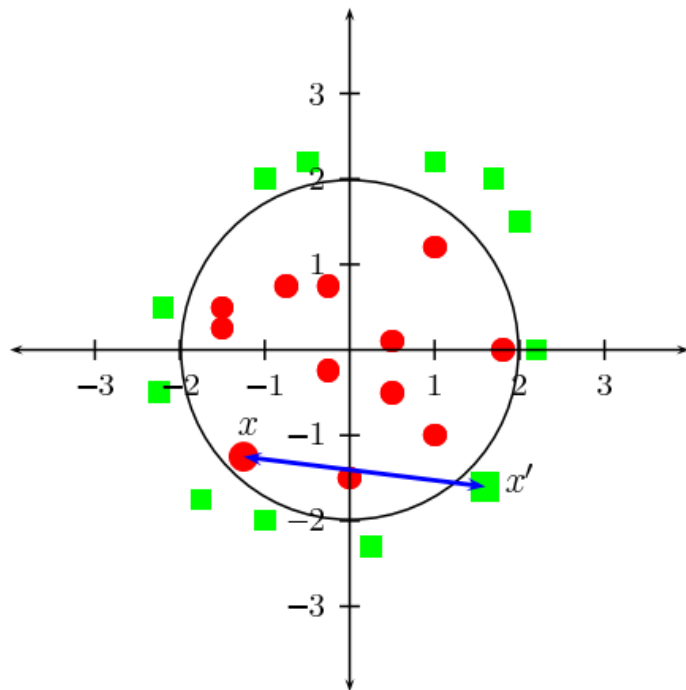
- Nondeterministic algorithm
- Generalizes well but doesn't have strong mathematical foundation
- Can easily be learned in incremental fashion
- To learn complex functions—use multilayer perceptron (nontrivial)



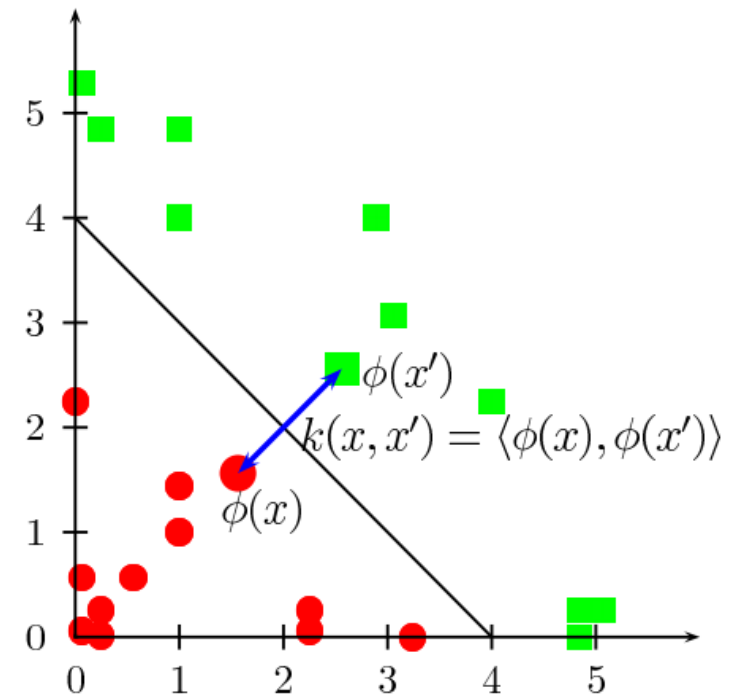
Kim, S. Y., Moon, S. K., Jung, D. C., Hwang, S. I., Sung, C. K., Cho, J. Y., Kim, S. H., Lee, J. & Lee, H. J. (2011) Pre-Operative Prediction of Advanced Prostatic Cancer Using Clinical Decision Support Systems: Accuracy Comparison between Support Vector Machine and Artificial Neural Network. *Korean J Radiol*, 12, 5, 588-594.



Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C.-H., Angelo, M., Ladd, C., Reich, M., Latulippe, E. & Mesirov, J. P. 2001. Multiclass cancer diagnosis using tumor gene expression signatures. *Proceedings of the National Academy of Sciences*, 98, (26), 15149-15154, doi:10.1073/pnas.211566398.

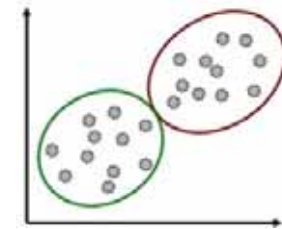


$$\mathbb{R}^2 \Rightarrow \mathcal{H}$$



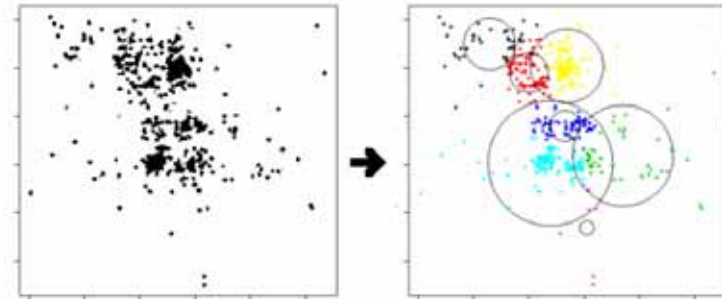
Borgwardt, K., Gretton, A., Rasch, J., Kriegel, H.-P., Schölkopf, B. & Smola, A. 2006. Integrating structured biological data by kernel max. mean discrepancy. Bioinformatics, 22, 14, e49-e57.

Why do we need Clustering in Health Informatics?



- 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

- Partite a set of n observations into k clusters so, that the intra-cluster variance is *argmin*
 - v ... variance (objective function)
 - S_i ... cluster
 - c_j ... mean (“centroid” for cluster j)
 - D ... set of all data points x_j
 - k ... number of clusters

Distance “medoid”

$$v(D) = \operatorname{argmin} \sum_{j=1}^k \sum_{i=1}^n \|(x_i - c_j)\|^2$$

Jain, A. K. 2010. Data clustering: 50 years beyond K-means. Pattern Recognition Letters, 31, (8), 651-666, doi:<http://dx.doi.org/10.1016/j.patrec.2009.09.011>.

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

1. 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 $s_i \in S$.
 2. Update each site s_i as the center of gravity of cluster C_i ; if $|C_i| = 0$, choose $s_i = x_l$ for a random $l \leq n$ with $x_l \neq s_j$ 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

Wu et al. (2008) Top 10 algorithms in data mining. *Knowledge & Information Systems*, 14, 1, 1-37.

- **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));
- **A-priori**
 - 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));*

02

Feature Engineering

“Applied Machine Learning is basically feature engineering”.



*Andrew Y. Ng, VP & Chief Scientist of Baidu;
Co-Chair/Founder of Coursera; Professor at Stanford University*

<http://www.andrewng.org>

- 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.

- Intuitively: a domain with a distance function
- Formally: Feature Space $\mathcal{F} = (\mathcal{D}, d)$
 - \mathcal{D} = ordered set of features
 - $d: D \times D \rightarrow \mathbb{R}_0^+$... 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

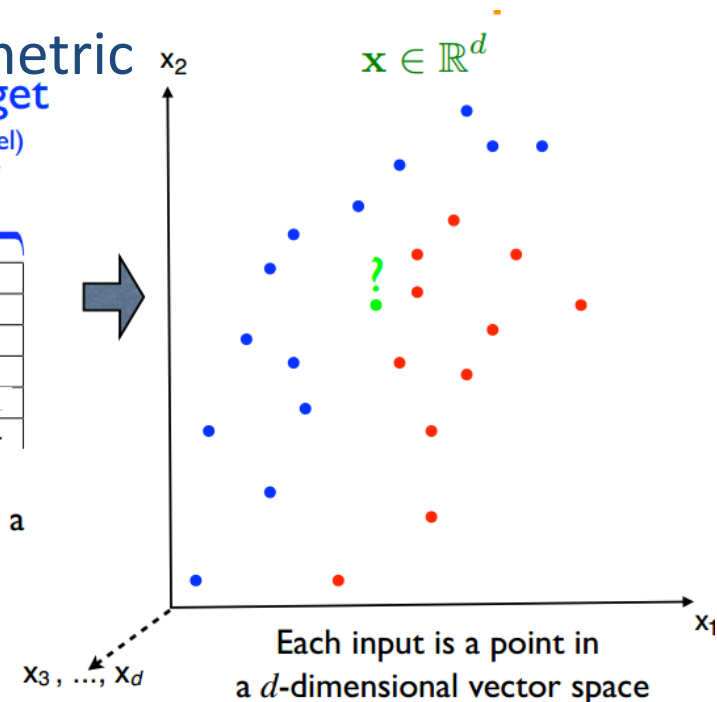
input $x \in \mathbb{R}^d$ target (label) y

n examples

x_1	x_2	x_3	x_4	x_5	
0.32	-0.27	+1	0	0.82	1
-0.12	0.42	-1	1	0.22	0
0.06	0.35	-1	1	-0.37	1
0.91	-0.72	+1	0	-0.63	1
...

Each example (row) is now a $d+1$ -dimensional vector

Image credit to Pascal Vincent



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.

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

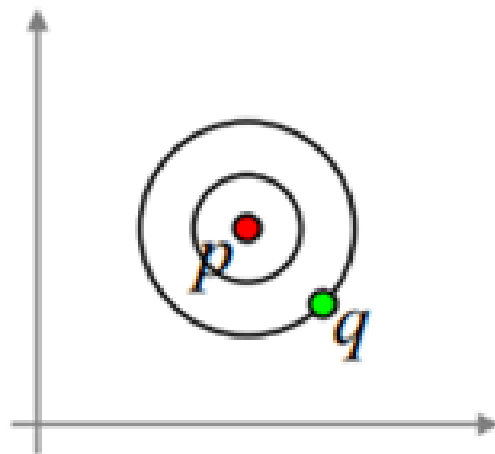
Remark 1. One example is \mathbb{R}^d with the Euclidean metric. Spheres S^n endowed with the spherical metric provide another example.

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

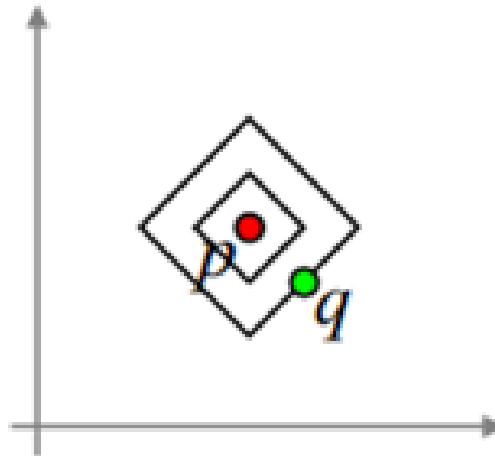
$$d(x, x) = 0$$

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

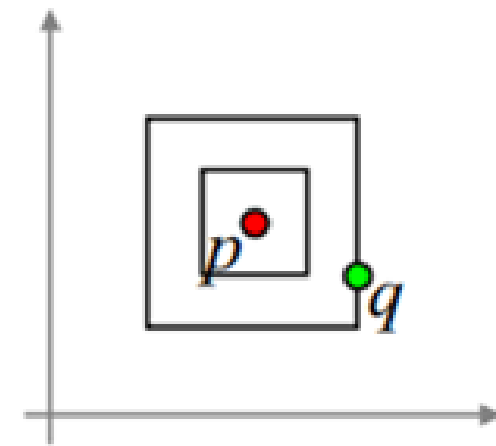
$$d(x^1, x^2) \leq d(x^1, x^3) + d(x^3, x^2) \text{ triangle inequality}$$



Euclidian norm



Manhattan norm



Maximums norm

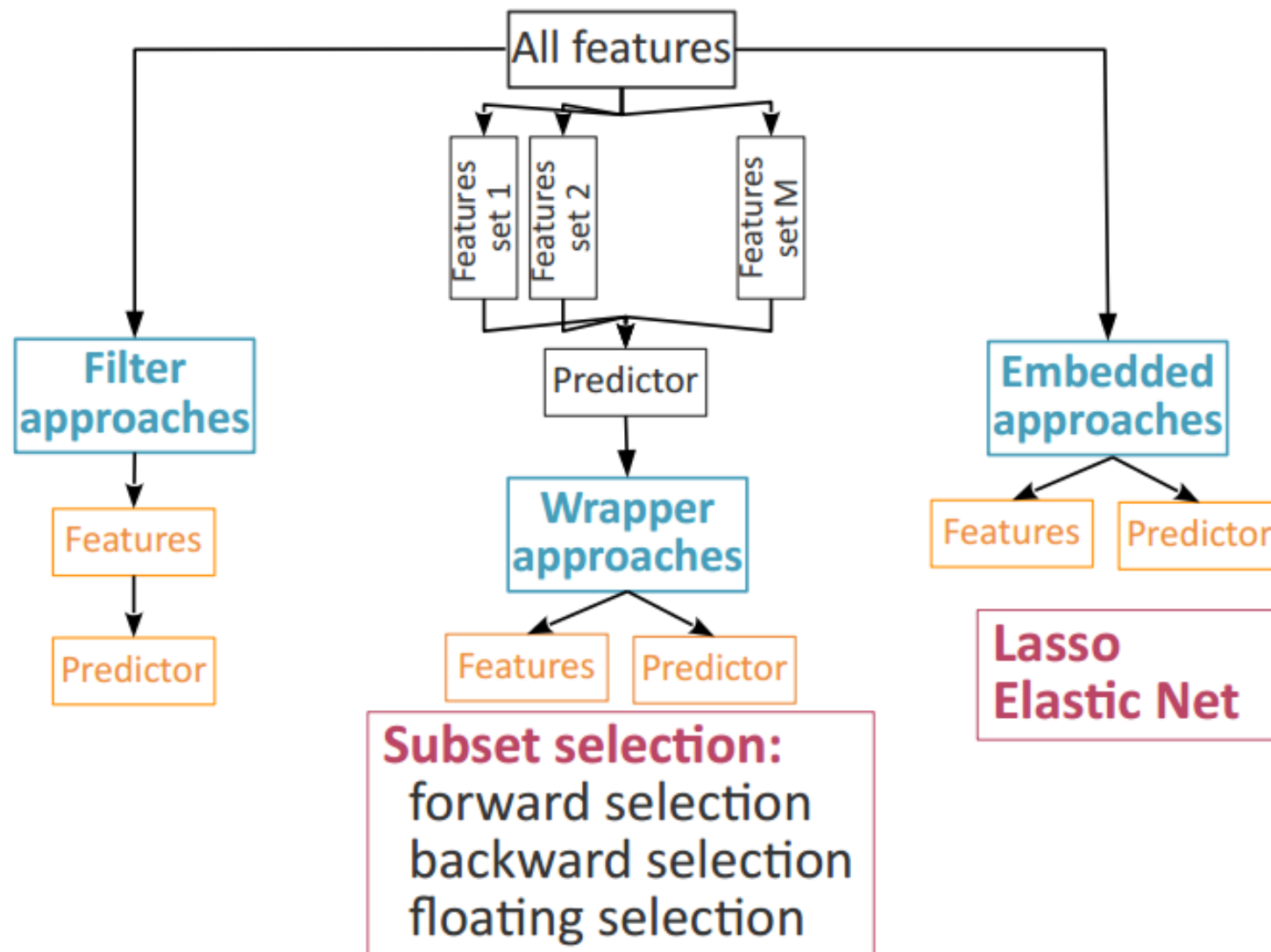
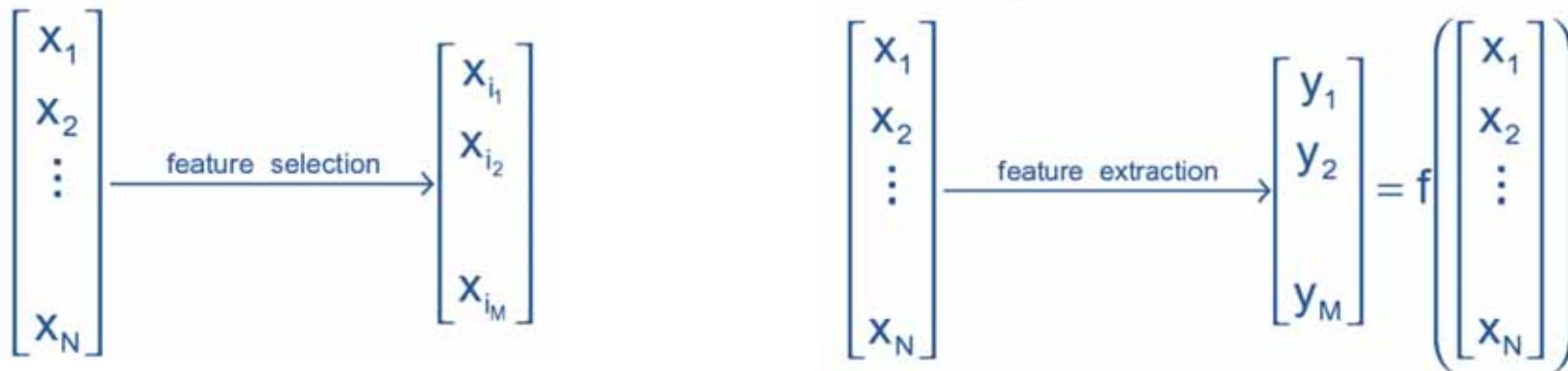


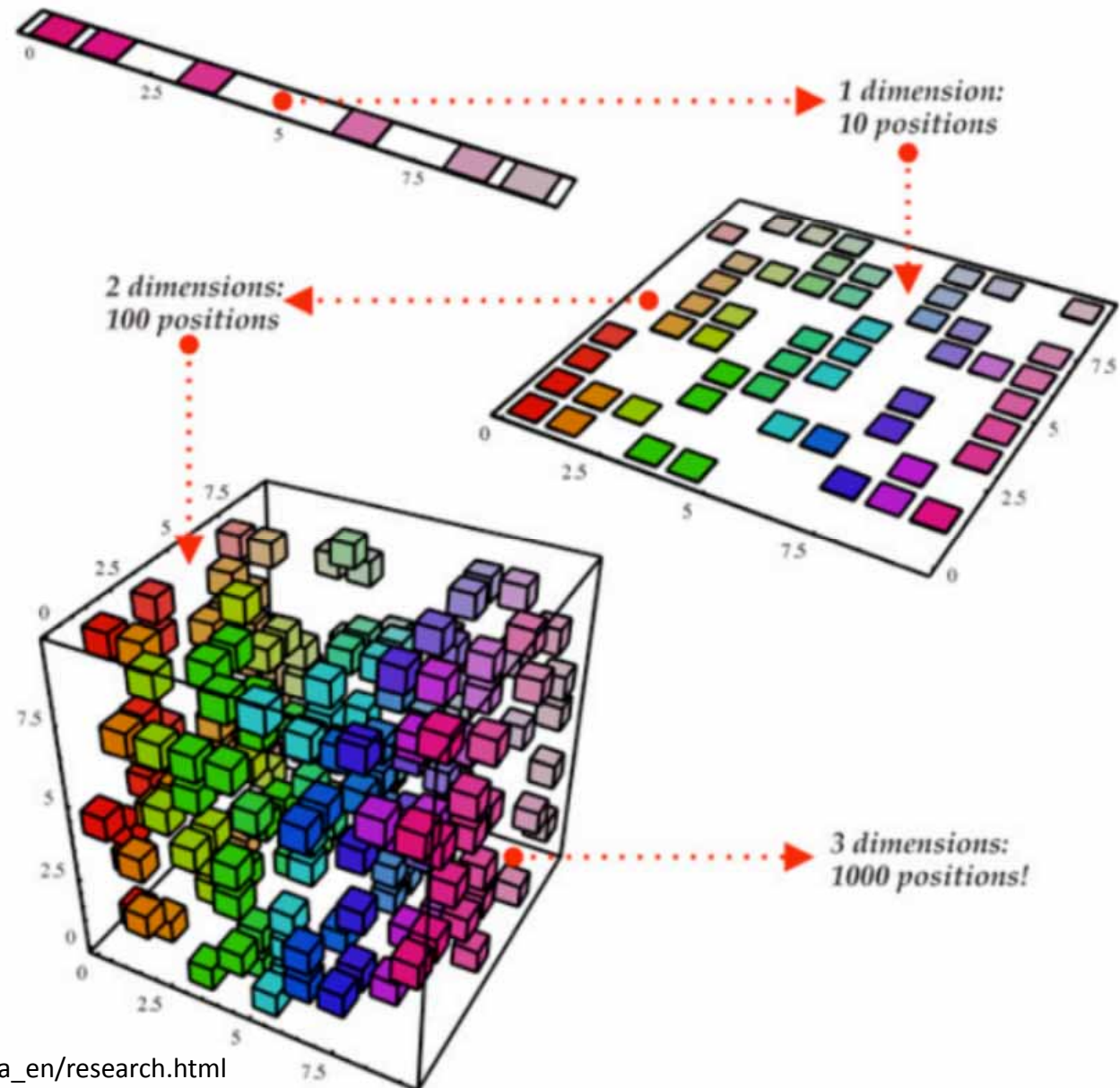
Image credit to Chloe Azencott

- 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



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

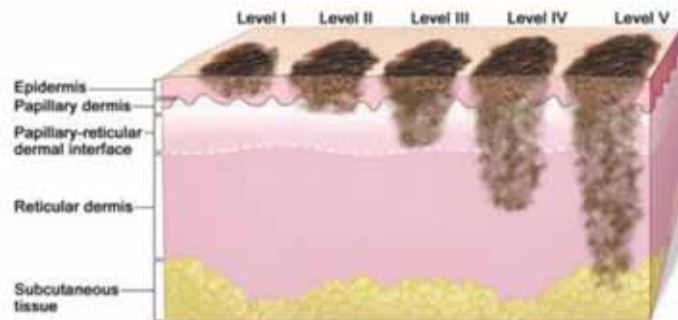
03 Curse of Dimensionality



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.iro.umontreal.ca/~bengioy/yoshua_en/research.html

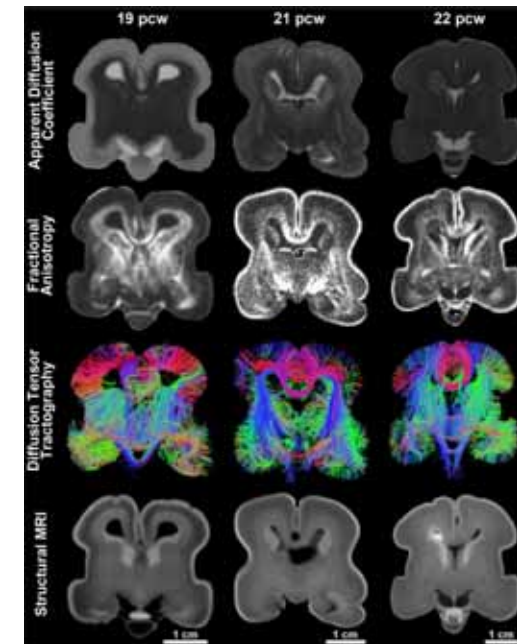
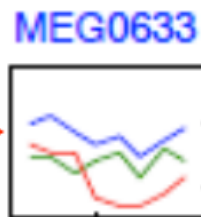
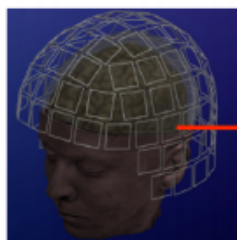
■ Medical Image Data (16 - 1000+ features)



<http://qsota.com/melanoma/>

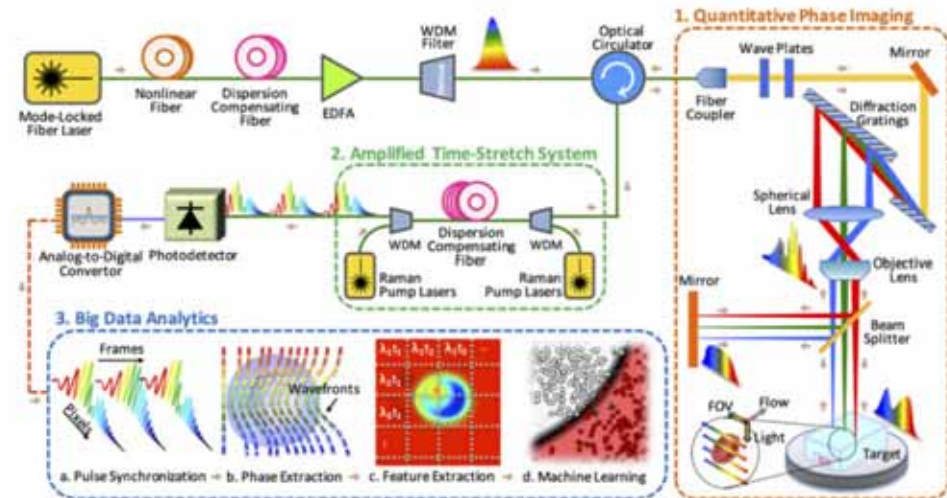
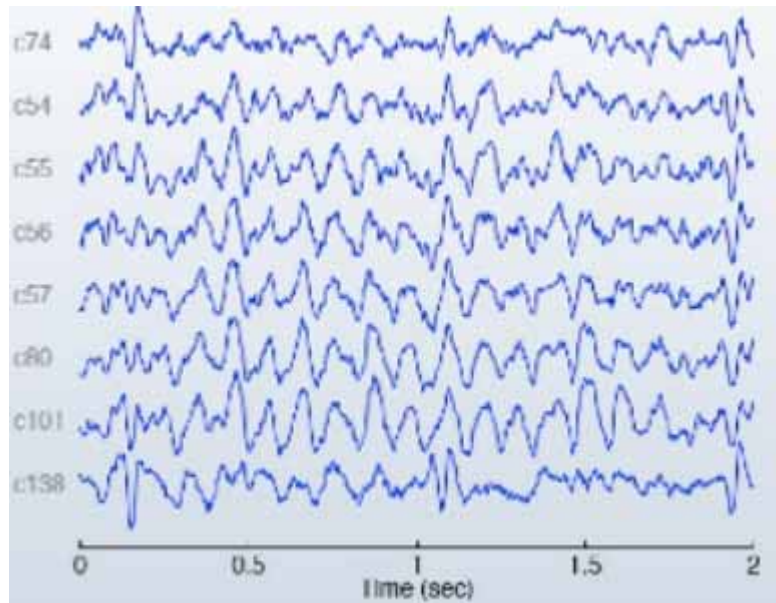
MEG Brain Imaging

120 locations x 500 time points
x 20 objects



Nature 508, 199–206
doi:10.1038/nature13185

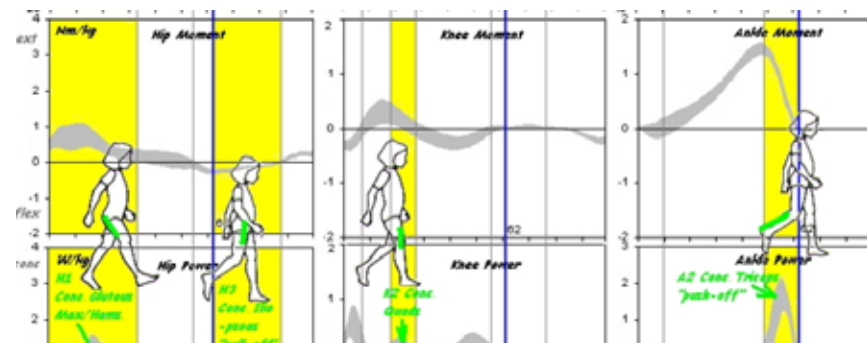
■ Biomedical Signal Data (10 - 1000+ features)



<http://www.nature.com/articles/srep21471#f1>

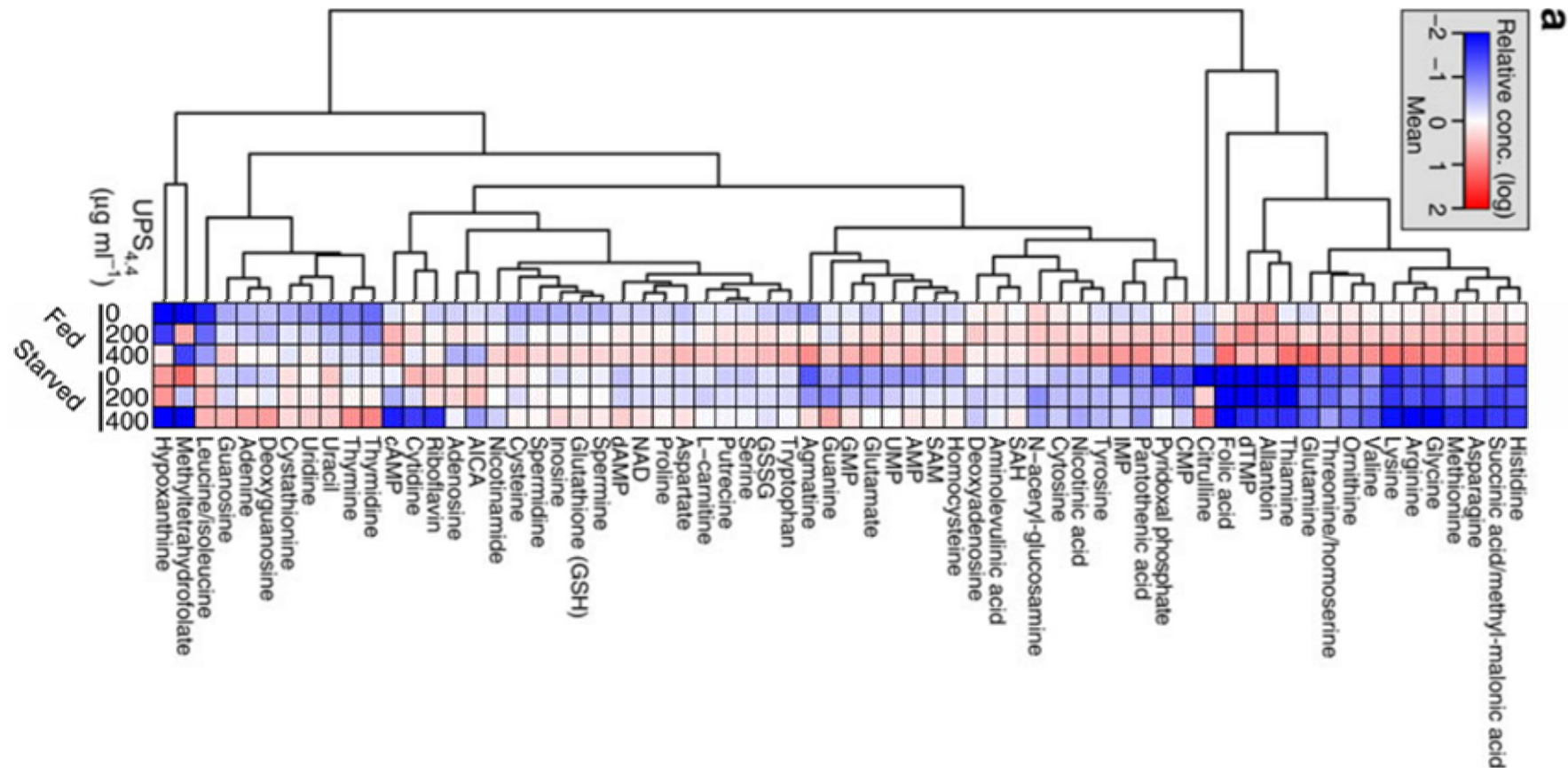


<http://www.mdpi.com/1424-8220/14/4/6124/htm>



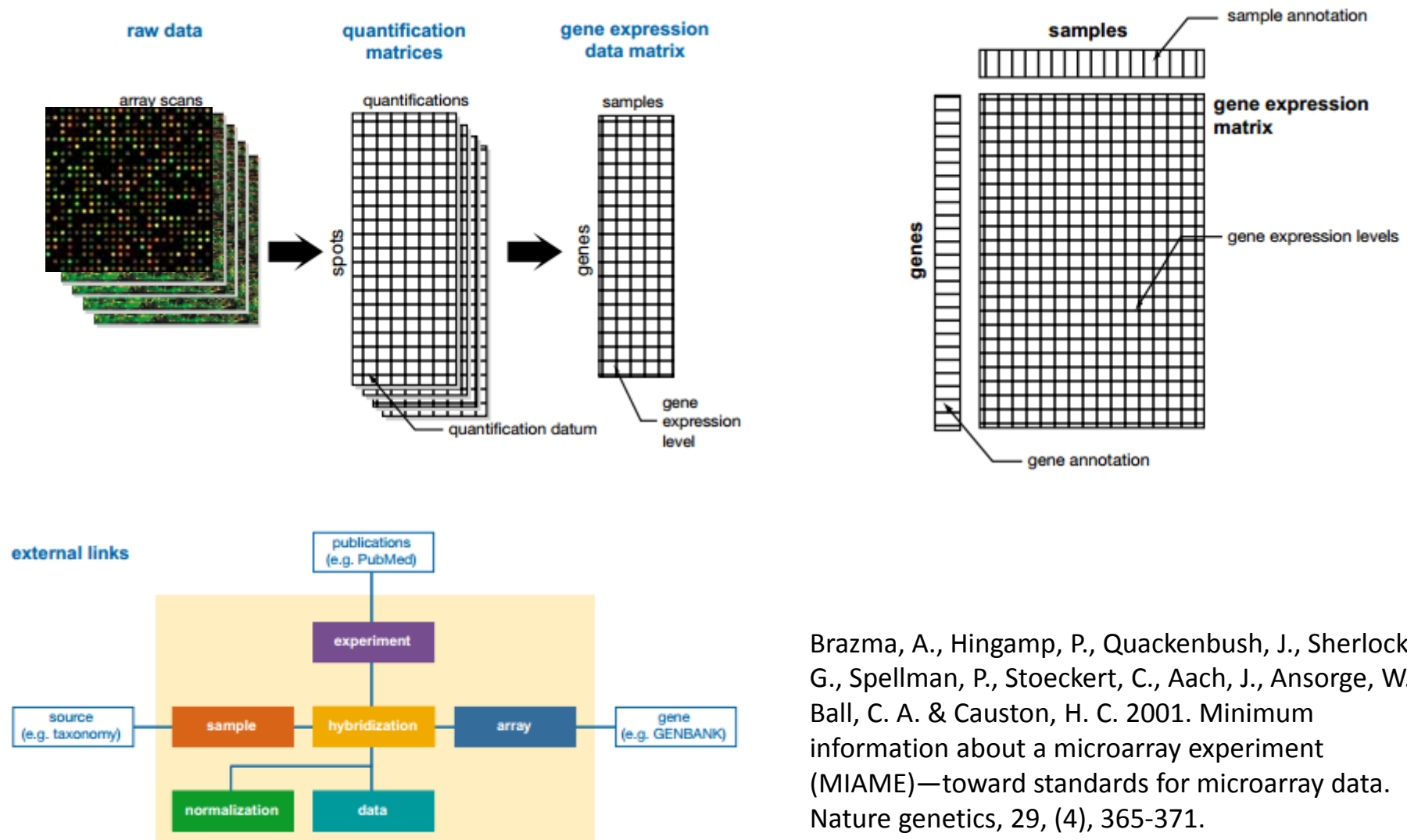
<http://www.clinicalgaitanalysis.com/data/>

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



http://www.nature.com/ncomms/2015/151005/ncomms9524/fig_tab/ncomms9524_F5.html

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

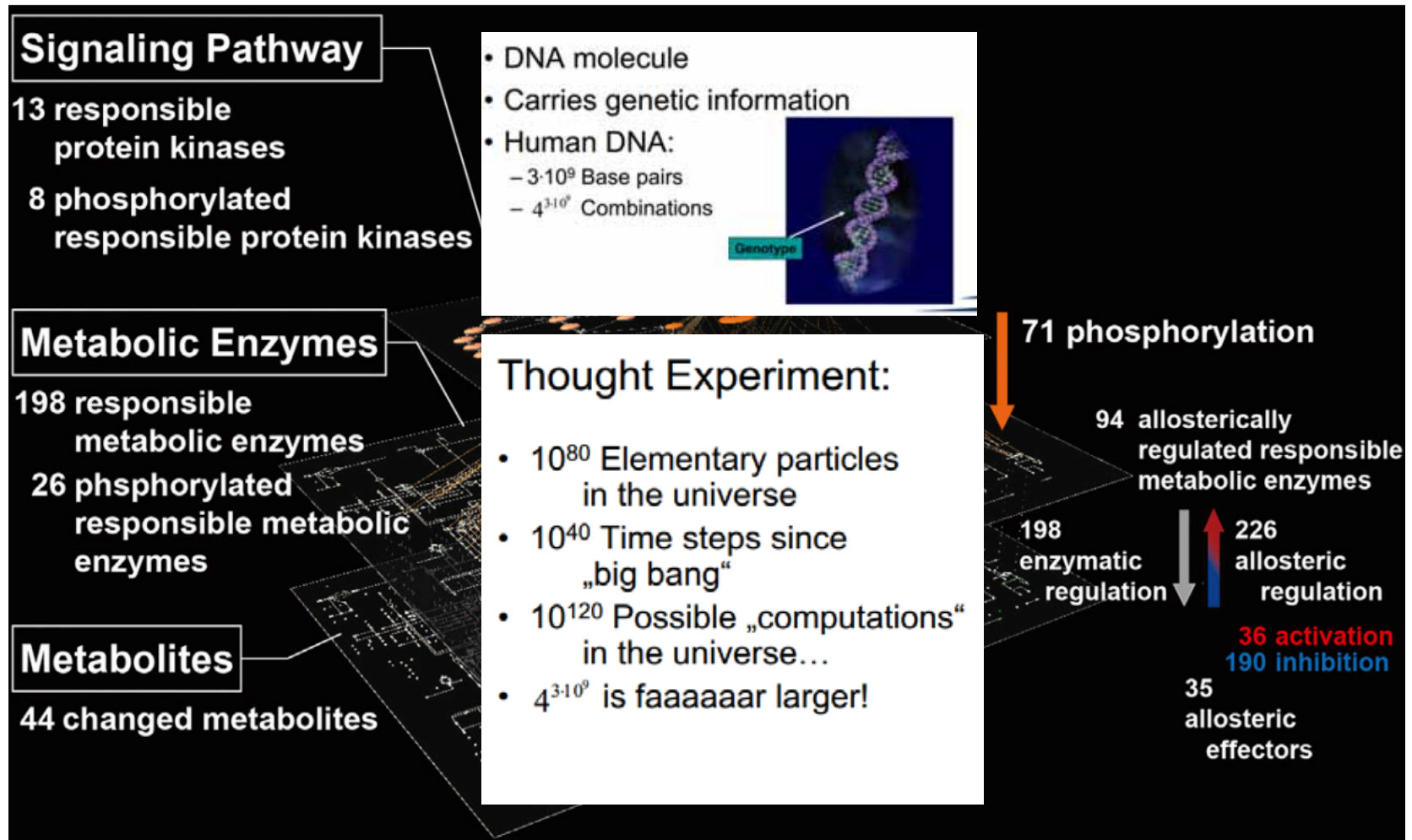


Brazma, A., Hingamp, P., Quackenbush, J., Sherlock, G., Spellman, P., Stoeckert, C., Aach, J., Ansorge, W., Ball, C. A. & Causton, H. C. 2001. Minimum information about a microarray experiment (MIAME)—toward standards for microarray data. *Nature genetics*, 29, (4), 365-371.

- 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 !!!

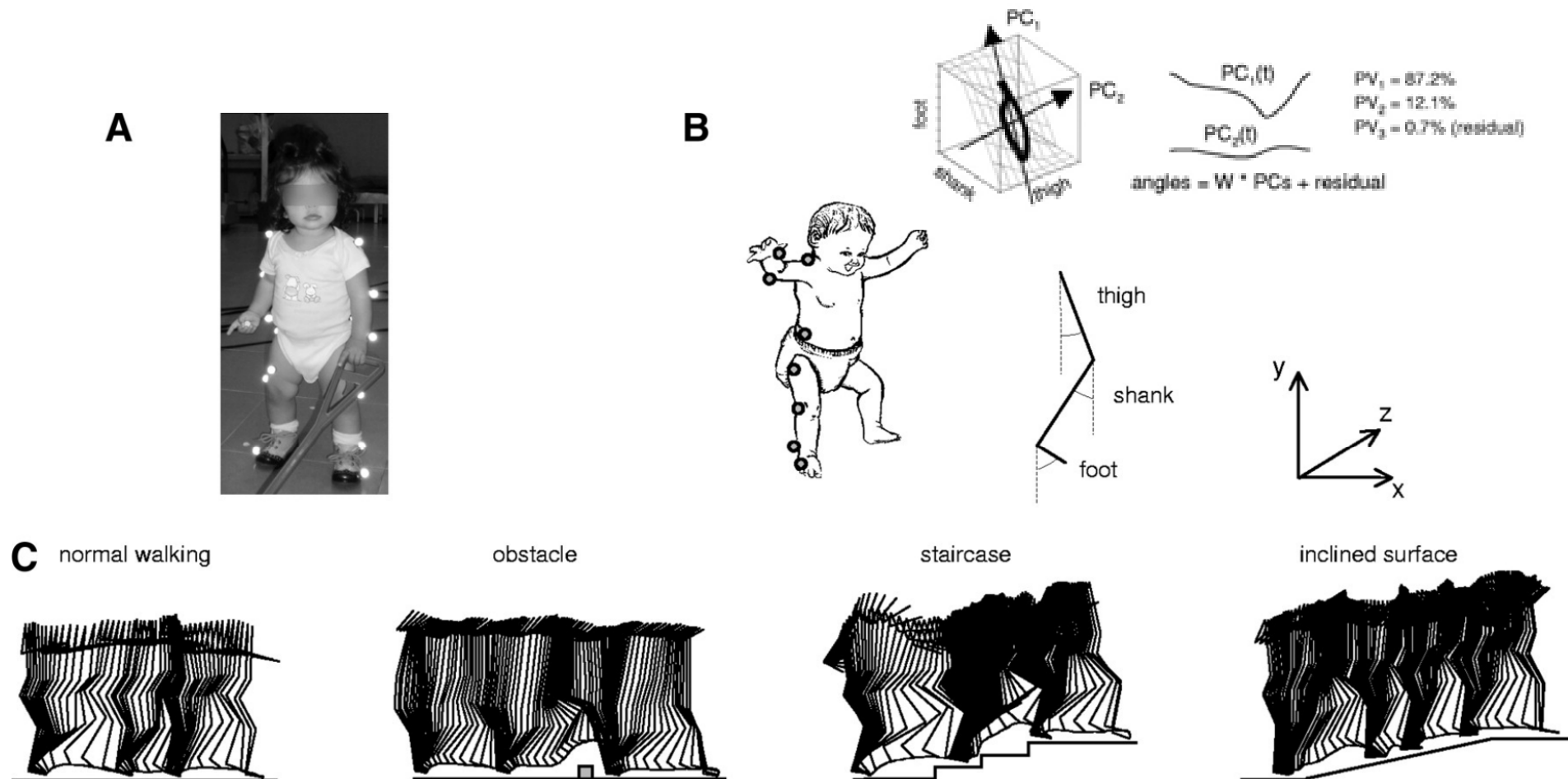


https://www.researchgate.net/publication/255723699_An_Answer_to_Who_Needs_a_Stylus_on_Handwriting_Recognition_on_Mobile_Devices



Yugi, K. et al. 2014. Reconstruction of Insulin Signal Flow from Phosphoproteome and Metabolome Data. Cell Reports, 8, (4), 1171-1183, doi:10.1016/j.celrep.2014.07.021.

- Hyperspace is large – all points are far apart
- Computationally challenging (both in time & 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 (humans are bound to R3/R2!)**
- 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

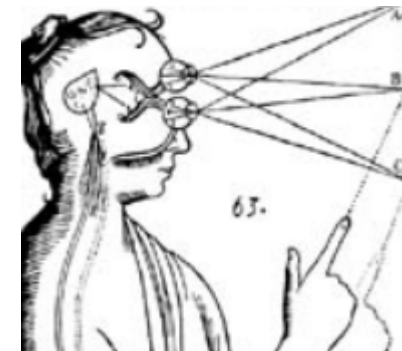
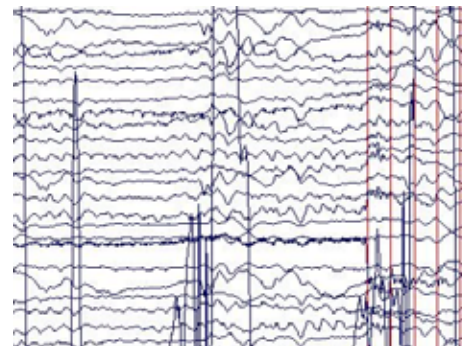
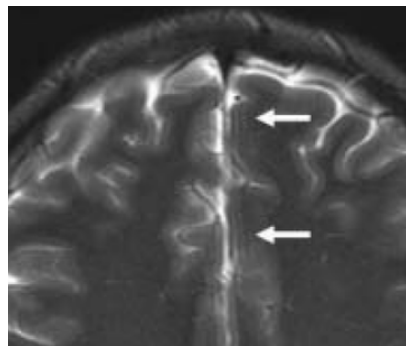


Dominici, N., Ivanenko, Y. P., Cappellini, G., Zampagni, M. L. & Lacquaniti, F. 2010. Kinematic Strategies in Newly Walking Toddlers Stepping Over Different Support Surfaces. *Journal of Neurophysiology*, 103, (3), 1673-1684, doi:10.1152/jn.00945.2009.

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.*

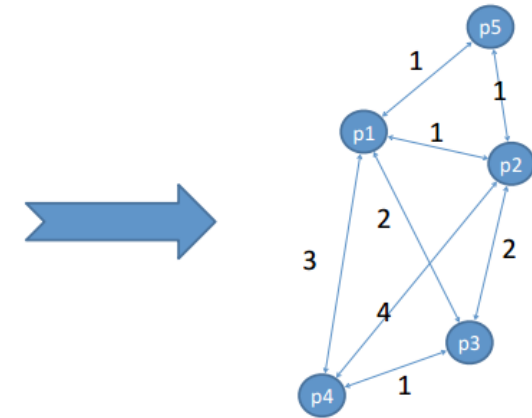
- Given n data points in d dimensions
 - Conversion to m data points in $r \ll d$ dimension
 - Challenge: **minimal loss of information *)**
-
- *) this is always a grand challenge, e.g. in k-Anonymization – see later
 - Very dangerous is the “modeling-of-artifacts”



- Linear methods (unsupervised):
 - PCA (Principal Component Analysis)
 - FA (Factor Analysis)
 - **MDS (Multi-dimensional Scaling)**
- Non-linear methods (unsupervised):
 - Isomap (Isometric feature mapping)
 - LLE (locally linear embedding)
 - Autoencoders
- Supervised methods:
 - LDA (Linear Discriminant Analysis)
- **Subspace Clustering with a human-in-the-loop**

- 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



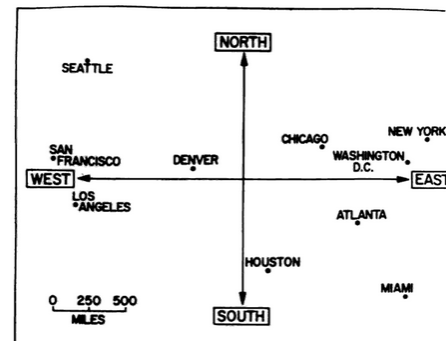
x_i Point in d dimensions

y_i Corresponding point in $r < d$ dimension

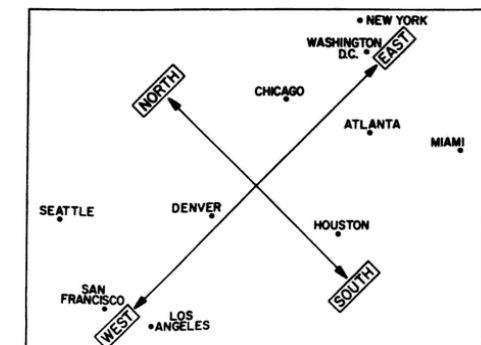
δ_{ij} Distance between x_i and x_j

d_{ij} Distance between y_i and y_j

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



CITIES	ATLA	CHIC	DENV	HOUS	L.A.	MIAMI	N.Y.	S.F.	SEAT	WASH D.C.
ATLANTA		587	1212	701	1936	604	748	2139	2182	543
CHICAGO	587		920	940	1745	1188	713	1858	1737	597
DENVER	1212	920		879	831	1726	1631	949	1021	1494
HOUSTON	701	940	879		1374	968	1420	1645	1891	1220
LOS ANGELES	1936	1745	831	1374		2339	2451	347	959	2300
MIAMI	604	1188	1726	968	2339		1092	2594	2734	923
NEW YORK	748	713	1631	1420	2451	1092		2571	2408	205
SAN FRANCISCO	2139	1858	949	1645	347	2594	2571		678	2442
SEATTLE	2182	1737	1021	1891	959	2734	2408	678		2329
WASHINGTON D.C.	543	597	1494	1220	2300	923	205	2442	2329	



Kruskal, J. B. 1964. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. Psychometrika, 29, (1), 1-27.

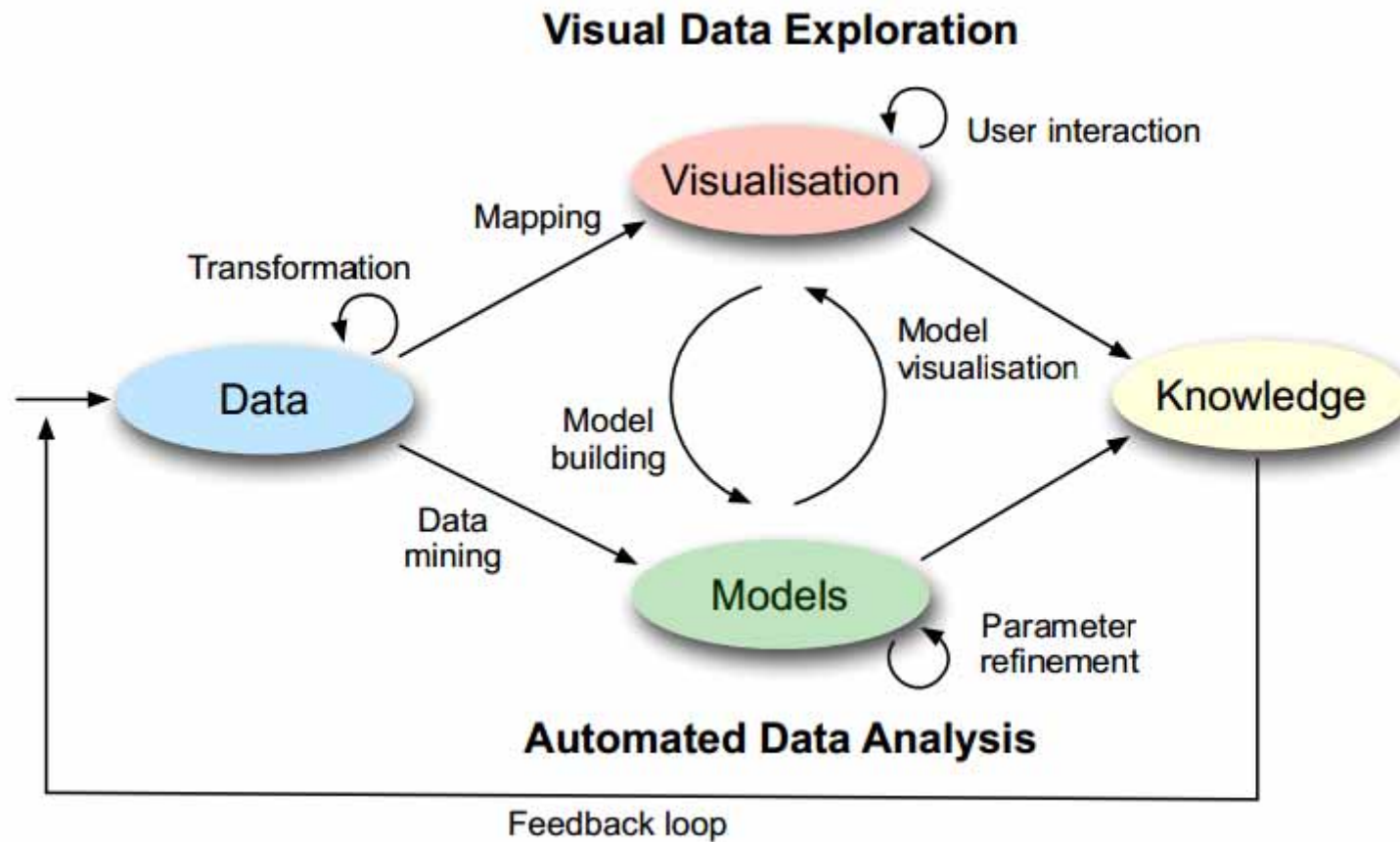
05 Subspace Clustering* and Analysis

* Two major issues

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

- Definitions:
- K clusters
- N data points
- D dimensions (original space)
- d dimensions (latent subspace)
- Subspace Clustering is the process of 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.

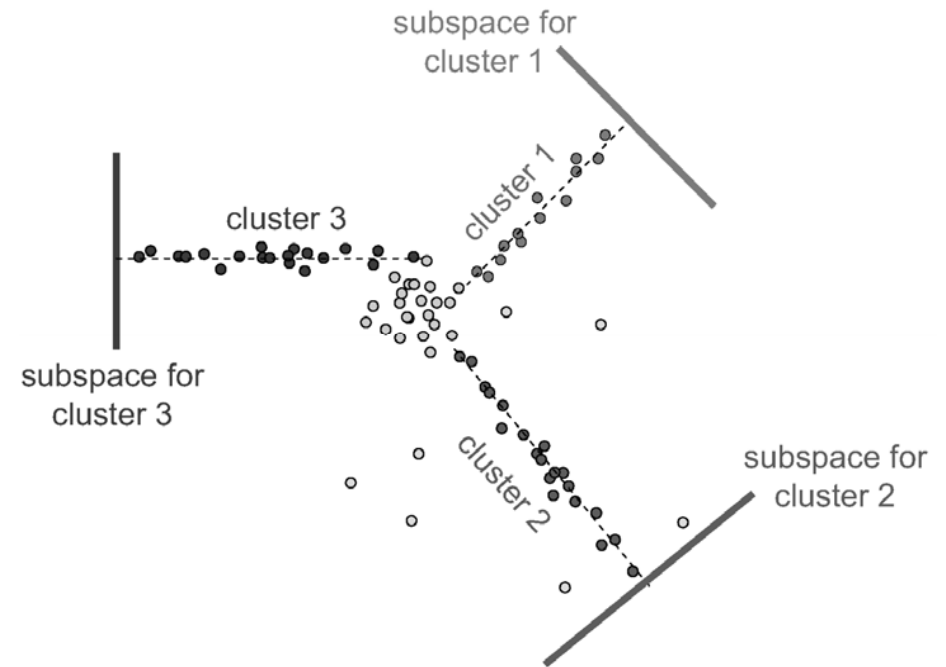
<http://www.vismaster.eu/wp-content/uploads/2010/11/VisMaster-book-lowres.pdf>

Large Amount of Dimensions

Large Amount of Records

Geography	Product Categories																																			
	Bike Racks	Bottles and Cages	Cleaners	Helmets	Hydration Packs	Looks	Pumps	Tires and Tubes	Bike Racks	Bottles and Cages	Cleaners	Helmets	Hydration Packs	Looks	Pumps	Tires and Tubes	Socks	Tights	Vests	Bottom Br	Brakes	Chains	Cranksets	Derailleur	Forks	Handlebars	Headsets	Mountain Pedals	Road Frame	Saddles	Touring Fr	Wheels	Road Bikes			
Virginia	c2,362	c45	c133	c2,021	c476	c180	c84	c12	c636	c39	c231	c1,243	c302	c18	c3	c1	c1,472	c42	c15	c832	c266	c1,684	c151	c31	c277	c1	c24	c709	c252	c7	c2	c1	c76	c763		
Arizona	c2,209	c61	c39	c1,881	c806	c75	c72	c19	c71	c716	c561	c894	c237	c80	c15	c5	c519	c20	c37	c145	c87	c194	c32	c180	c68	c44	c79	c758	c152	c59	c5	c1	c25	c727		
Colorado	c4,153	c148	c262	c4,326	c1,631	c165	c228	c12	c239	c372	c1,430	c1,017	c1,352	c136	c10	c10	c2,808	c117	c139	c1,500	c149	c1,447	c1,706	c81	c19	c1	c225	c326	c1,194	c53	c5	c3	c1,477	c727		
Florida	c4,422	c182	c206	c3,848	c1,068	c180	c144	c33	c1,941	c889	c1,208	c1,133	c987	c109	c23	c6	c2,128	c270	c383	c3,843	c119	c406	c1,029	c315	c764	c147	c54	c101	c72	c21	c1	c4	c703	c533		
Illinois	c576	c27	c33	c489	c237	c45				c12	c1,389	c195	c187	c672	c549	c581	c309	c18	c332	c220	c130	c3,032	c1,131	c2,410	c1,239	c188	c1,958	c26	c68	c990	c1,766	c4,598	c2,714	c194	c8,000	c577
Indiana	c1,250	c33	c32	c1,330	c474	c45	c24	c14	c334	c48	c458	c649	c136	c23	c44	c22	c278	c36	c167	c153	c231	c82	c176	c62	c78	c9	c86	c925	c207	c36	c6	c1	c87	c578		
Maine	c2,069	c69	c137	c1,948	c507	c60	c132	c12	c372	c259	c701	c476	c324	c49	c242	c8	c1,375	c42	c1	c2,326	c460	c545	c463	c119	c270	c21	c162	c40	c445	c40	c86	c11	c5,000	c7,500		
Michigan	c2,421	c88	c140	c2,842	c691	c60	c94	c22	c478	c24	c816	c351	c723	c33	c57	c30	c1,589	c164	c6	c2,512	c669	c631	c252	c326	c194	c0	c478	c1,186	c1,509	c62	c75	c17	c4,159	c9,500		
Missouri	c1,368	c63	c91	c1,140	c660	c75	c60			c483	c193	c406	c309	c25	c23	c13	c111	c6	c58	c2,170	c1,106	c580	c297	c228	c397	c26	c255	c623	c651	c2	c24	c50	c340	c6,800		
Nevada	c1,656	c122	c149	c1,621	c738					c12	c1,389	c195	c187	c672	c549	c581	c309	c18	c332	c220	c130	c3,032	c1,131	c2,410	c1,239	c188	c1,958	c26	c68	c990	c1,766	c4,598	c2,714	c194	c8,000	c577
New Mexico	c1,531	c56	c133	c1,396	c594	c105	c48	c14	c337	c129	c742	c136	c323	c64	c48	c6	c291	c3	c212	c1,904	c108	c571	c368	c159	c240	c3	c348	c183	c823	c525	c105	c79	c9,311	c2,219		
New York	c3,217	c185	c312	c4,317	c2,070	c165	c108	c43	c1,571	c829	c429	c3,362	c1,461	c101	c163	c14	c3,228	c201	c127	c2,265	c2,325	c731	c780	c509	c97	c46	c128	c4,048	c4,593	c1,416	c1,500	c253	c4,003	c2,219		
Ohio	c1,656	c51	c67	c1,091	c462					c1	c1,249	c194	c286	c487	c266	c0	c0	c1	c454	c101	c31	c146	c361	c0	c0	c0	c850	c4	c30	c1,020	c466	c0	c0	c11	c2,450	c1,767
Virginia	c289	c24	c70	c1,799	c518	c328	c74	c31	c91	c126	c274	c334	c111	c73	c18	c20	c576	c34	c44	c1,187	c273	c1,106	c245	c120	c7	c8	c63	c287	c249	c464	c42	c421	c3,532	c4,289		
Arizona	c1,527	c23	c90	c2,326	c178	c183	c178	c40	c960	c132	c81	c125	c31	c22	c75	c3	c517	c27	c132	c2,368	c19	c253	c257	c40	c844	c3	c15	c97	c63	c253	c630	c987	c3,334	c3,379		
Colorado	c169	c143	c101	c1,225	c1,420	c102	c160	c34	c18	c378	c524	c2,038	c240	c26	c55	c6	c72	c23	c39	c3,055	c602	c98	c437	c113	c12	c2	c48	c2,188	c312	c40	c193	c490	c421	c451		
Florida	c3,567	c93	c366	c3,442	c2,180	c402	c42	c89	c2,737	c386	c1,302	c392	c711	c116	c12	c16	c3,234	c118	c287	c42	c616	c2,382	c1,116	c166	c730	c3	c64	c496	c764	c841	c966	c775	c3,532	c68		
Illinois	c1,376	c150	c145	c1,721	c219	c738	c81	c51	c1,006	c3	c682	c36	c35	c43	c46	c23	c1,405	c1	c36	c168	c72	c742	c57	c55	c75	c0	c246	c136	c57	c211	c235	c275	c4,286	c1		
Indiana	c38	c8	c20	c334	c1,075	c18	c4	c15	c19	c45	c82	c38	c197	c3	c0	c3	c8	c2	c0	c53	c541	c83	c6	c49	c5	c7	c23	c11	c2	c2	c3	c19	c19	c763		
Maine	c430	c9	c22	c558	c742	c79	c5	c10	c214	c43	c49	c43	c130	c20	c3	c1	c100	c3	c28	c203	c448	c109	c2	c9	c33	c7	c9	c33	c37	c34	c2	c1	c181	c727		
Michigan	c1,615	c596	c0	c2,498	c533	c48	c165	c32	c1,473	c65	c3	c133	c149	c2	c26	c6	c595	c16	c0	c221	c76	c113	c11	c75	c467	c10	c1	c81	c154	c1	c54	c28	c1,060	c727		
Missouri	c687	c7	c54	c1,241	c348	c151	c87	c22	c467	c2	c166	c397	c24	c16	c20	c4	c689	c4	c46	c393	c91	c220	c136	c63	c445	c0	c1	c632	c30	c12	c44	c1	c1	c2,572	c533	
Nevada	c372	c115	c29	c3,375	c84	c2,099	c465	c17	c355	c412	c34	c36	c27	c307	c33	c3	c116	c4	c5	c366	c13	c3,842	c1,427	c50	c24	c52	c3	c35	c10	c200	c29	c2	c2,275	c3,997		
New Mexico	c1,531	c56	c133	c1,396	c594	c105	c48	c14	c337	c129	c742	c136	c323	c64	c48	c6	c291	c3	c212	c1,904	c108	c571	c368	c159	c240	c3	c348	c183	c823	c525	c105	c79	c9,311	c2,219		
New York	c3,217	c185	c312	c4,317	c2,070	c165	c108	c43	c1,571	c829	c429	c3,362	c1,461	c101	c163	c14	c3,228	c201	c127	c2,265	c2,325	c731	c780	c509	c97	c46	c128	c4,048	c4,593	c1,416	c1,500	c253	c4,003	c2,219		
Ohio	c1,656	c51	c67	c1,091	c462					c1	c1,249	c194	c286	c487	c266	c0	c0	c1	c454	c101	c31	c146	c361	c0	c0	c0	c850	c4	c30	c1,020	c466	c0	c0	c11	c2,450	c1,767
Virginia	c289	c24	c70	c1,799	c518	c328	c74	c31	c91	c126	c274	c334	c111	c73	c18	c20	c576	c34	c44	c1,187	c273	c1,106	c245	c120	c7	c8	c63	c287	c249	c464	c42	c421	c3,532	c4,289		
Arizona	c1,527	c23	c90	c2,326	c178	c183	c178	c40	c960	c132	c81	c125	c31	c22	c75	c3	c517	c27	c132	c2,368	c19	c253	c257	c40	c844	c3	c15	c97	c63	c253	c630	c987	c3,334	c3,379		
Colorado	c169	c143	c101	c1,225	c1,420	c102	c160	c34	c18	c378	c524	c2,038	c240	c26	c55	c6	c72	c23	c39	c3,055	c602	c98	c437	c113	c12	c2	c48	c2,188	c312	c40	c193	c490	c421	c451		
Florida	c3,567	c93	c366	c3,442	c2,180	c402	c42	c89	c2,737	c386	c1,302	c392	c711	c116	c12	c16	c3,234	c118	c287	c42	c616	c2,382	c1,116	c166	c730	c3	c64	c496	c764	c841	c966	c775	c3,532	c68		
Illinois	c1,376	c150	c145	c1,721	c219	c738	c81	c51	c1,006	c3	c682	c36	c35	c43	c46	c23	c1,405	c1	c36	c168	c72	c742	c57	c55	c75	c0	c246	c136	c57	c211	c235	c275	c4,286	c1		
Indiana	c38	c8	c20	c334	c1,075	c18	c4	c15	c19	c45	c82	c38	c197	c3	c0	c3	c8	c2	c0	c53	c541	c83	c6	c49	c5	c7	c23	c11	c2	c2	c3	c19	c19	c763		
Maine	c430	c9	c22	c558	c742	c79	c5	c10	c214	c43	c49	c43	c130	c20	c3	c1	c100	c3	c28	c203	c448	c109	c2	c9	c33	c7	c9	c33	c37	c34	c2	c1	c181	c727		
Michigan	c1,615	c596	c0	c2,498	c533	c48	c165	c32	c1,473	c65	c3	c133	c149	c2	c26	c6	c595	c16	c0	c221	c76	c113	c11	c75	c467	c10	c1	c81	c154	c1	c54	c28	c1,060	c727		
Missouri	c687	c7	c54	c1,241	c348	c151	c87	c22	c467	c2	c166	c397	c24	c16	c20	c4	c689	c4	c46	c393	c91	c220	c136	c63	c445	c0	c1	c632	c30	c12	c44	c1	c1	c2,572	c533	
Nevada	c372	c115	c29	c3,375	c84	c2,099	c465	c17	c355	c412	c34	c36	c27	c307	c33	c3	c116	c4	c5	c366	c13	c3,842	c1,427	c50	c24	c52	c3	c35	c10	c200	c29	c2	c2,275	c3,997		
New Mexico	c1,531	c56	c133	c1,396	c594	c105	c48	c14	c337	c129	c742	c136	c323	c64	c48	c6	c291	c3	c212	c1,904	c108	c571	c368	c159	c240	c3	c348	c183	c823	c525	c105	c79	c9,311	c2,219		
New York	c3,217	c185	c312	c4,317	c2,070	c165	c108	c43	c1,571	c829	c429	c3,362	c1,461	c101	c163	c14	c3,228	c201	c127	c2,265	c2,325	c731	c780	c509	c97	c46	c128	c4,048	c4,593	c1,416	c1,500	c253	c4,003	c2,219		
Ohio	c1,656	c51	c67	c1,091	c462					c1	c1,249	c194	c286	c487	c266	c0	c0	c1	c454	c101	c31	c146	c361	c0	c0	c0	c850	c4	c30	c1,020	c466	c0	c0	c11	c2,450	c1,767
Virginia	c289	c24	c70	c1,799	c518	c328	c74	c31	c91	c126	c274	c334	c111	c73	c18	c20	c576	c34	c44	c1,187	c273	c1,106	c245	c120	c7	c8	c63	c287	c249	c464	c42	c421	c3,532	c4,289		
Arizona	c1,527	c23	c90	c2																																

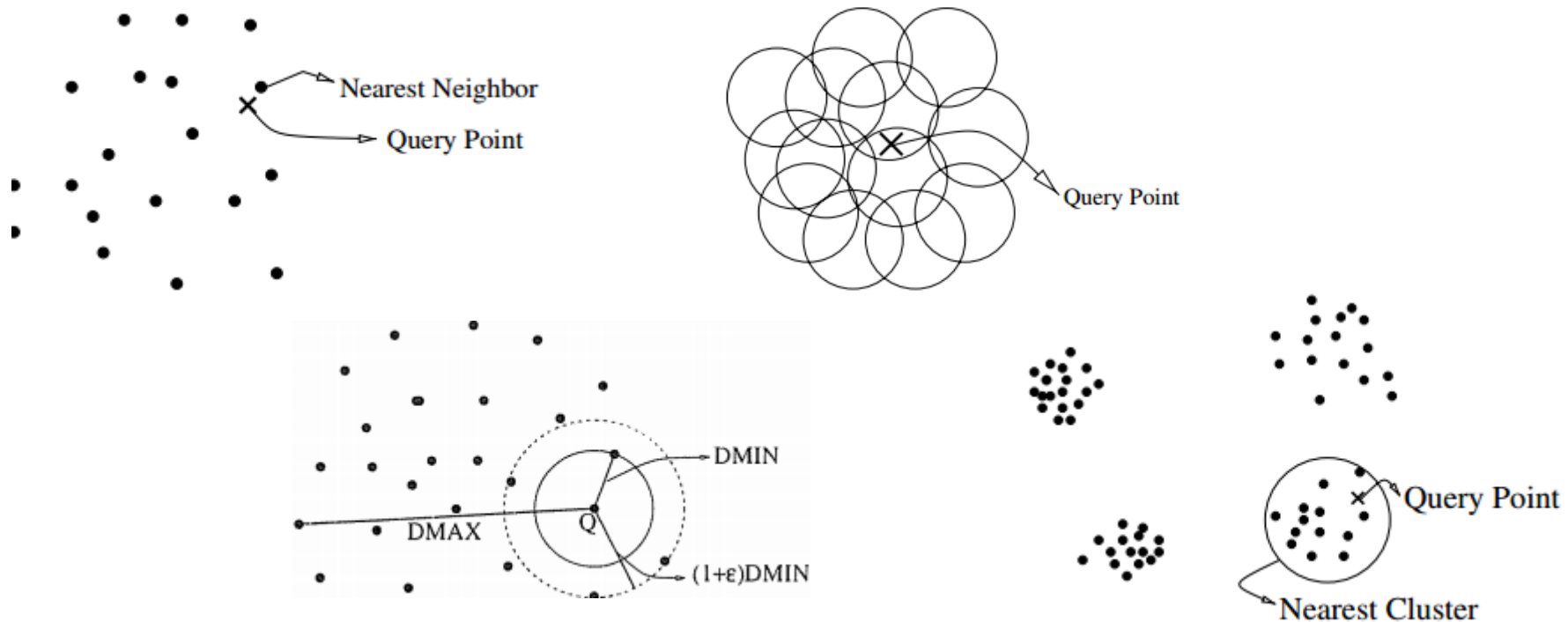
- Many irrelevant dimensions
- Correlated and redundant dimensions
- Conflicting dimensions
- Wrong Interpretation of global analysis results



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.

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.

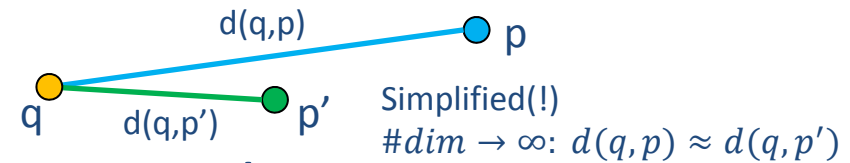
- NN problem: Given n data points and a query point in an m –dimensional metric space
- find the data point closest to the query point.



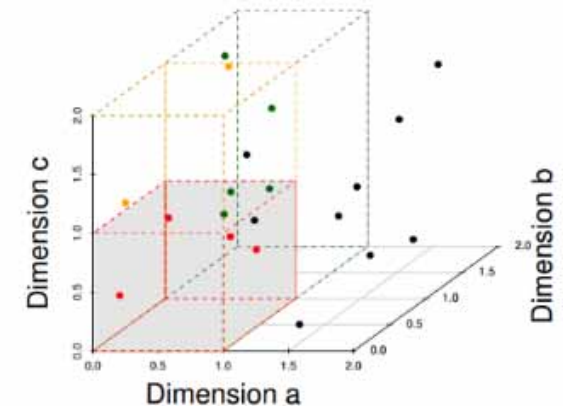
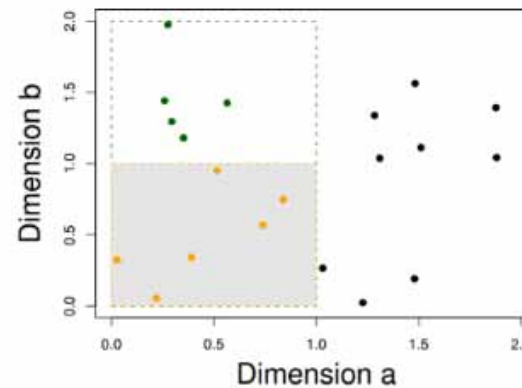
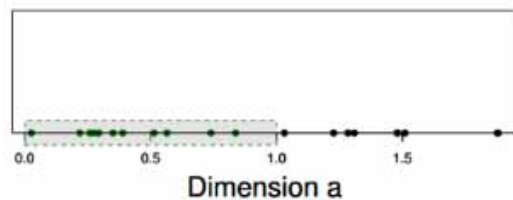
Beyer, K., Goldstein, J., Ramakrishnan, R. & Shaft, U. 1999. When is "nearest neighbor" meaningful? *In: Beer, C. & Buneman, P. (eds.) Database Theory ICDT 99, LNCS 1540.* Berlin: Springer, pp. 217-235.

- Concentration Effect

- Discriminability of similarity gets lost
- Impact on usefulness of a similarity measure



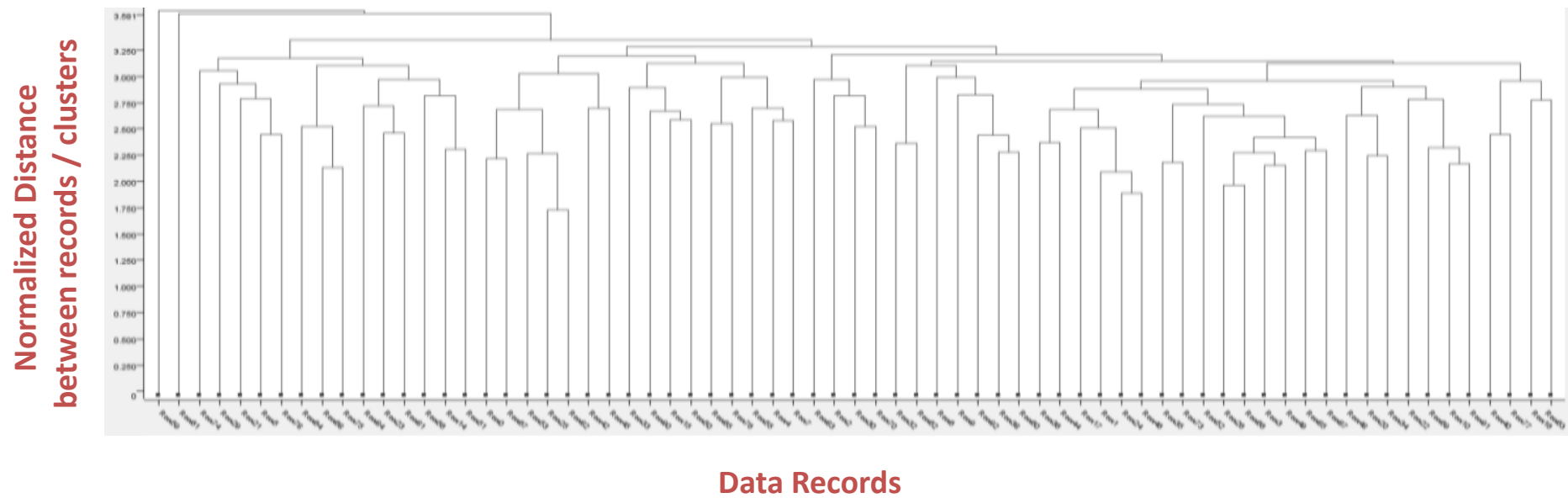
- High-Dimensional Data is Sparse

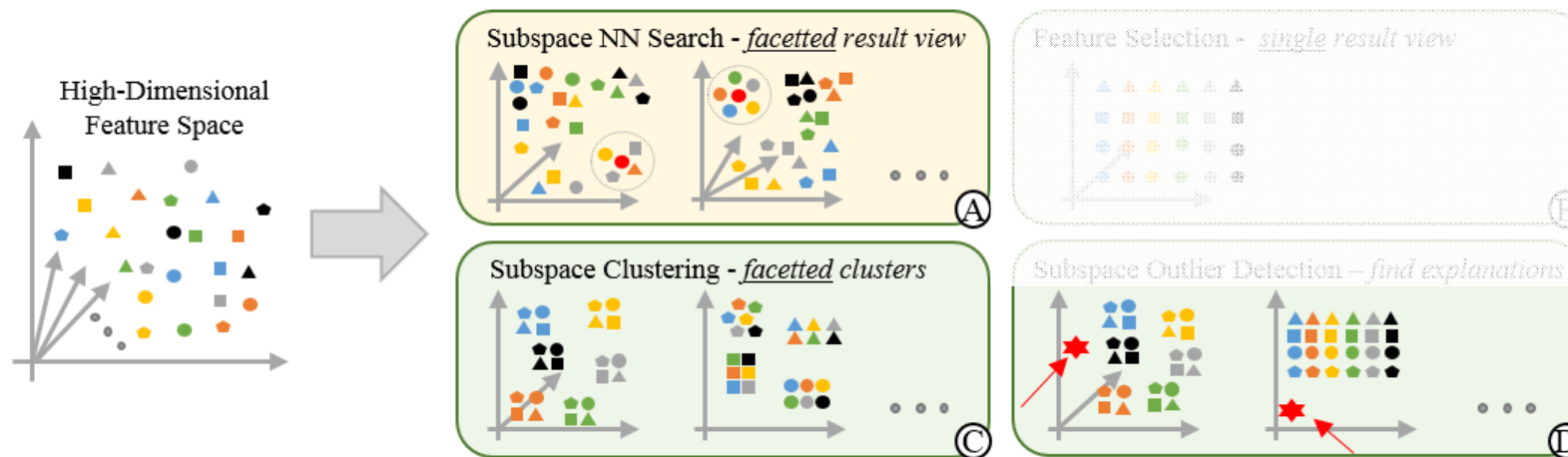
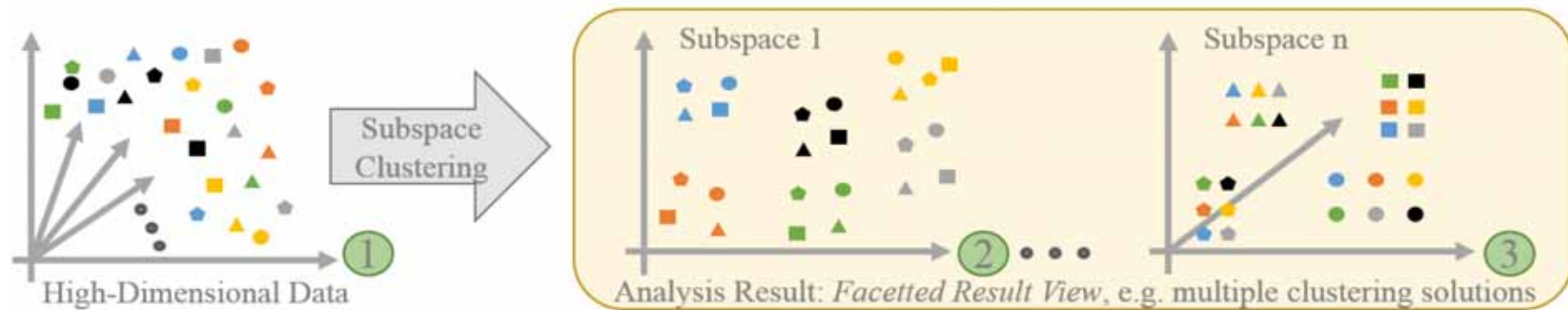


Optimization Problem and Combinatorial Issues

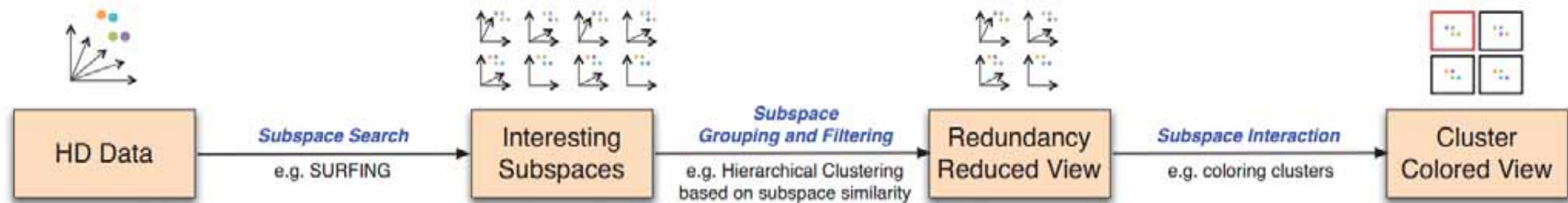
Feature selection and dimension reduction

$2^d - 1$ possible subsets of dimensions (\rightarrow subspaces)

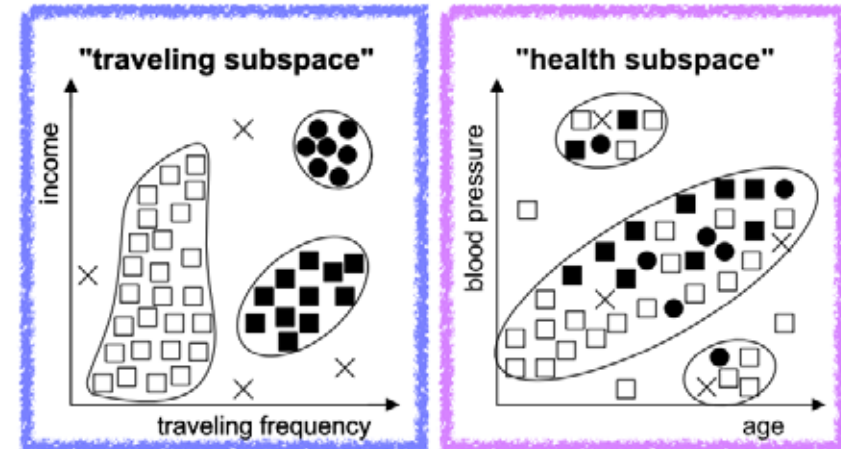




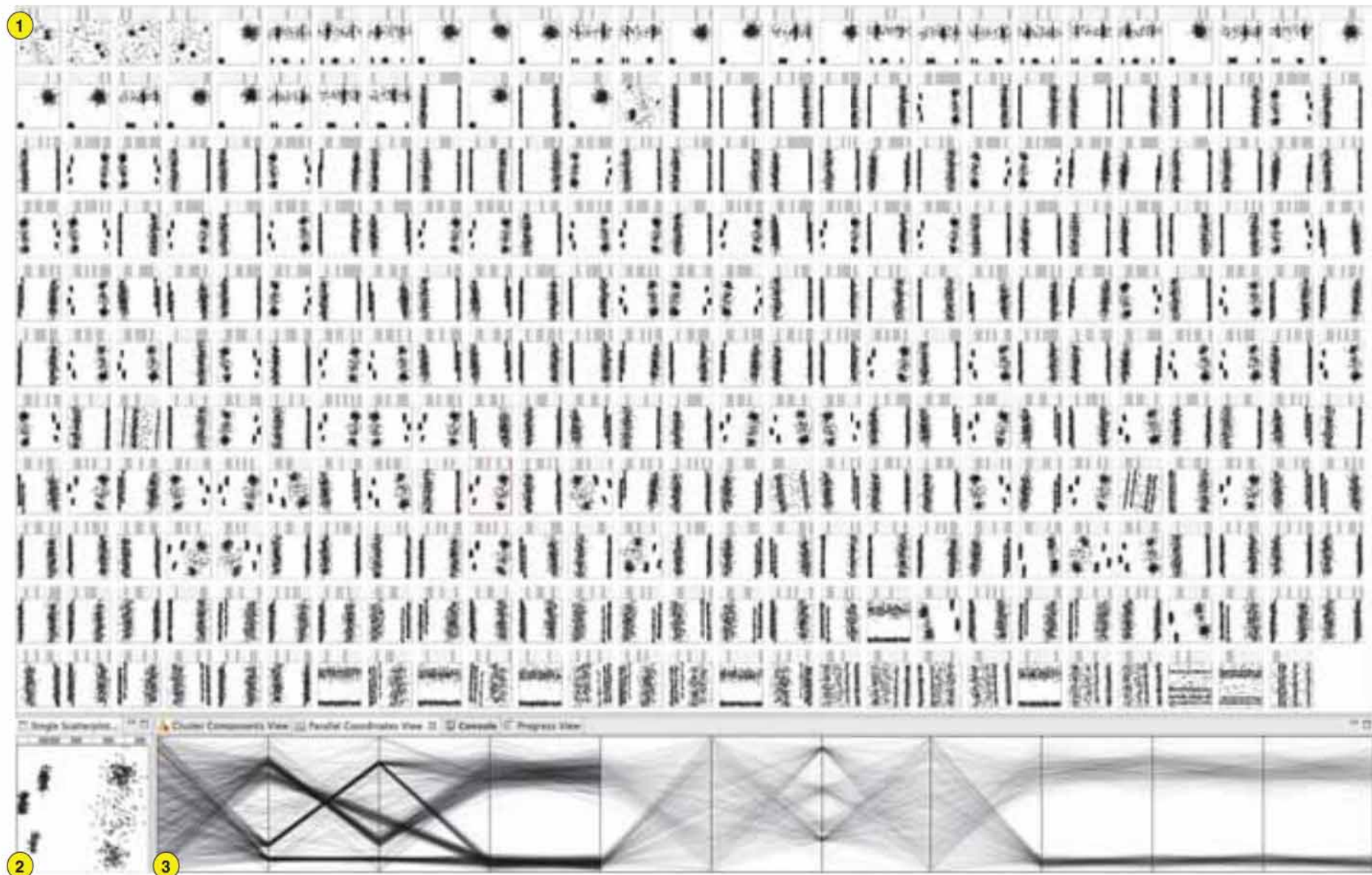
- Patterns may be found in subspaces (dimension combinations)
- Patterns may be complementary or redundant to each other

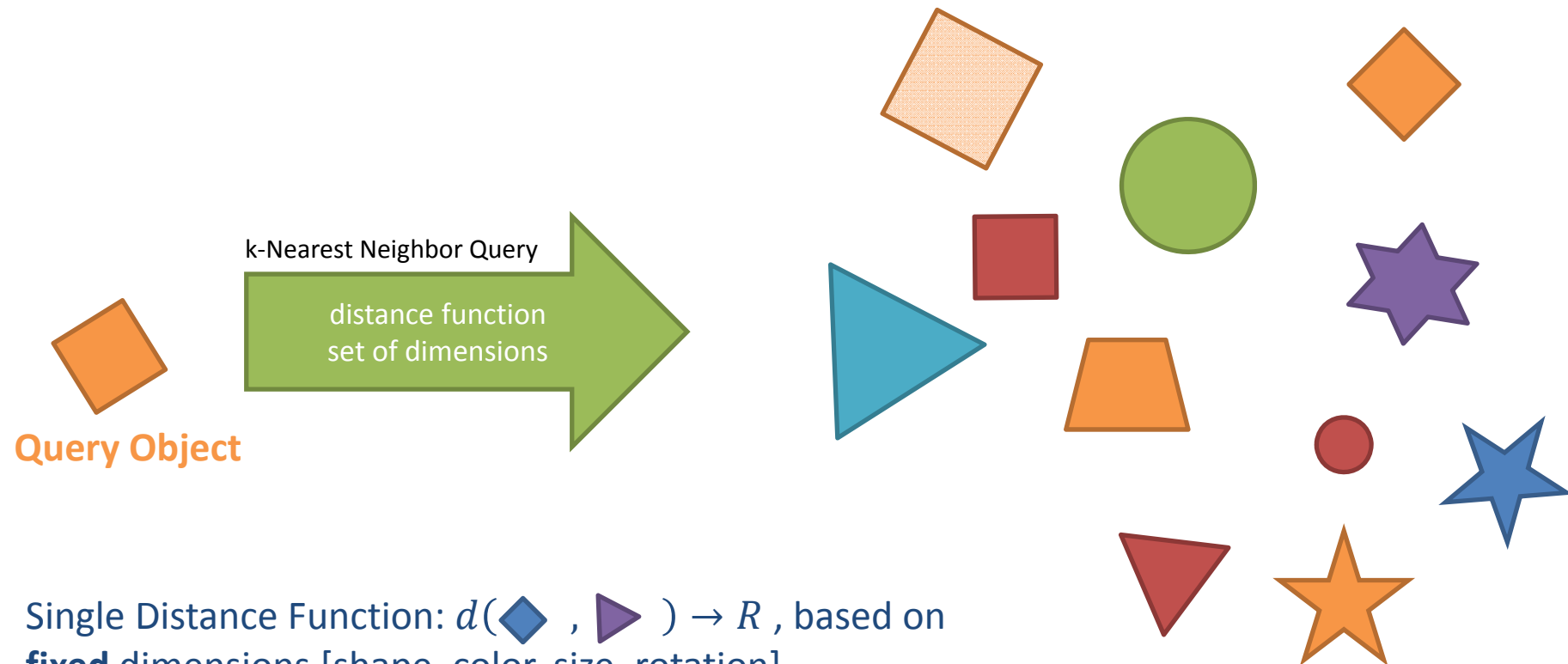


objectID	age	blood pres.	sportactiv	income	trav. freq.
1	ABC	ABC	ABC	ABC	ABC
2	ABC	ABC	ABC	ABC	ABC
3	ABC	ABC	ABC	ABC	ABC
4	ABC	ABC	ABC	ABC	ABC
5	ABC	ABC	ABC	ABC	ABC
6	ABC	ABC	ABC	ABC	ABC
7	ABC	ABC	ABC	ABC	ABC
8	ABC	ABC	ABC	ABC	ABC
9	ABC	ABC	ABC	ABC	ABC

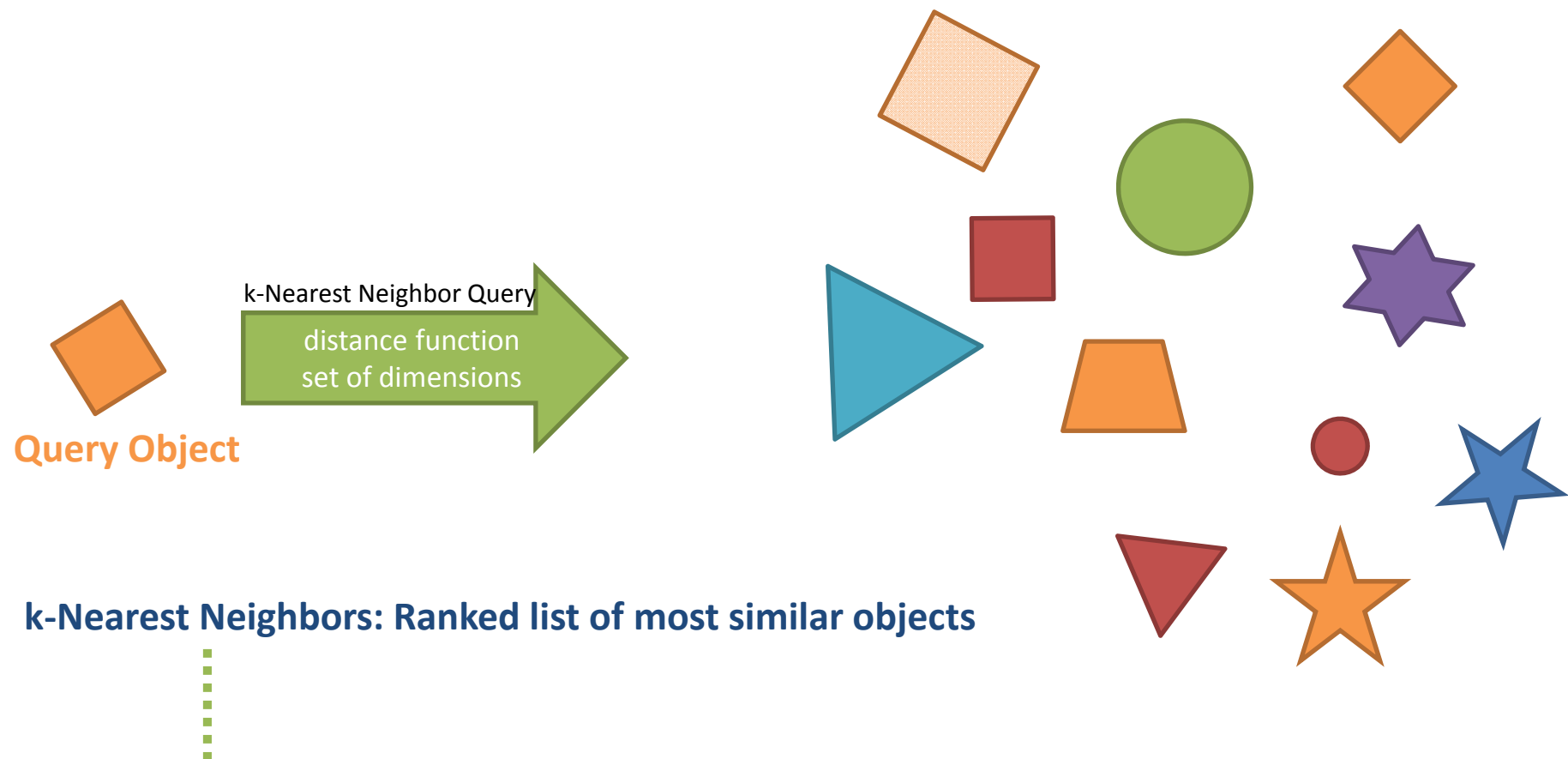


Tatu, A., Maass, F., Faerber, I., Bertini, E., Schreck, T., Seidl, T. & Keim, D. Subspace search and visualization to make sense of alternative clusterings in high-dimensional data. IEEE Symposium on Visual Analytics Science and Technology (VAST), 2012 Seattle. IEEE, 63-72, doi:10.1109/VAST.2012.6400488.

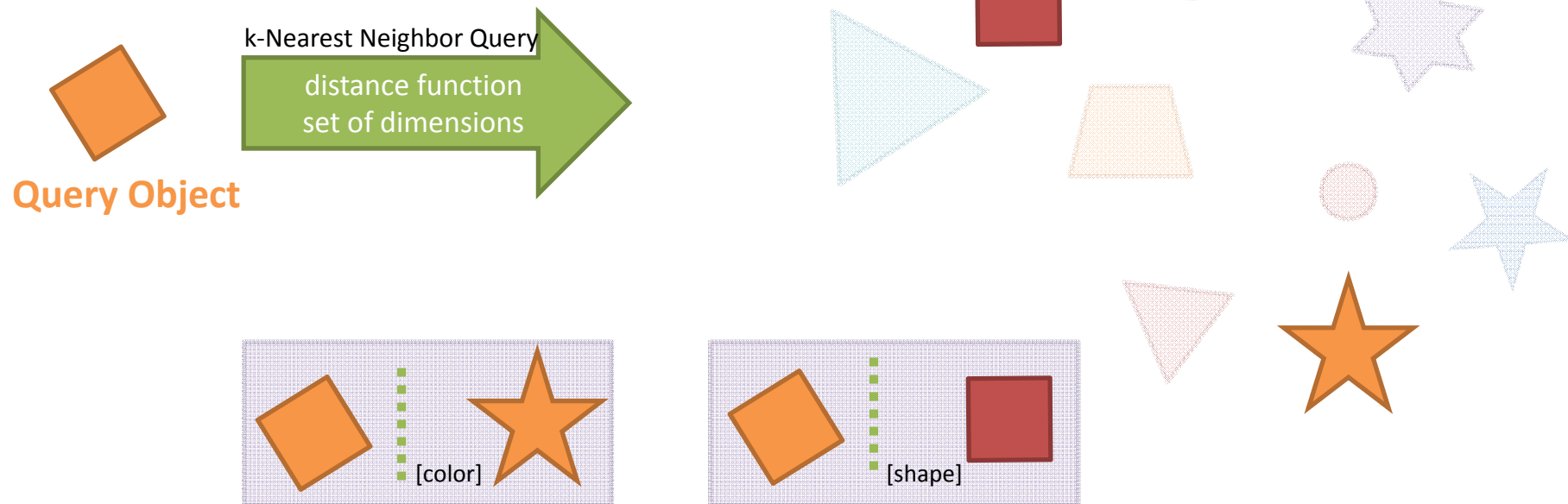




Hund, M., Behrisch, M., Färber, I., Sedlmair, M., Schreck, T., Seidl, T. & Keim, D. 2015. Subspace Nearest Neighbor Search-Problem Statement, Approaches, and Discussion. *Similarity Search and Applications*. Springer, pp. 307-313.



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





Nearest Neighbor
Search

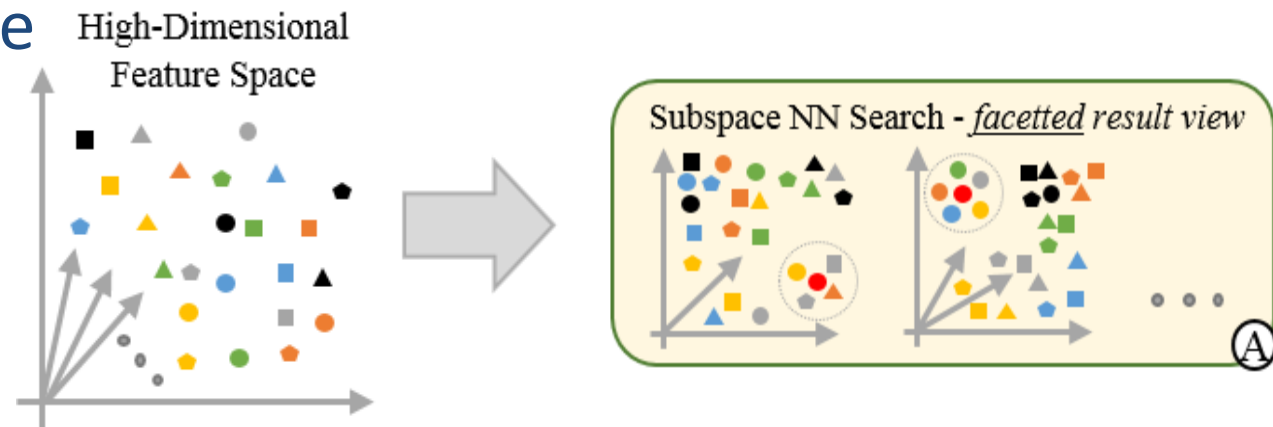


- (1) Relevant subspaces *depend on the patient* and are *unknown* beforehand
- (2) *Multiple* subspaces might be relevant
- (3) Subspaces helps to *interpret* the nearest neighbors (*semantic* meaning)

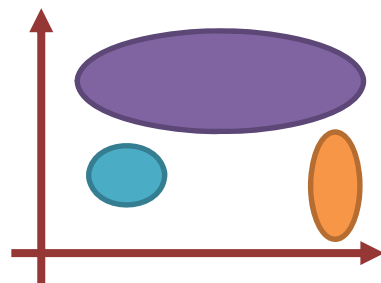
Sex, Age, Blood Type,
Blood Pressure,
Former Diseases,
Medication, ...

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

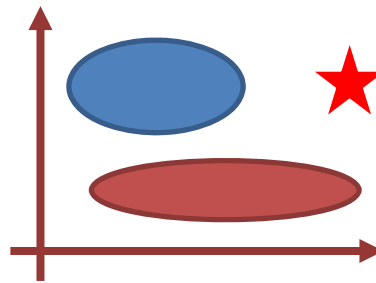
Characteristics:



- 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



Subspace Outlier Detection



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.

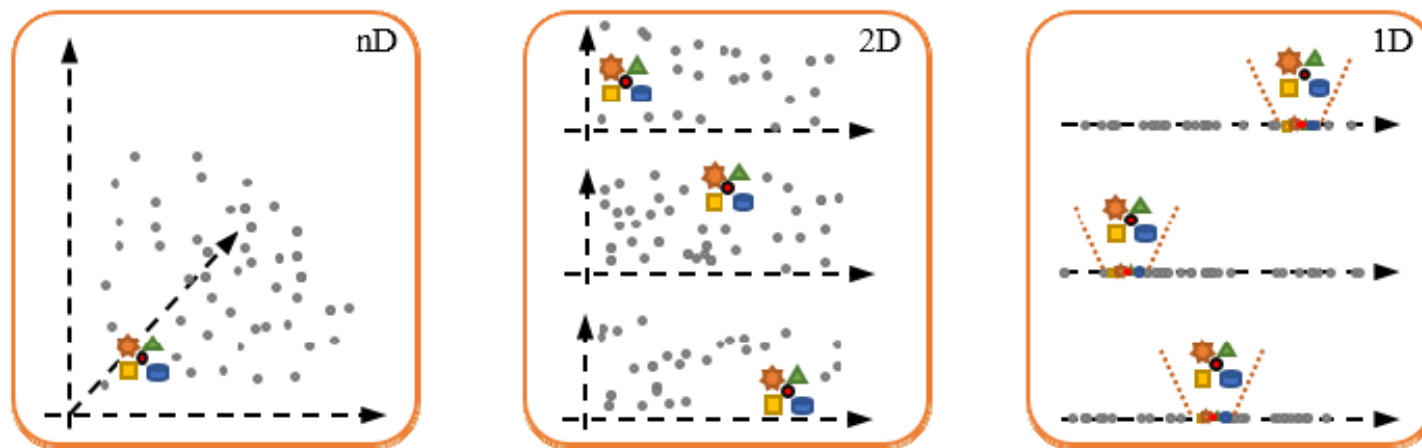
[2] Zimek, A., Schubert, E. & Kriegel, H. P. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5, (5), 363-387.

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

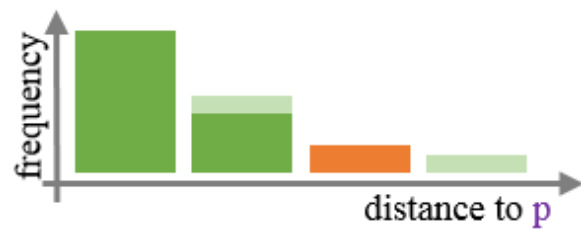
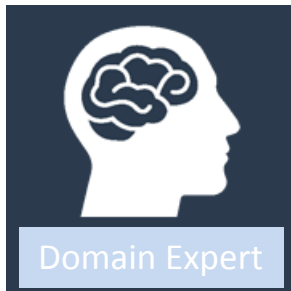
A subspace is considered **relevant**, iff it contains relevant nearest neighbors



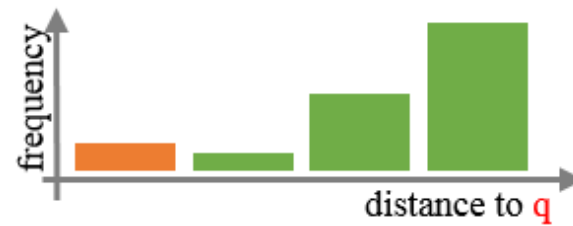
Dimensionality

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.

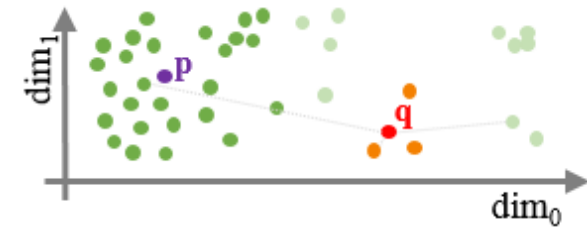
- **Interpretability: reflects the semantic meaning**
 - In which way are NN's similar to the query?
 - → In all dimensions of the subspace
- **Fulfills the downward-closure property**
 - Make use of *Apriori-like algorithms* for subspace search
- **No global distance function necessary**
 - Heterogeneous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



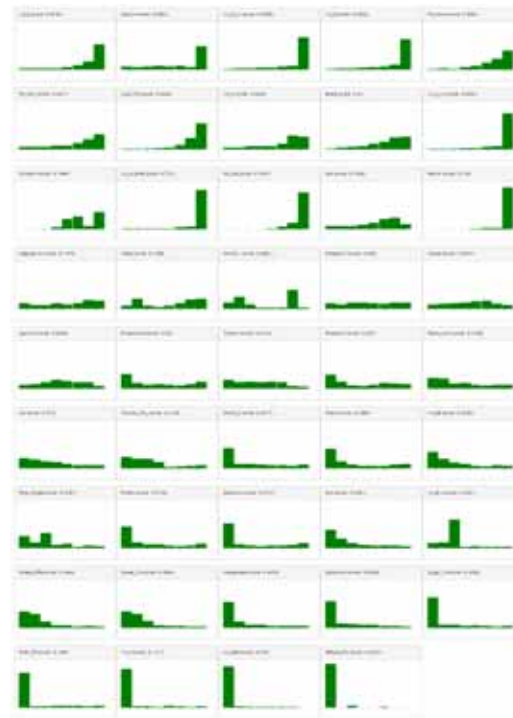
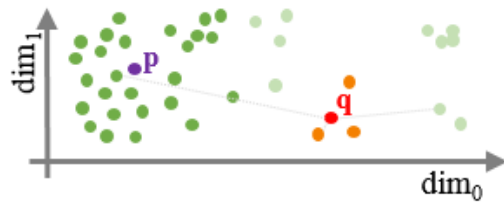
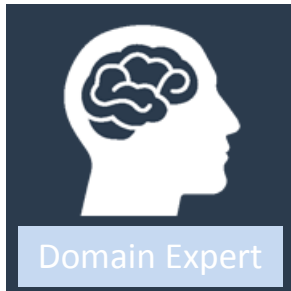
Non-Characteristic
Dimension



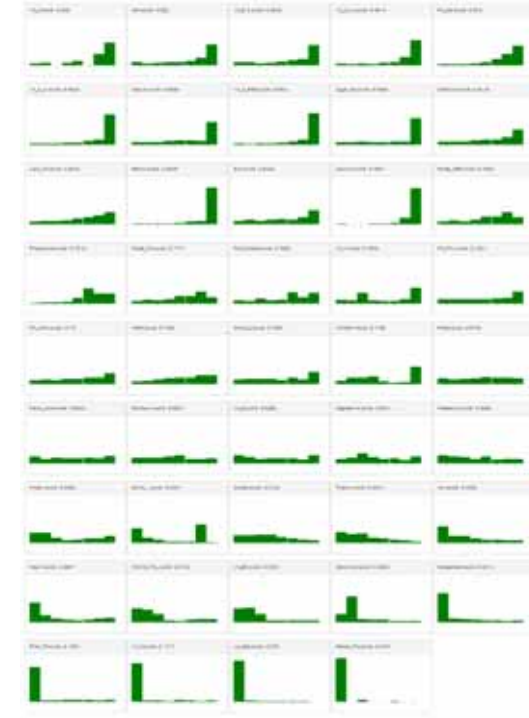
Characteristic
Dimension



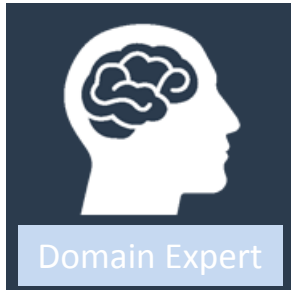
Data Distribution



query A



query B



(1) Determine Nearest Neighbors per Dimension

(2) Efficient Search Strategy



(3) Query-Based Interestingness for Dimensions

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

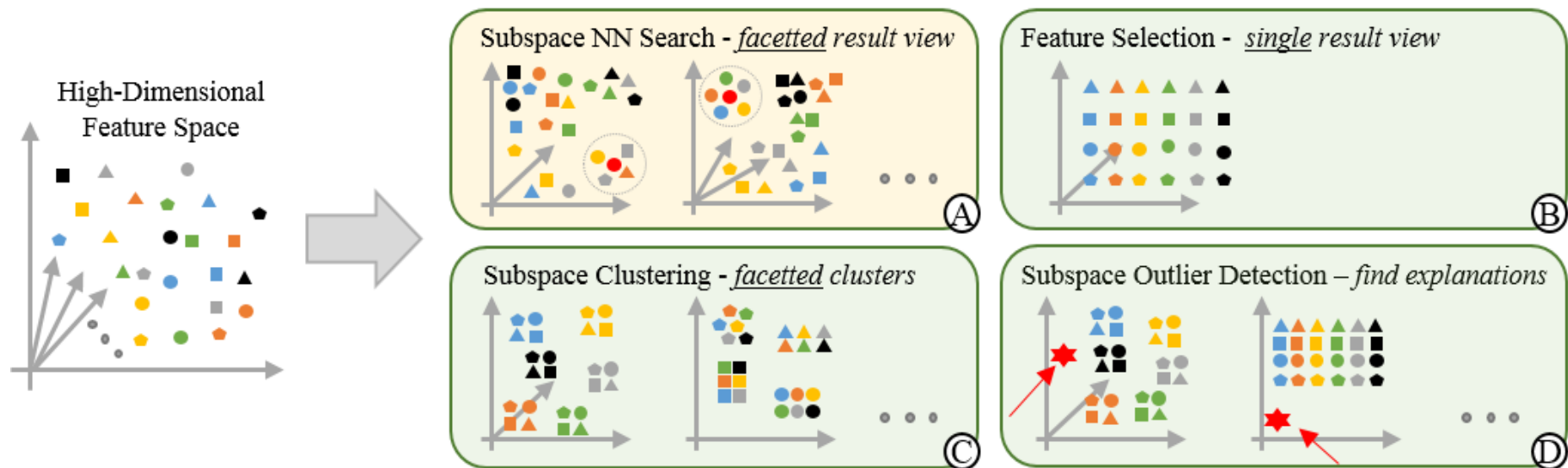


(5) Evaluation Methods and Development of Benchmark Datasets



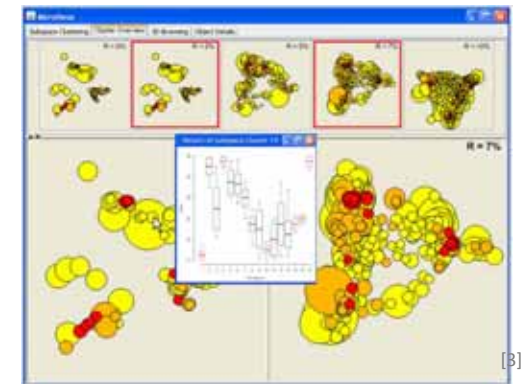
(6) Multi-input Subspace Nearest Neighbor Search

(7) Visualization and User Interaction

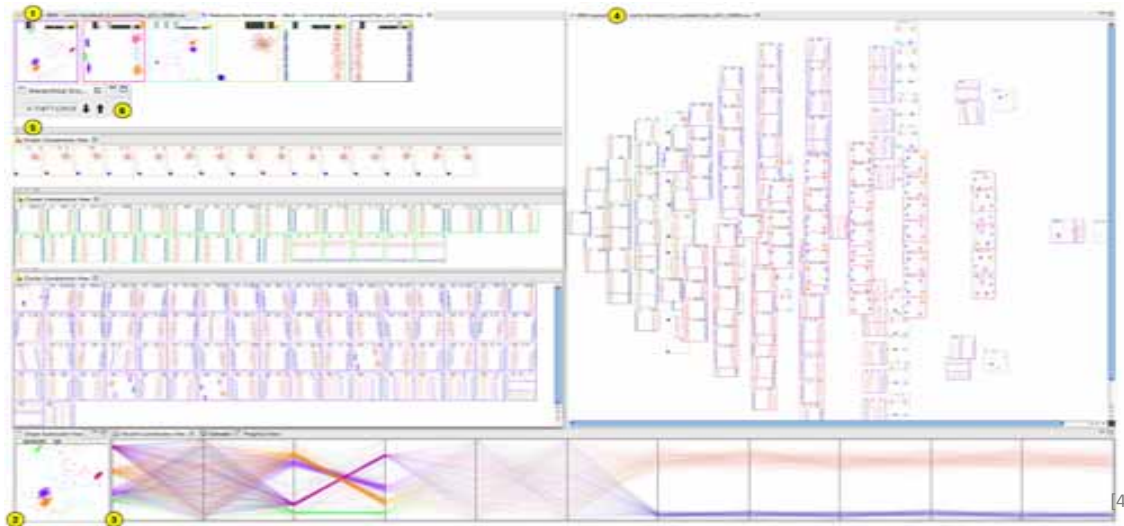


Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majnaric, L. & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Results Based on Subspace Clustering. In: Guo, Y., Friston, K., Aldo, F., Hill, S. & Peng, H. (eds.) Brain Informatics and Health, Lecture Notes in Artificial Intelligence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

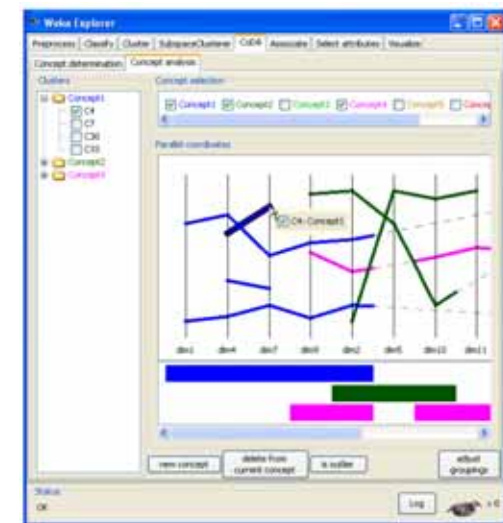
- VISA by Assent et al. (2007)
- CoDa by Günnemann et al (2010)
- Morpheus by Müller et al. (2008)
- Visual Analytics Framework by Tatu et al. (2012)



[3]



[4]



[2]

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

Is the parameter setting appropriate for the data?

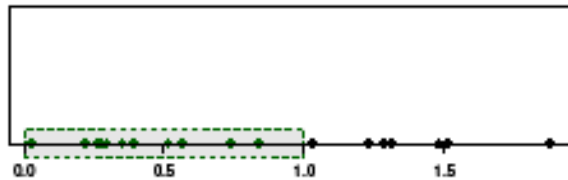
What happens if algorithms cannot scale with the #dimensions?

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

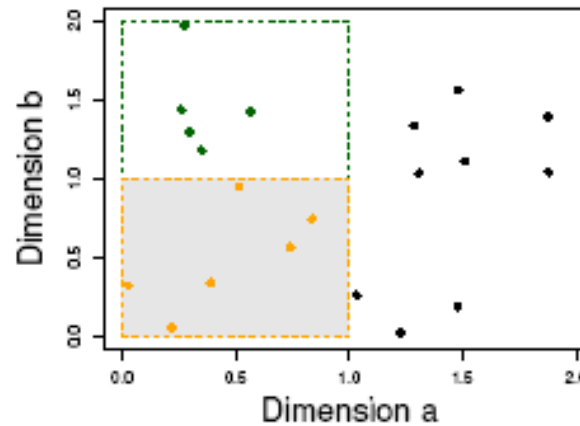
Domain Expert



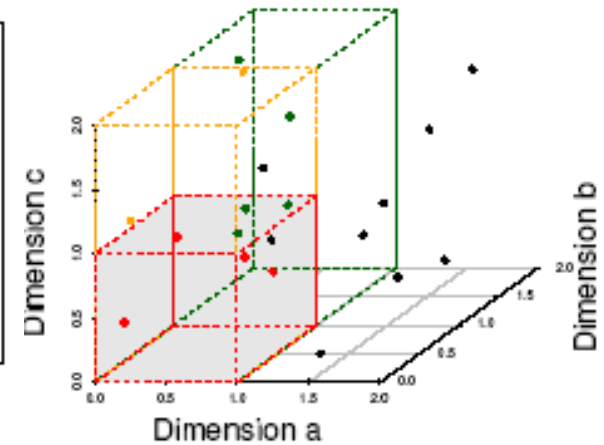
Hund, M., Boehm, D., Sturm, W., Sedlmair, M., Schreck, T., Ullrich, T., Keim, D. A., Majnarić, L. & Holzinger, A. 2016. Visual analytics for concept exploration in subspaces of patient groups: Making sense of complex datasets with the Doctor-in-the-loop. Brain Informatics, 3, (4), 233-247, doi:10.1007/s40708-016-0043-5.



(a) 11 Objects in One Unit Bin



(b) 6 Objects in One Unit Bin

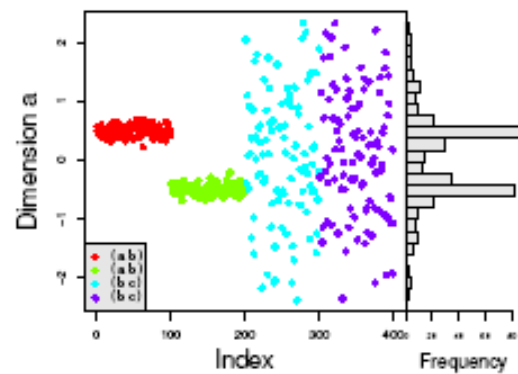
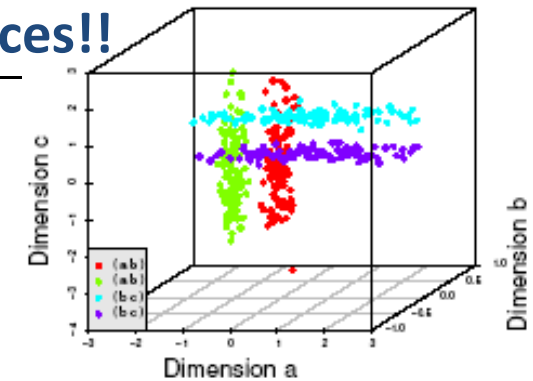


(c) 4 Objects in One Unit Bin

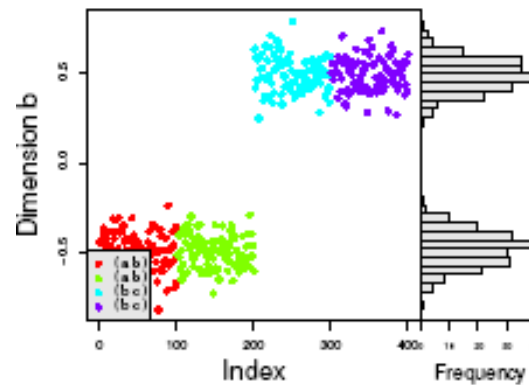
- Data in only one dimension is relatively packed
- Adding a dimension “stretch” the points across that dimension, making them further apart
- Adding more dimensions will make the points further apart—high dimensional data is extremely sparse
- Distance measure becomes meaningless—due to equidistance

- **Data set** - 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.

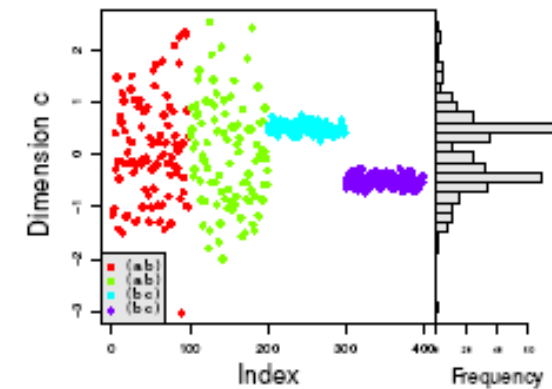
Parsons, L., Haque, E. & Liu, H. 2004. Subspace clustering for high dimensional data: a review. SIGKDD Explorations 6, (1), 90-105.



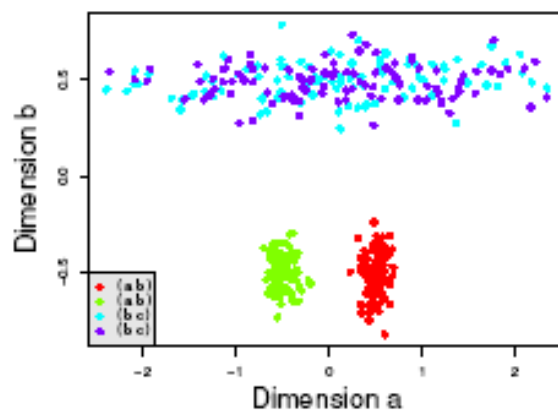
(a) Dimension *a*



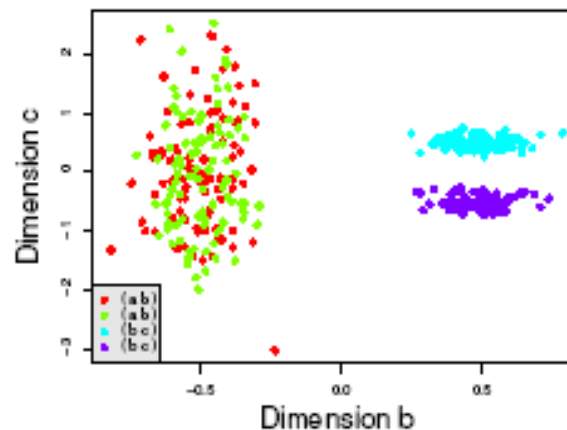
(b) Dimension *b*



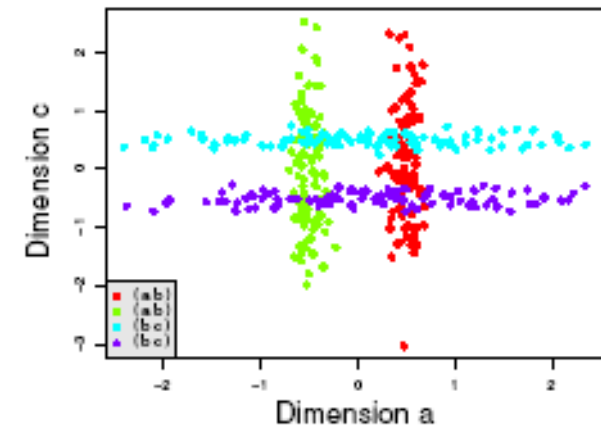
(c) Dimension *c*



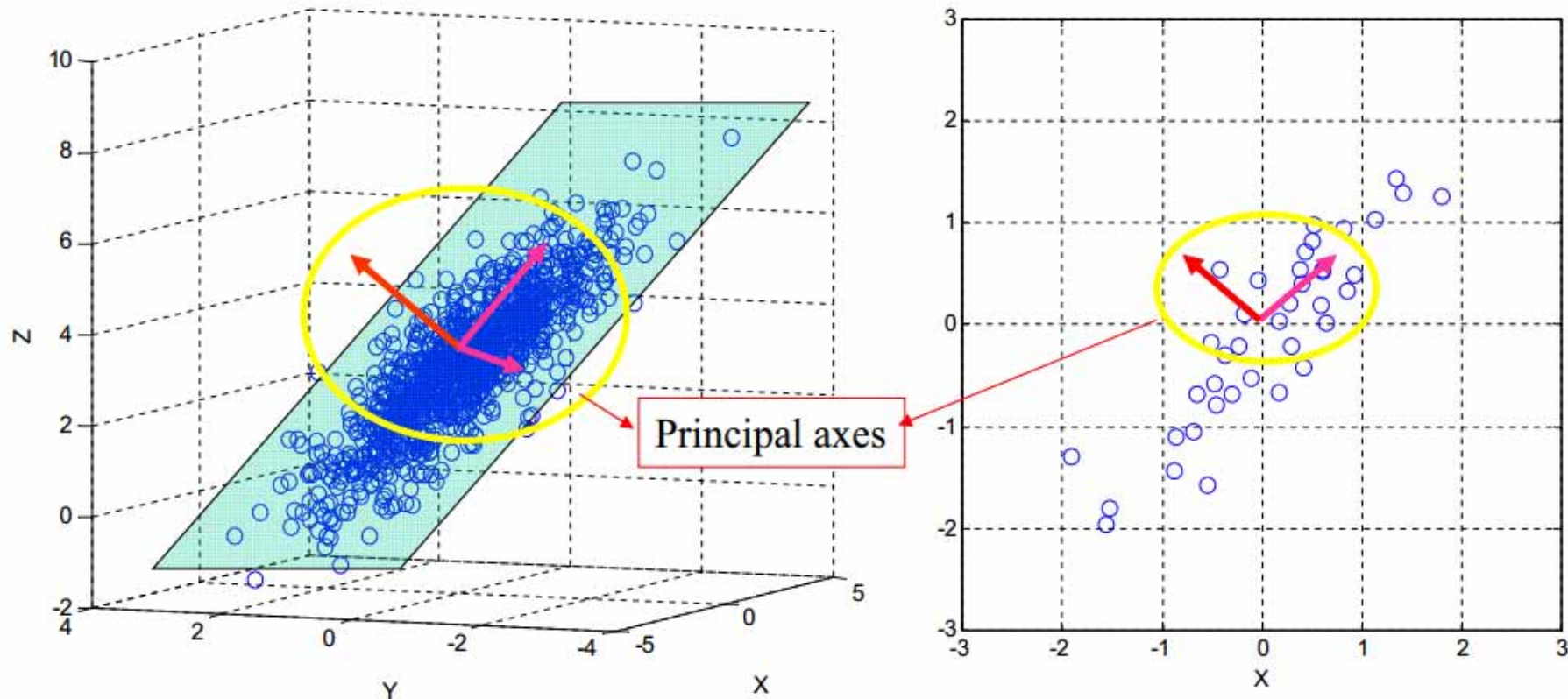
(a) Dims *a* & *b*



(b) Dims *b* & *c*



(c) Dims *a* & *c*



- We assume that
- 1) data sets concentrate to a low d-dim. linear subspace
- 2) axes of the subspaces are representations of the data
- 3) identifying the axes can be done by PCA

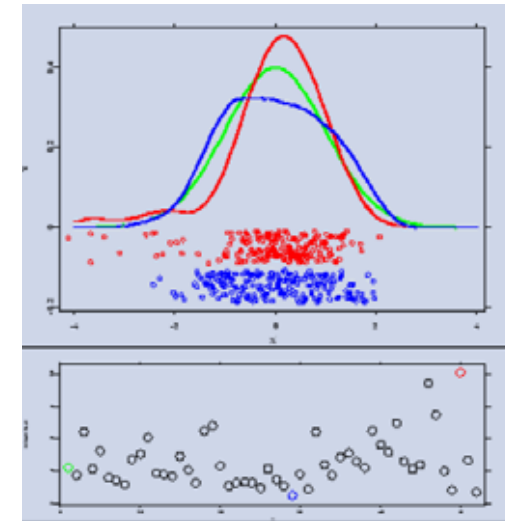
06 “What is interesting?” Projection Pursuit



- **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.

- Remember: Gaussian distribution maximizes the entropy!
- Now the objective is to minimize the entropy:
- $\min H(t)$ for $t = \omega^T x$
- (i.e. t is normalized)

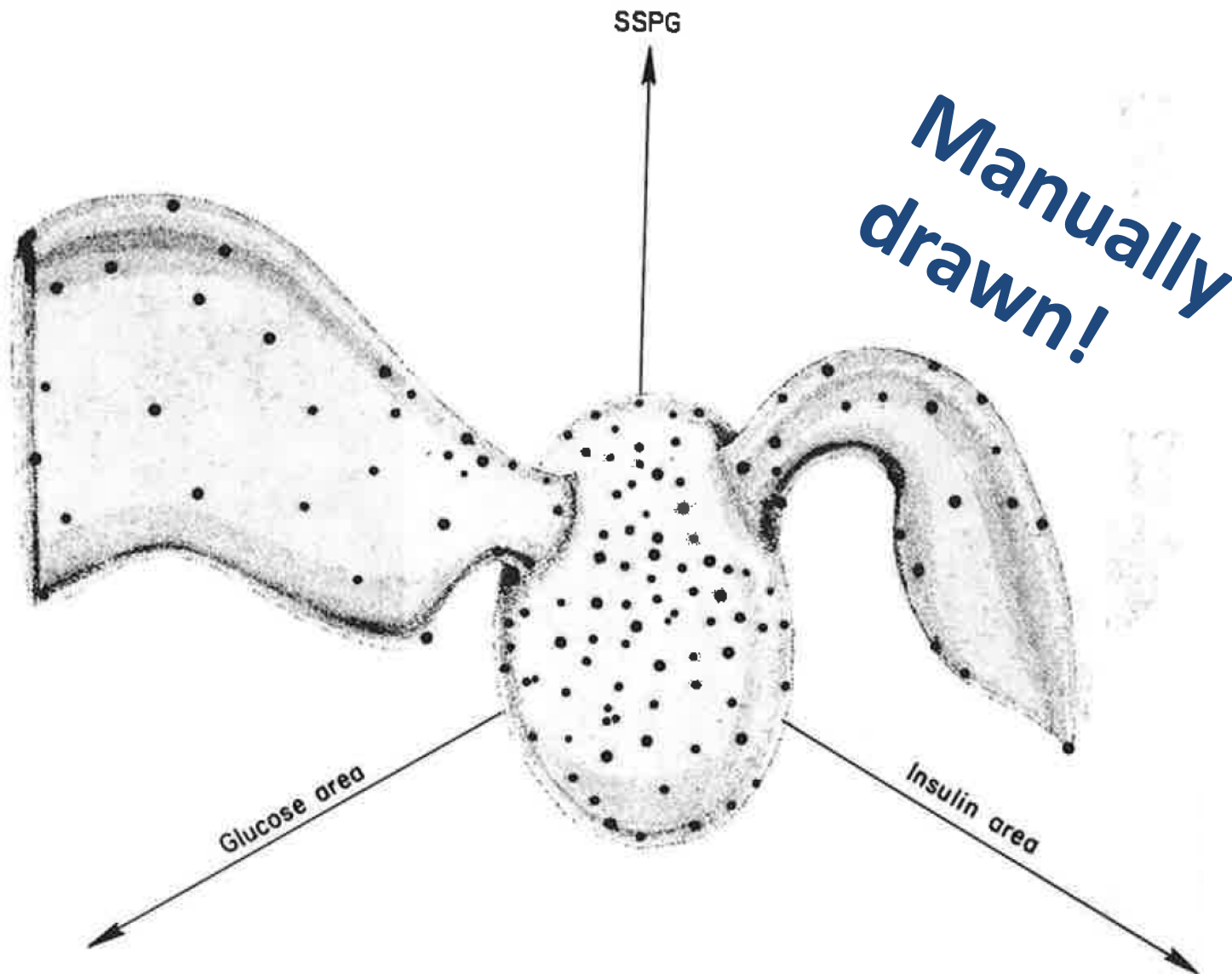


<http://fedc.wiwi.hu-berlin.de/xplore/tutorials/mvahtmlnode115.html>

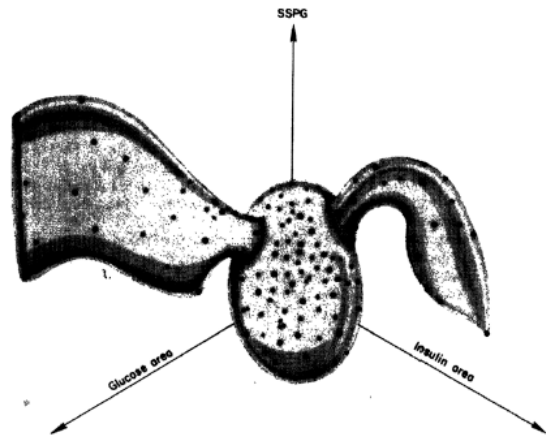
Friedman, J. H. & Tukey, J. W. 1974. A projection pursuit algorithm for exploratory data analysis. IEEE Transactions on Computers, 100, (9), 881-890.

- 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.

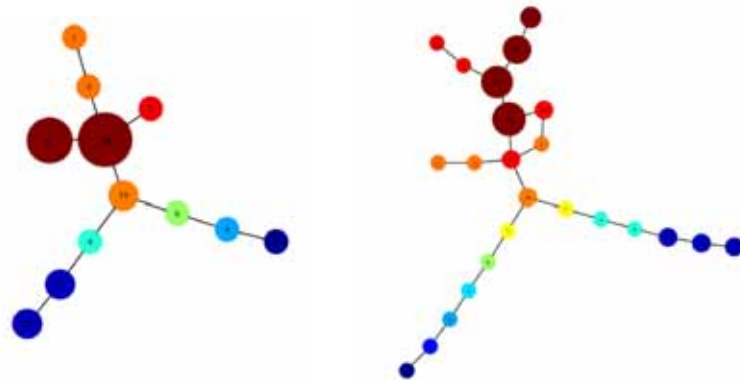


Given a point cloud data set X and a covering U
 \Rightarrow *simplicial complex*

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

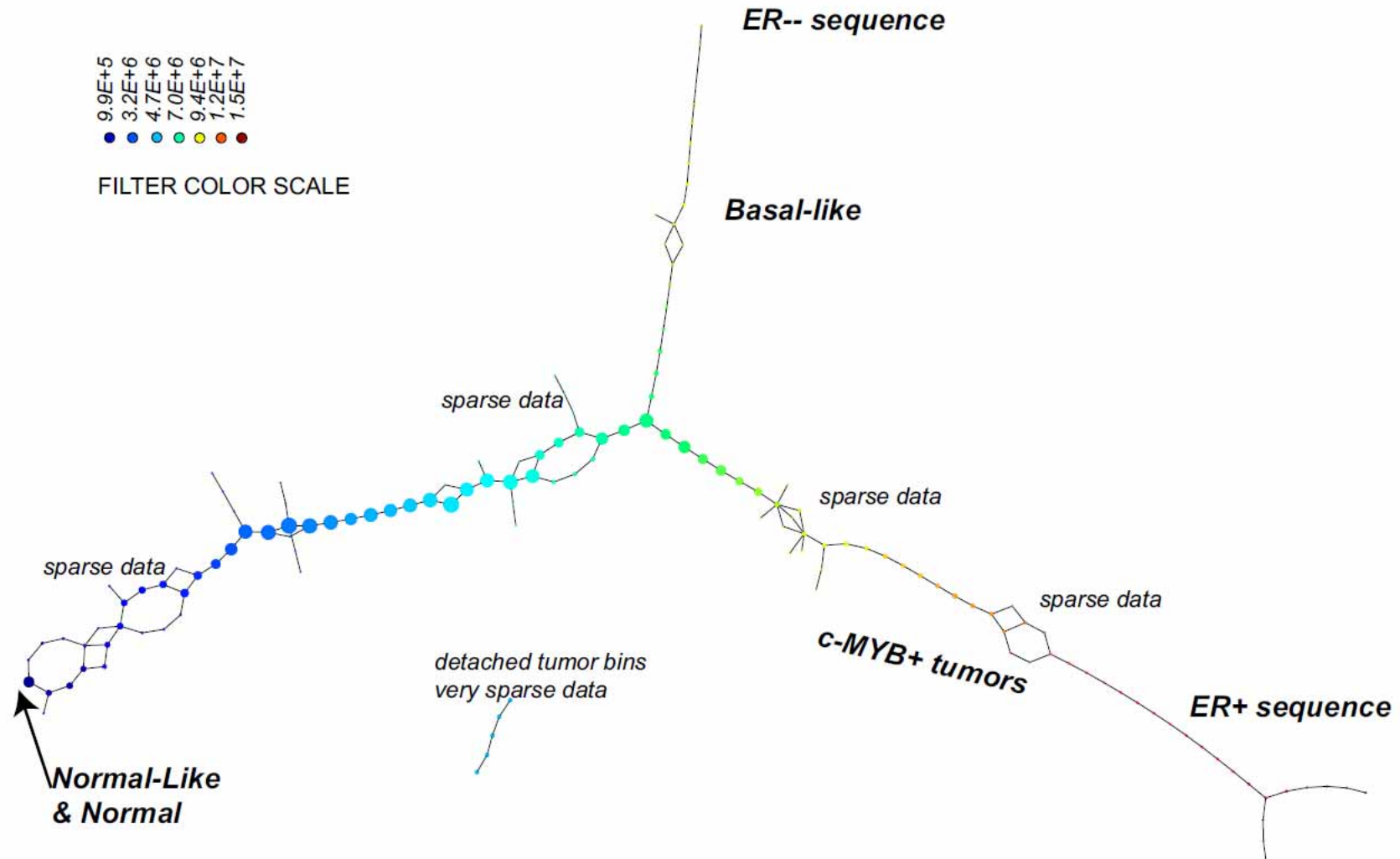
$$f: X \rightarrow Z$$

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



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

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.

Conclusion and Future Challenges

- Sometimes we have
 - A small number of data sets
 - Rare events – “little data”
 - NP-hard problems (e.g. k-Anonymization, Protein-Folding, Graph Coloring, Subspace Clustering, ...)
- Then we still need the “human-in-the-loop”



- 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: “little data”
- Integration in “real-world” (e.g. Hospital context)
- 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?

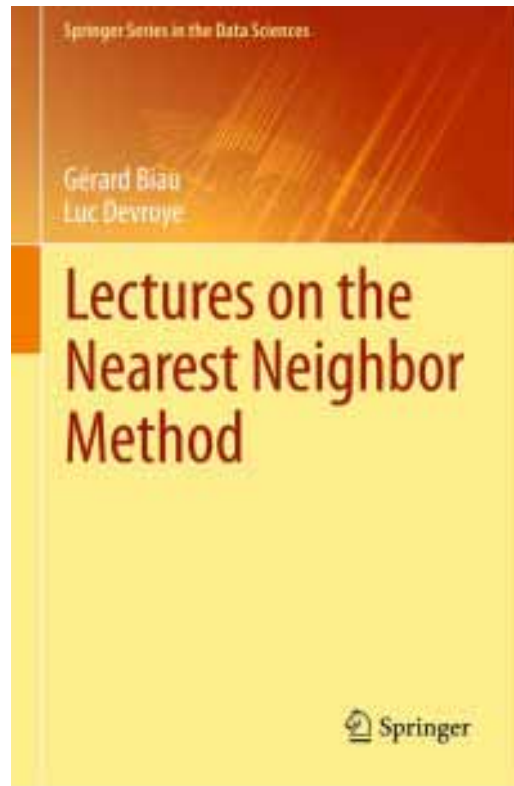


Thank you!

Questions

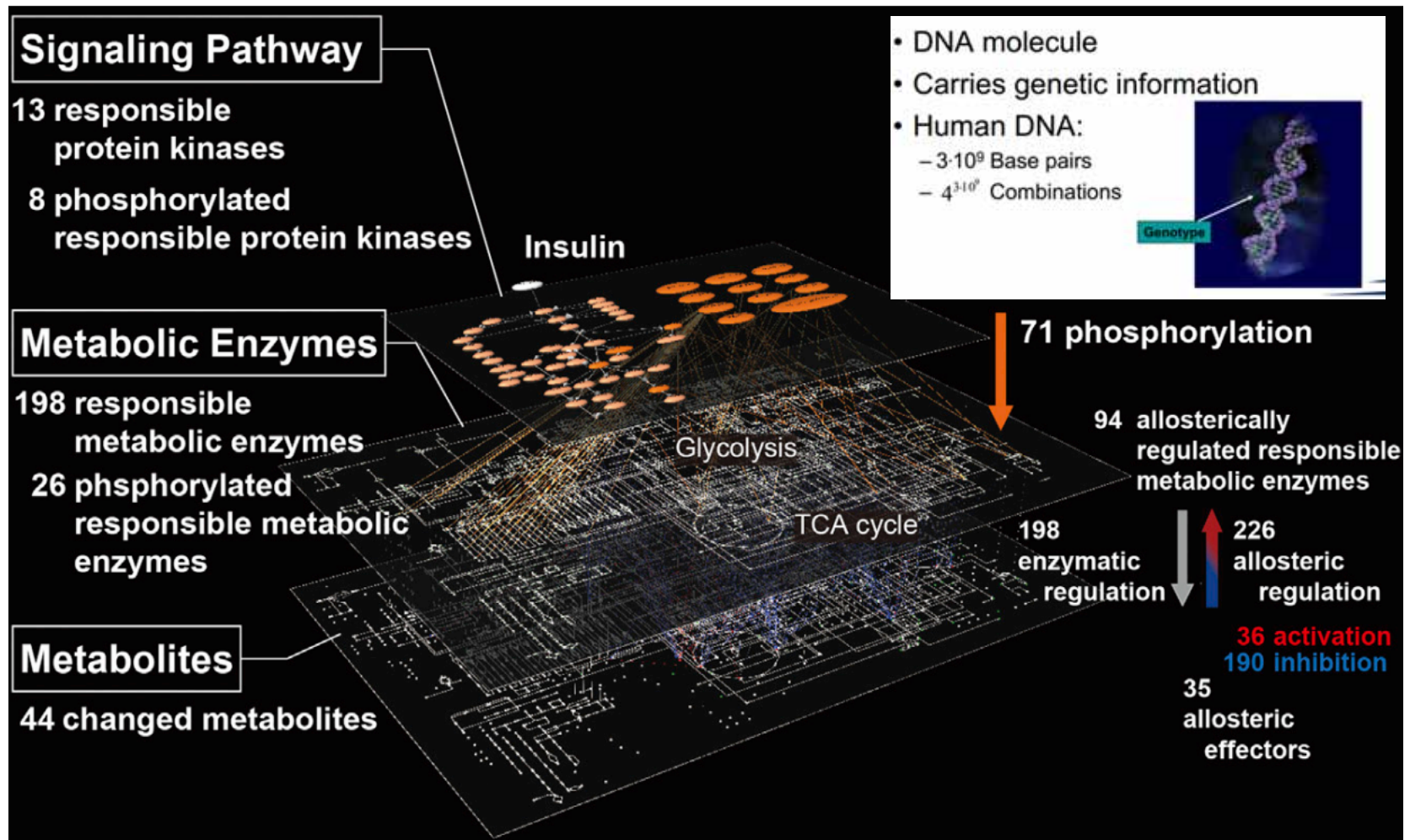
- Why would we wish at all to reduce the dimensionality of a data set?
- Why is feature selection so important? What is the difference between feature selection and feature extraction?
- What types of feature selection do you know?
- Can Neural Networks also be used to select features?
- Why do we need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

Appendix



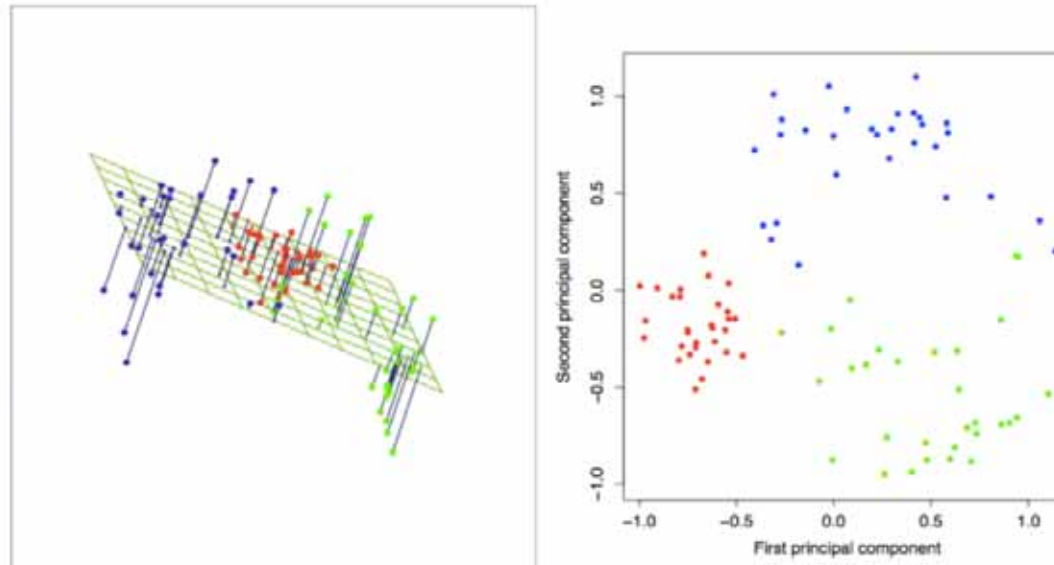
“Children learn effortlessly by example and exhibit a remarkable capacity of generalization. The field of machine learning, on the other hand, stumbles along clumsily in search of algorithms and methods, but nothing available today comes even close to an average two-year-old toddler ... “

Biau, G. & Devroye, L. 2016. Lectures on the nearest neighbor method, Springer, doi:10.1007/978-3-319-25388-6.

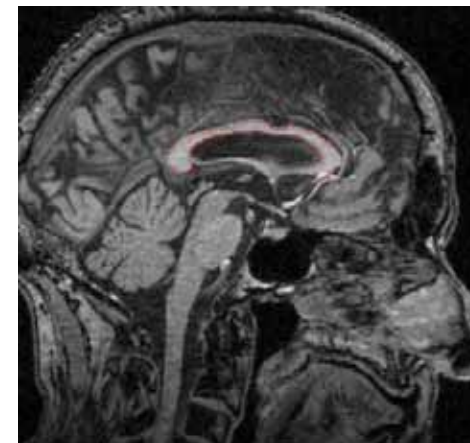


Yugi, K. et al. 2014. Reconstruction of Insulin Signal Flow from Phosphoproteome and Metabolome Data. Cell Reports, 8, (4), 1171-1183, doi:10.1016/j.celrep.2014.07.021.

- Linear methods (unsupervised):
 - PCA (Principal Component Analysis)
 - FA (Factor Analysis)
 - MDS (Multi-dimensional Scaling)
- Non-linear methods (unsupervised):
 - Isomap (Isometric feature mapping)
 - LLE (locally linear embedding)
 - Autoencoders
- Supervised methods:
 - LDA (Linear Discriminant Analysis)
- Subspace Clustering with a human-in-the-loop



- 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
$$S = \frac{1}{N} X^T X$$
- Compute k largest eigenvectors of S
- These eigenvectors are the k principal components



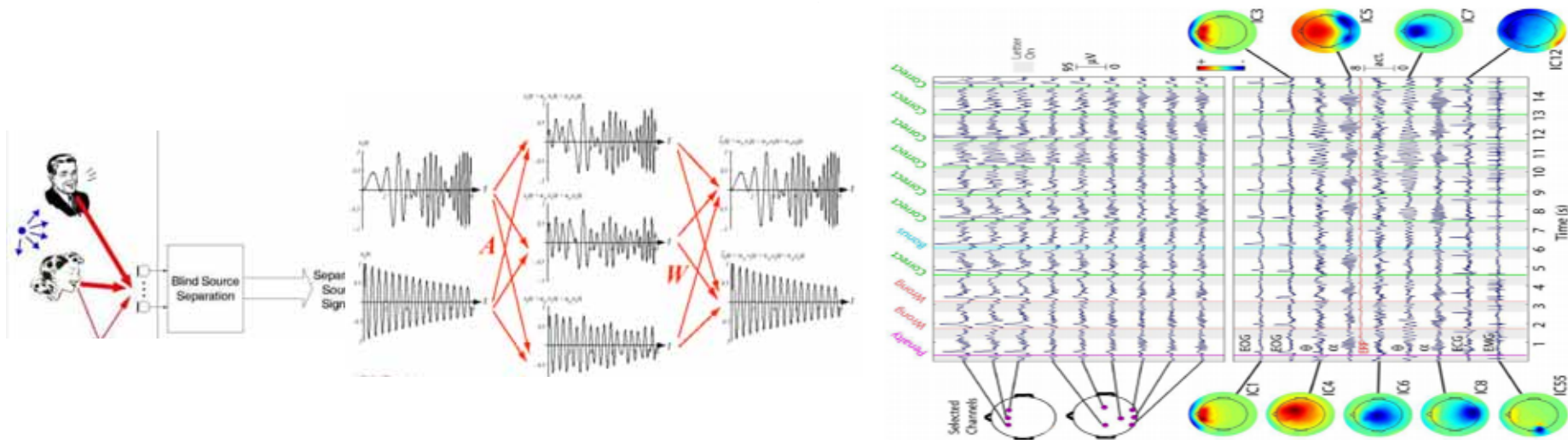
Hastie, T., Tibshirani, R. & Friedman, J. 2009. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second Edition, New York, Springer, doi:10.1007/978-0-387-84858-7

- 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

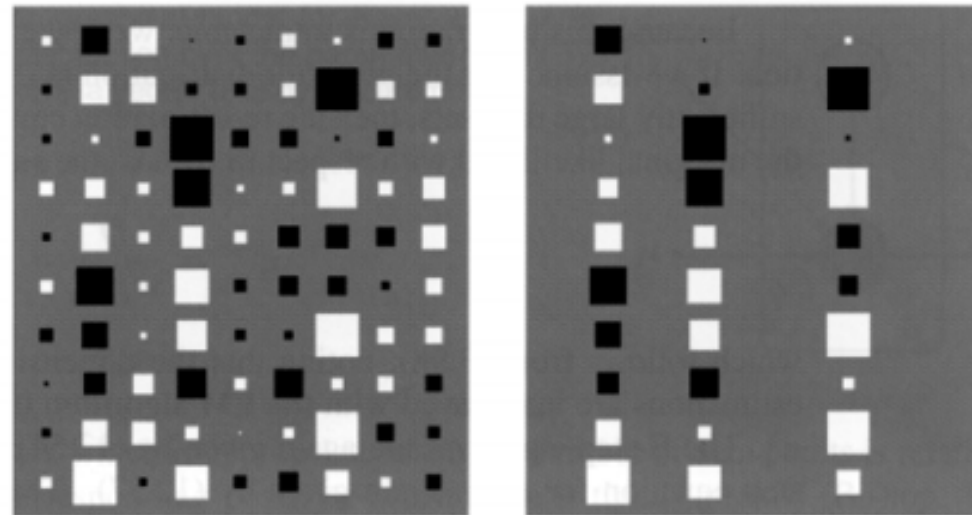


Holzinger, A., Scherer, R., Seeber, M., Wagner, J. & Müller-Putz, G. 2012. Computational Sensemaking on Examples of Knowledge Discovery from Neuroscience Data: Towards Enhancing Stroke Rehabilitation. In: Böhm, C., Khuri, S., Lhotská, L. & Renda, M. (eds.) Information Technology in Bio- and Medical Informatics,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

- FA describes variability of observations given unobserved **latent variables = factors**.
- Factors explain correlation between variables
- Similar to PCA, the difference is the conditional probability of the data (ψ = diagonal matrix):

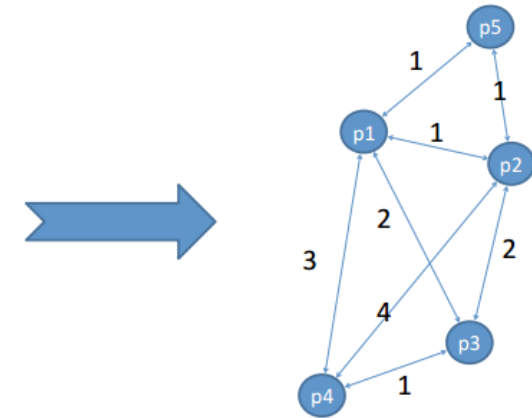
$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

Bishop, C. M. 2006. Pattern Recognition and Machine Learning, Heidelberg, Springer, Chapter 12.2.4



- 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



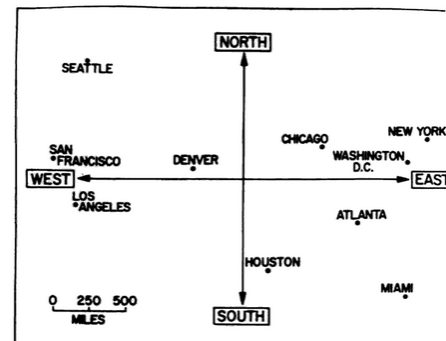
x_i Point in d dimensions

y_i Corresponding point in $r < d$ dimension

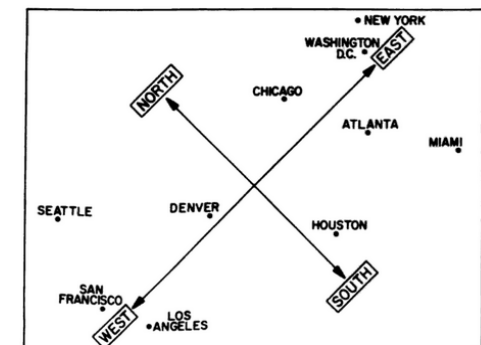
δ_{ij} Distance between x_i and x_j

d_{ij} Distance between y_i and y_j

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



CITIES	ATLA	CHIC	DENV	HOUS	L.A.	MIAMI	N.Y.	S.F.	SEAT	WASH D.C.
ATLANTA		587	1212	701	1936	604	748	2139	2182	543
CHICAGO	587		920	940	1745	1188	713	1858	1737	597
DENVER	1212	920		879	831	1726	1631	949	1021	1494
HOUSTON	701	940	879		1374	968	1420	1645	1891	1220
LOS ANGELES	1936	1745	831	1374		2339	2451	347	959	2300
MIAMI	604	1188	1726	968	2339		1092	2594	2734	923
NEW YORK	748	713	1631	1420	2451	1092		2571	2408	205
SAN FRANCISCO	2139	1858	949	1645	347	2594	2571		678	2442
SEATTLE	2182	1737	1021	1891	959	2734	2408	678		2329
WASHINGTONDC	543	597	1494	1220	2300	923	205	2442	2329	



Kruskal, J. B. 1964. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. Psychometrika, 29, (1), 1-27.

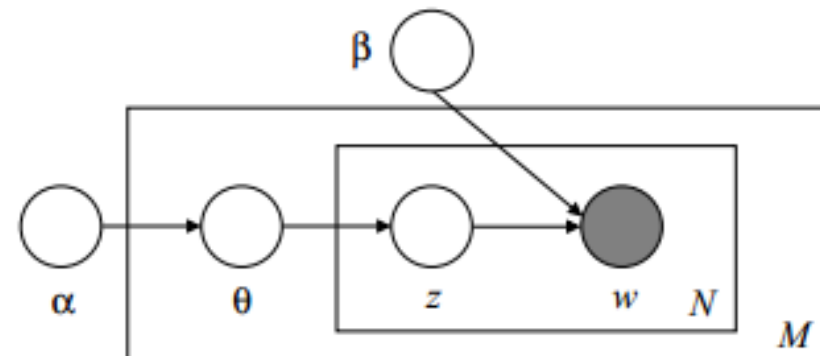
"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

- (1) initialize $\phi_{ni}^0 := 1/k$ for all i and n
- (2) initialize $\gamma_i := \alpha_i + N/k$ for all i
- (3) **repeat**
- (4) **for** $n = 1$ **to** N
- (5) **for** $i = 1$ **to** k
- (6) $\phi_{ni}^{t+1} := \beta_{ni} \exp(\Psi(\gamma_i))$
- (7) normalize ϕ_{ni}^{t+1} to sum to 1.
- (8) $\gamma^{t+1} := \alpha + \sum_{n=1}^N \phi_n^{t+1}$
- (9) **until** convergence

$$p(D|\alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left(\prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d\theta_d$$

Blei, D. M., Ng, A. Y. & Jordan, M. I. 2003.
Latent Dirichlet allocation. Journal of
Machine Learning Research, 3, (4-5), 993-
1022.



A Global Geometric Framework for Nonlinear Dimensionality Reduction

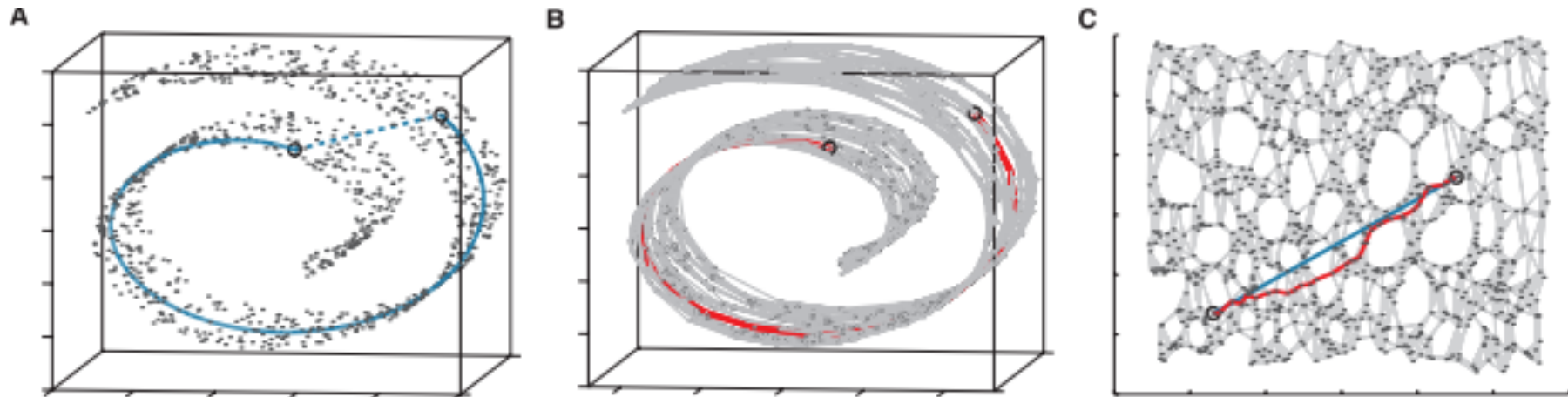
Joshua B. Tenenbaum,^{1*} Vin de Silva,² John C. Langford³

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 high-dimensional observations. The human brain confronts the same problem in everyday perception, extracting from its high-dimensional sensory inputs—30,000 auditory nerve fibers or 10^6 optic nerve fibers—a manageably small number of perceptually relevant features. Here we describe an approach to solving dimensionality reduction problems that uses easily measured local metric information to learn the underlying global geometry of a data set. Unlike classical techniques such as principal component analysis (PCA) and multidimensional scaling (MDS), our approach is capable of discovering the nonlinear degrees of freedom that underlie complex natural observations, such as human handwriting or images of a face under different viewing conditions. In contrast to previous algorithms for nonlinear dimensionality reduction, ours efficiently computes a globally optimal solution, and, for an important class of data manifolds, is guaranteed to converge asymptotically to the true structure.

Goal: Find projection onto *nonlinear* manifold

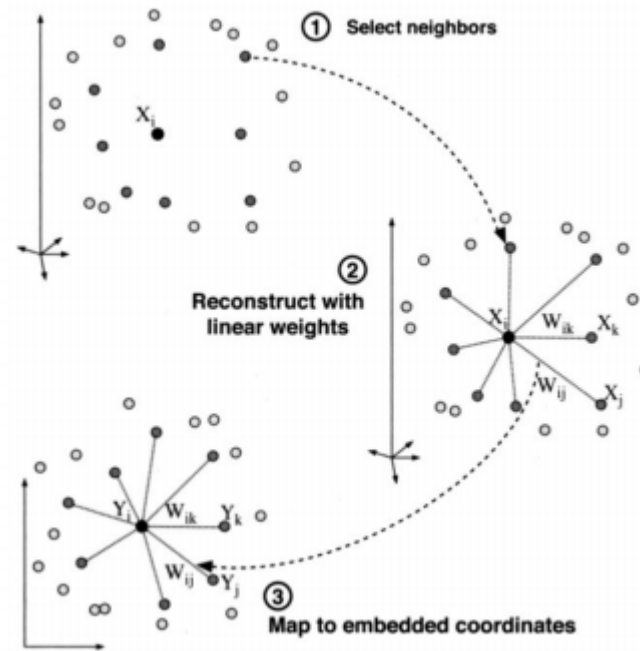
1. Construct neighborhood graph G :
For all x_i, x_j
If $\text{distance}(x_i, x_j) < \epsilon$
Then add edge (x_i, x_j) to G
2. Compute shortest distances along graph $\delta_G(x_i, x_j)$
(e.g., by Floyd's algorithm)
3. Apply multidimensional scaling to $\delta_G(x_i, x_j)$

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



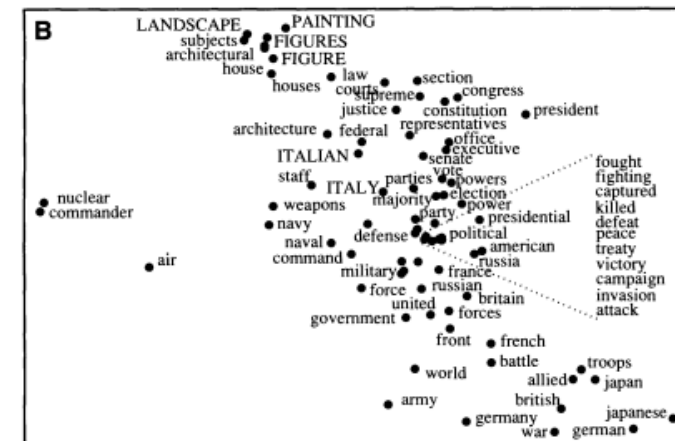
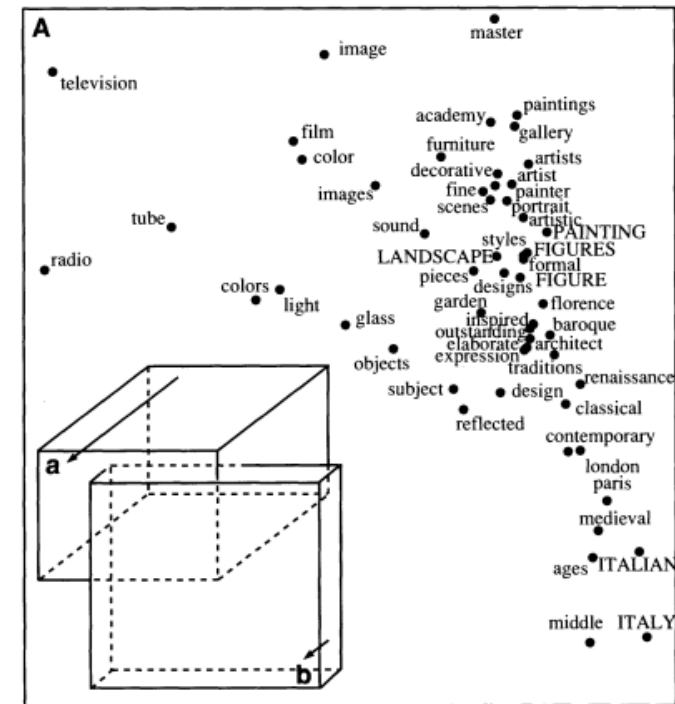
Tenenbaum, J. B., De Silva, V. & Langford, J. C. 2000. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290, (5500), 2319-2323, doi:10.1126/science.290.5500.2319.

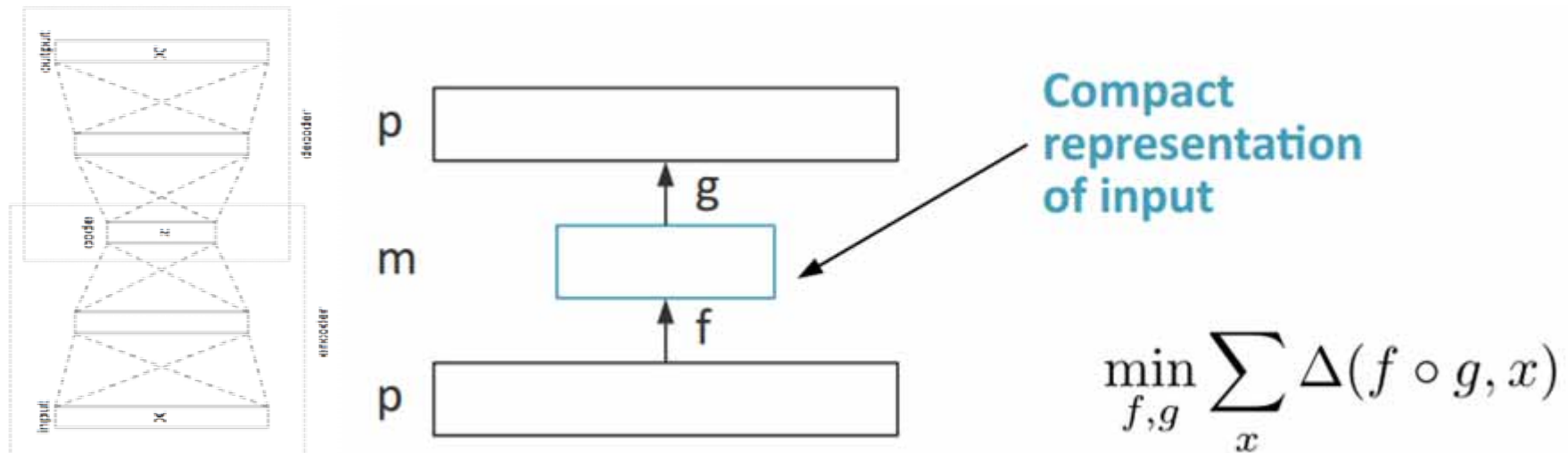
Roweis, S. T. & Saul, L. K. 2000. Nonlinear dimensionality reduction by locally linear embedding. Science, 290, (5500), 2323-2326, doi:10.1126/science.290.5500.2323.



$$\epsilon(W) = \sum_i \left| \vec{X}_i - \sum_j W_{ij} \vec{X}_j \right|^2$$

$$\Phi(Y) = \sum_i \left| \vec{Y}_i - \sum_j W_{ij} \vec{Y}_j \right|^2$$





- 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.

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

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

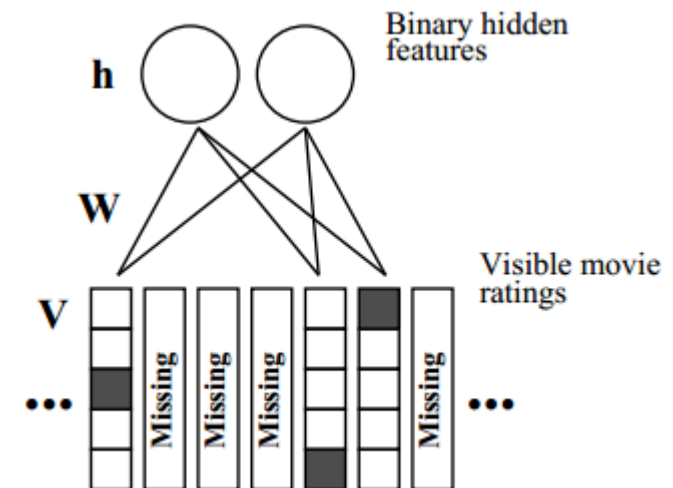
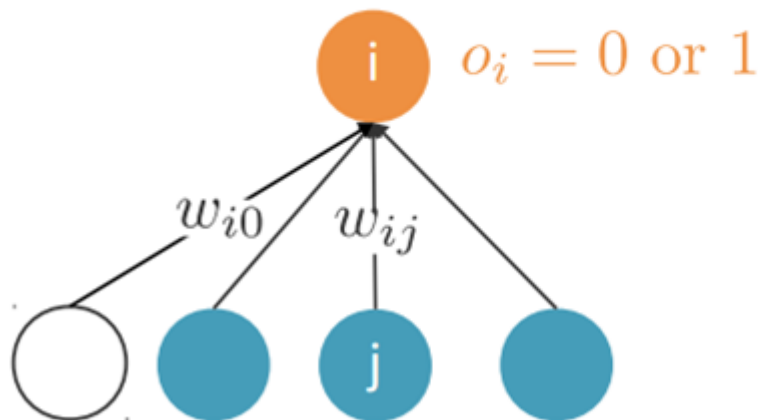
- **Linear autoencoders:** Baldi, P. & Hornik, K. 1989. Neural networks and principal component analysis: Learning from examples without local minima. Neural networks, 2, (1), 53-58.

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

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

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

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



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.

- 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)