

185.A83 Machine Learning for Health Informatics

2016S, VU, 2.0 h, 3.0 ECTS

Week 17 - 26.04.2016 17:00-20:00

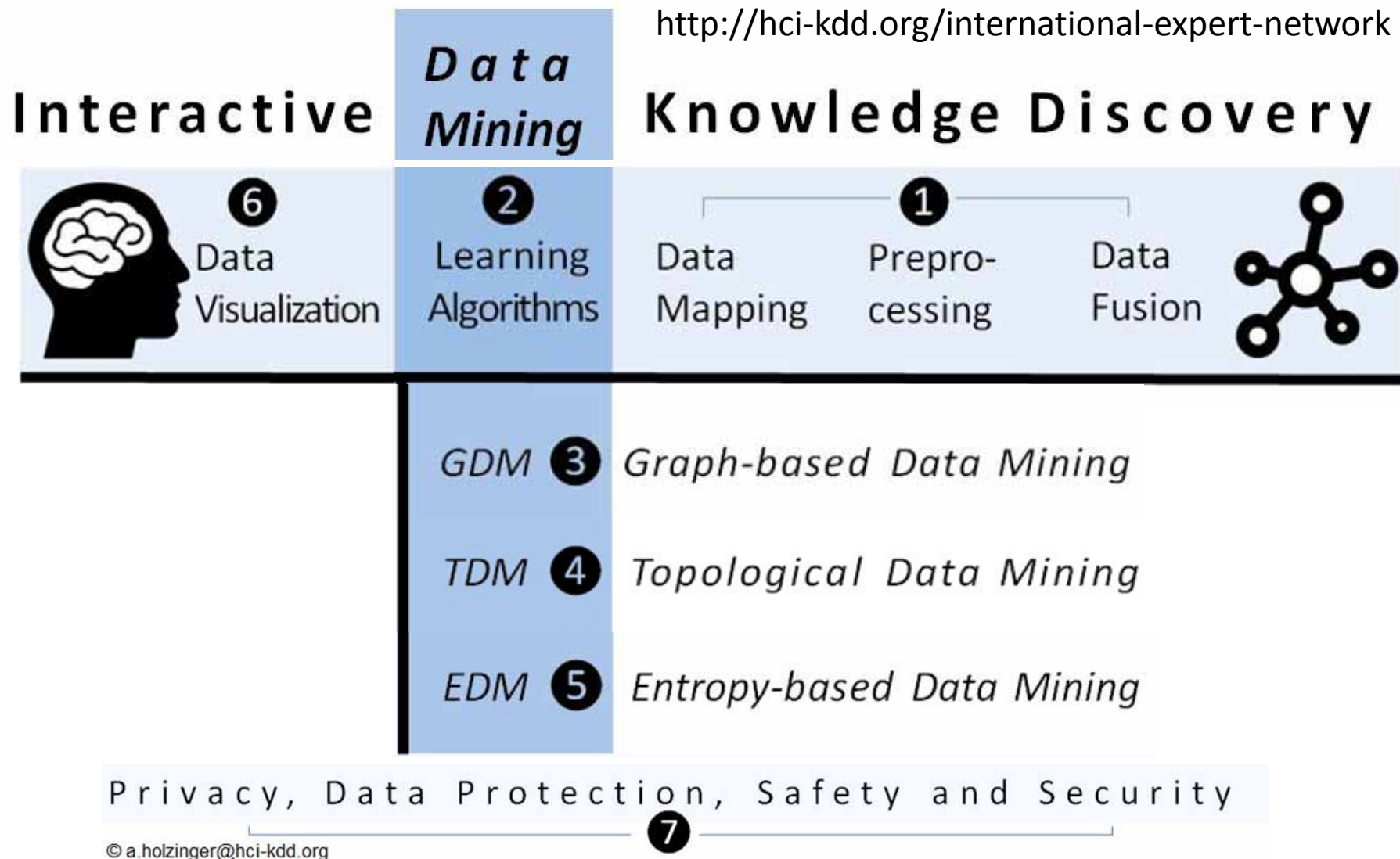
# Dimensionality Reduction and Subspace Clustering: Example for the Expert-in-the-Loop

a.holzinger@hci-kdd.org

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

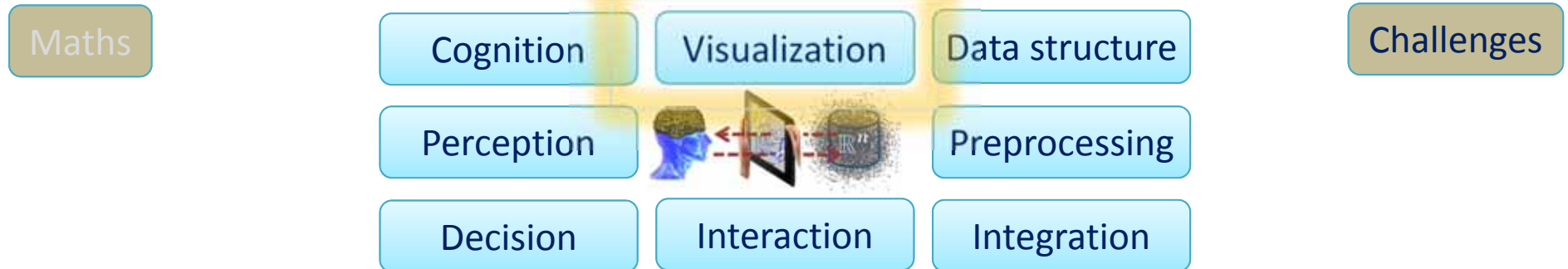


<http://hci-kdd.org/international-expert-network>

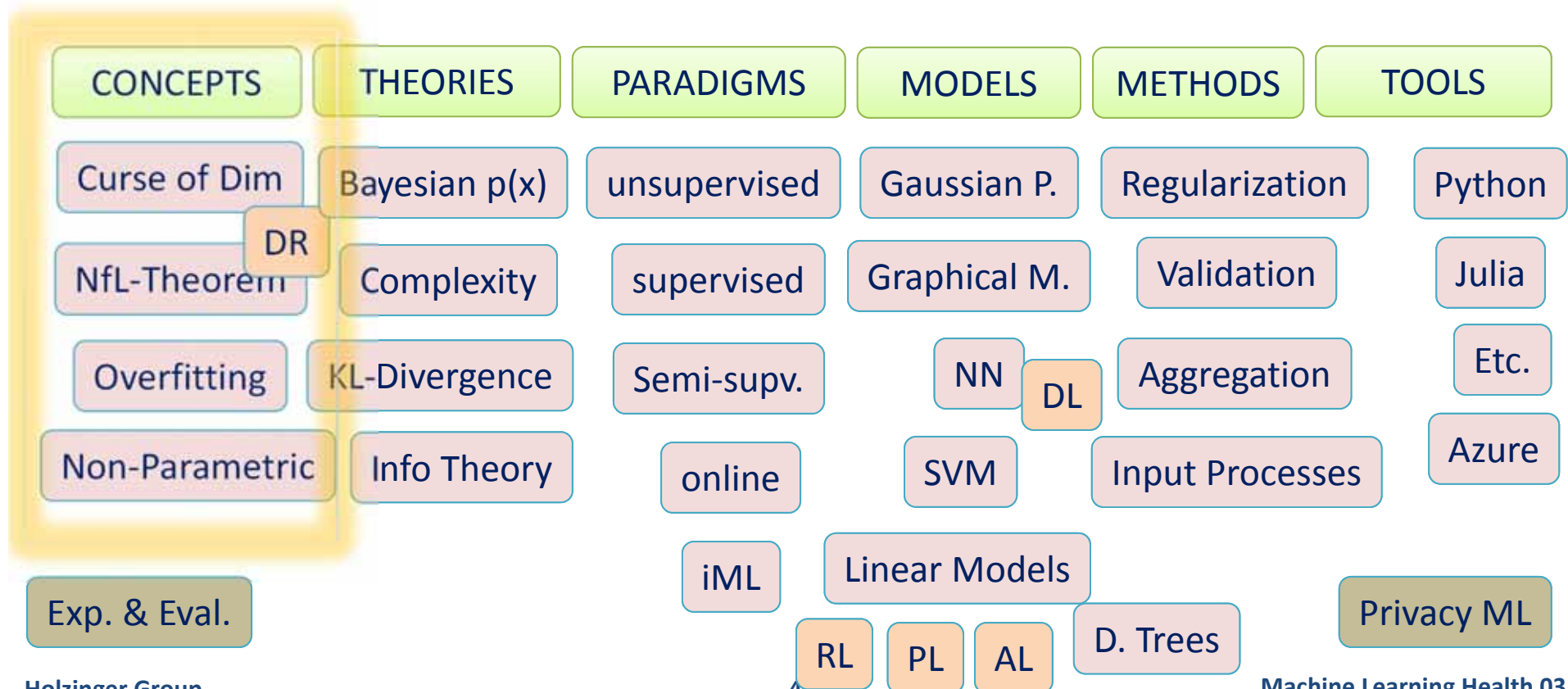


Holzinger, A. 2014. Trends in Interactive Knowledge Discovery for Personalized Medicine: **Cognitive Science meets Machine Learning**. IEEE Intelligent Informatics Bulletin, 15, (1), 6-14.

- 1) Classification vs Clustering
- 2) Feature spaces, feature engineering
  - Feature selection, feature extraction
- 3) The curse of dimensionality
- 4) Dimensionality reduction
  - PCA, ICA, FA, MDS, LDA – Isomap, LLE, Autoencoder
- 5) Subspace clustering and analysis
- 6) Projection Pursuit: “What is interesting?”



Always with a focus/application in health informatics



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

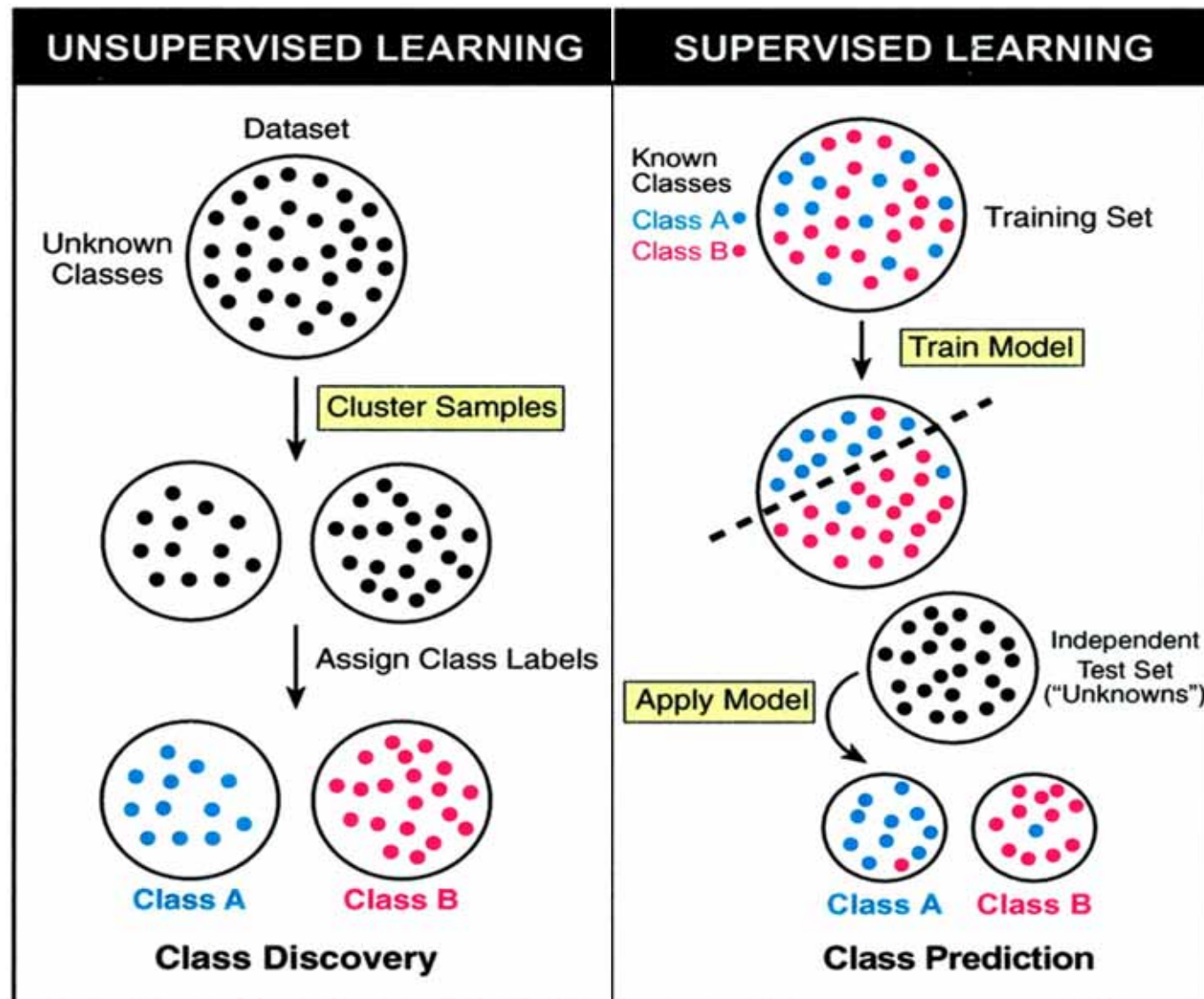
# 1) Classification vs. Clustering



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

- **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;



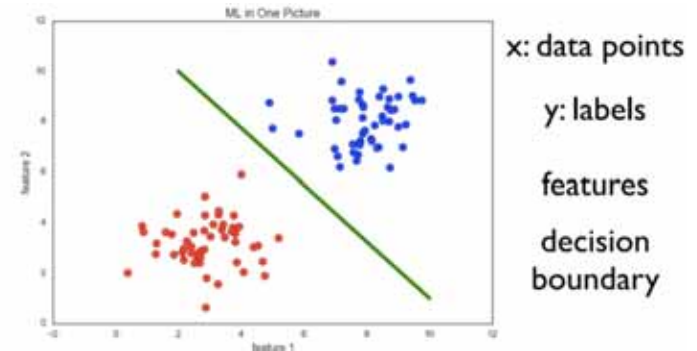


Ramaswamy, S. & Golub, T. R. (2002) DNA Microarrays in Clinical Oncology. *Journal of Clinical Oncology*, 20, 7, 1932-1941.

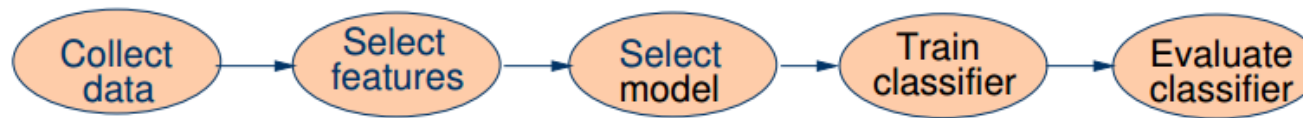


$x$  -- set of pixel intensities

$C_1$ : Cancer present  
 $C_2$ : Cancer absent



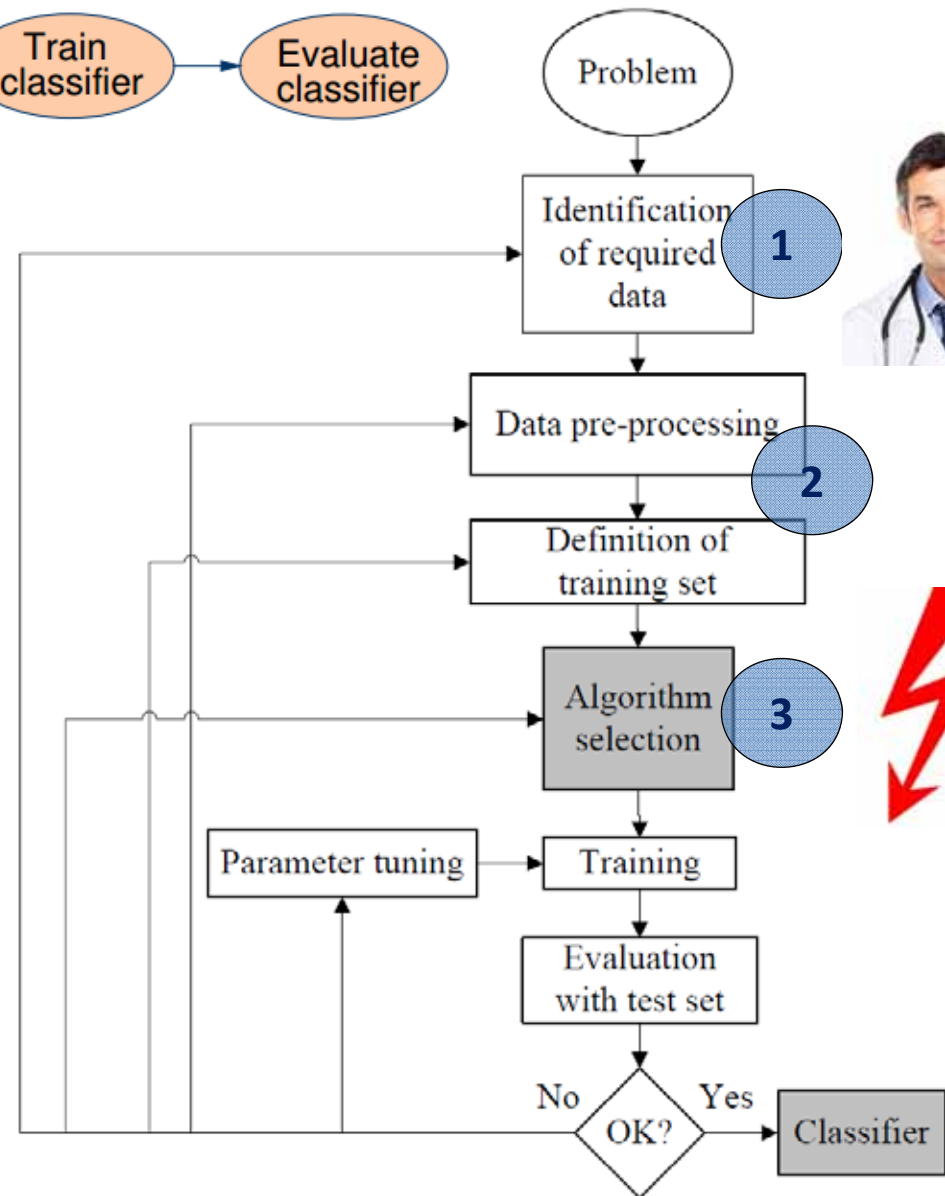
- Typical questions include:
  - Is this protein functioning as an enzyme?
  - Does this gene sequence contain a splice site?
  - Is this melanoma malign?
- Given object  $x$  – predict the class label  $y$ 
  - If  $y \in \{0,1\} \rightarrow$  binary classification problem
  - If  $y \in \{1, \dots, n\}$  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  $x_2 \Rightarrow y_1 = y_2$

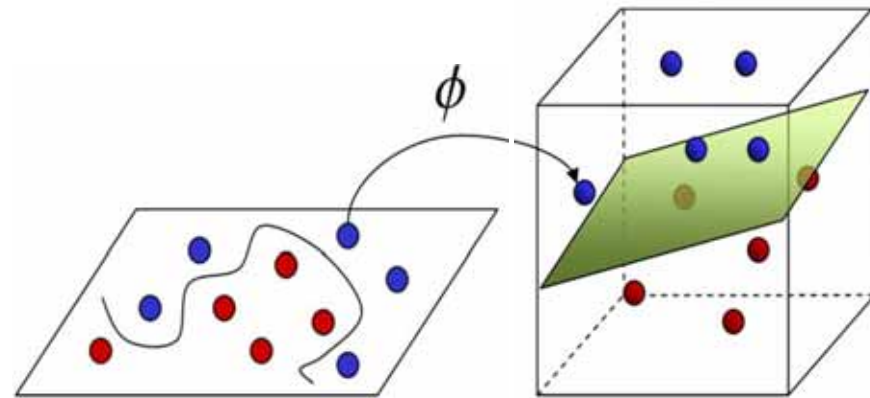
$$x_j = \operatorname{argmin}_{x \in D} ||x - x_i||^2 \Rightarrow y_i = y_j$$

- SVM – a plane/hyperplane separates two classes of data – very versatile for classification and clustering – also via the Kernel trick in high-dimensions

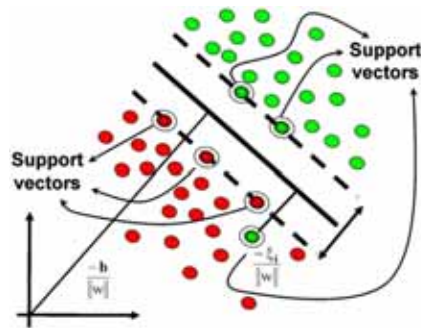
```
1: Input:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n), C, \epsilon$ 
2:  $S_i \leftarrow \emptyset$  for all  $i = 1, \dots, n$ 
3: repeat
4:   for  $i = 1, \dots, n$  do
5:      $H(\mathbf{y}) \equiv \Delta(\mathbf{y}_i, \mathbf{y}) + \mathbf{w}^T \Psi(\mathbf{x}_i, \mathbf{y}) - \mathbf{w}^T \Psi(\mathbf{x}_i, y_i)$ 
6:     compute  $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} H(\mathbf{y})$ 
7:     compute  $\xi_i = \max\{0, \max_{\mathbf{y} \in S_i} H(\mathbf{y})\}$ 
8:     if  $H(\hat{\mathbf{y}}) > \xi_i + \epsilon$  then
9:        $S_i \leftarrow S_i \cup \{\hat{\mathbf{y}}\}$ 
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)

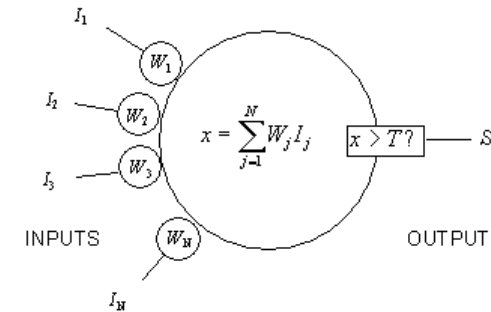


- = classification method for both linear and nonlinear data;
- Within the new dimension, it searches for the linear optimal separating **hyperplane** (i.e., “decision boundary”);
- By nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated with a hyperplane;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



## ■ SVM

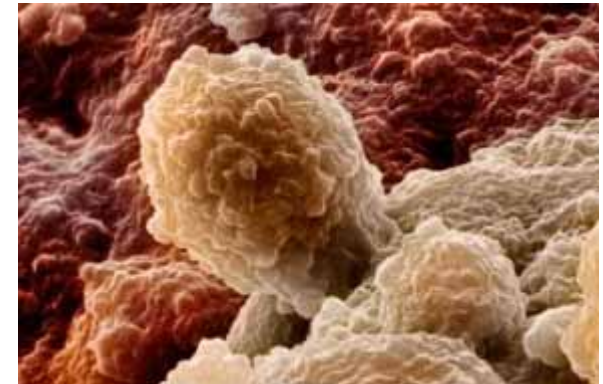
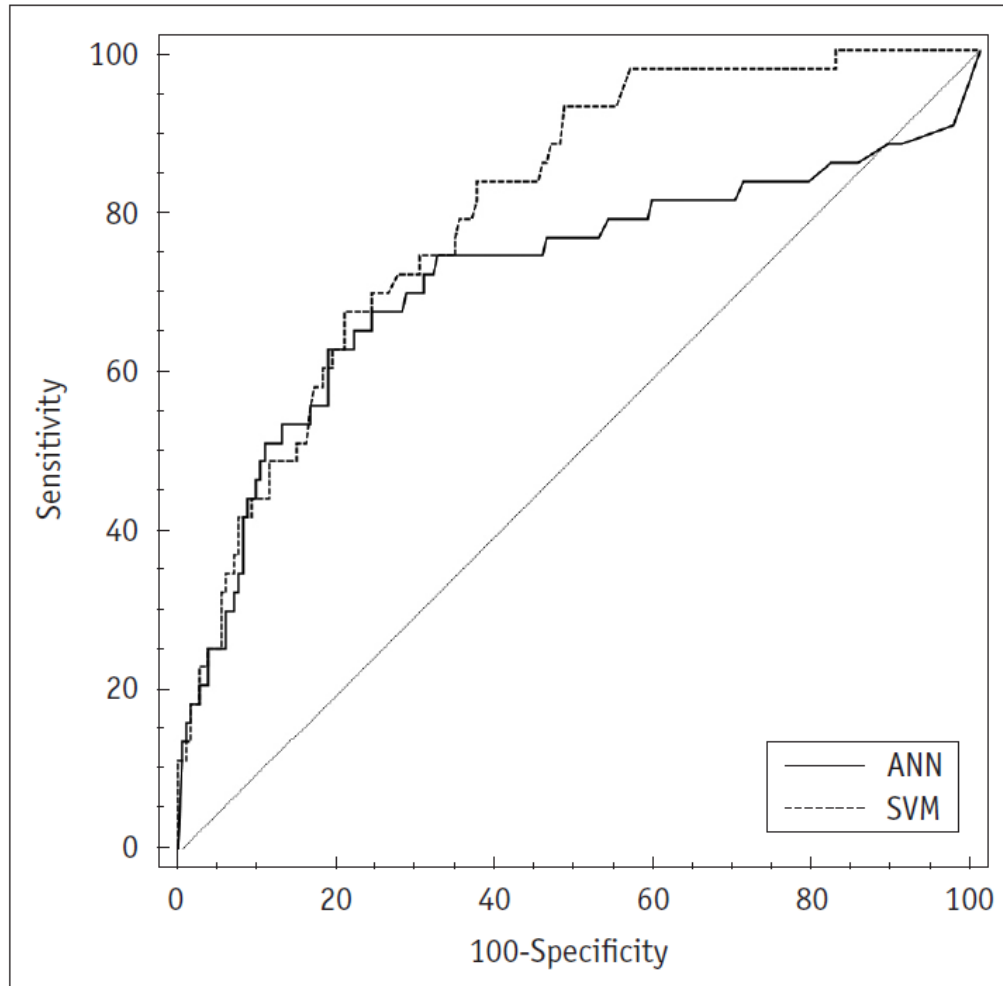
- Deterministic algorithm
- Nice generalization properties
- Hard to learn – learned in batch mode using quadratic programming techniques
- Using kernels can learn very complex functions



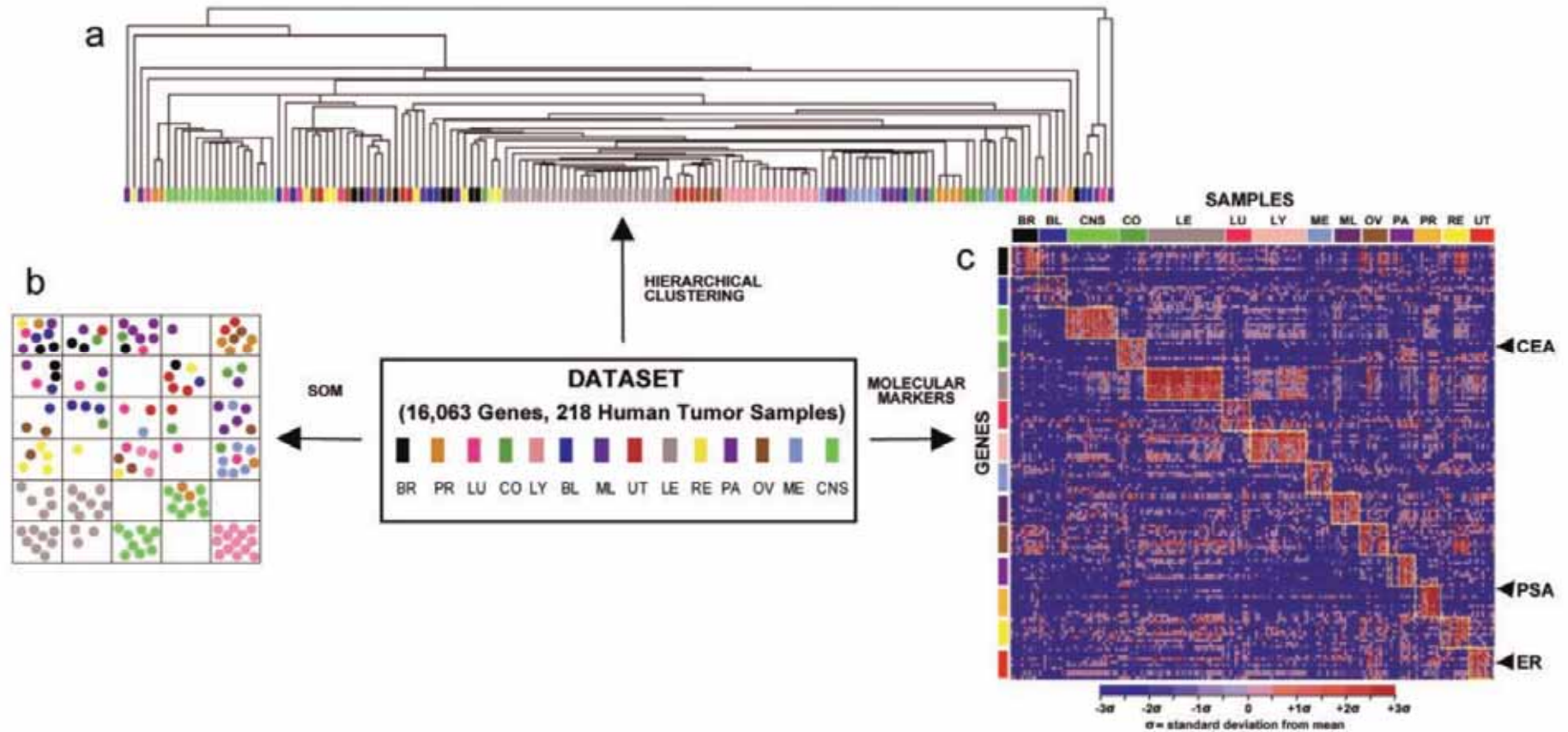
## ■ ANN

- 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)



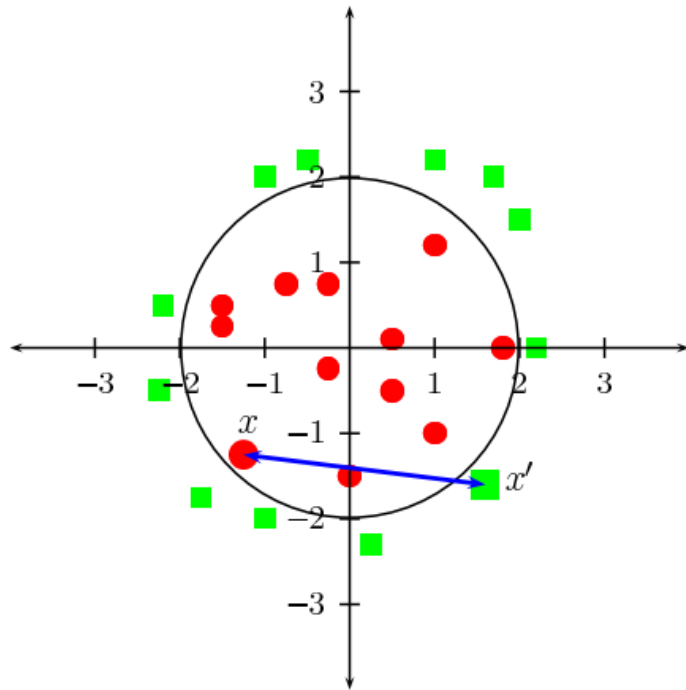


Kim, S. Y., Moon, S. K., Jung, D. C., Hwang, S. I., Sung, C. K., Cho, J. Y., Kim, S. H., Lee, J. & Lee, H. J. (2011) Pre-Operative Prediction of Advanced Prostatic Cancer Using Clinical Decision Support Systems: Accuracy Comparison between Support Vector Machine and Artificial Neural Network. *Korean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

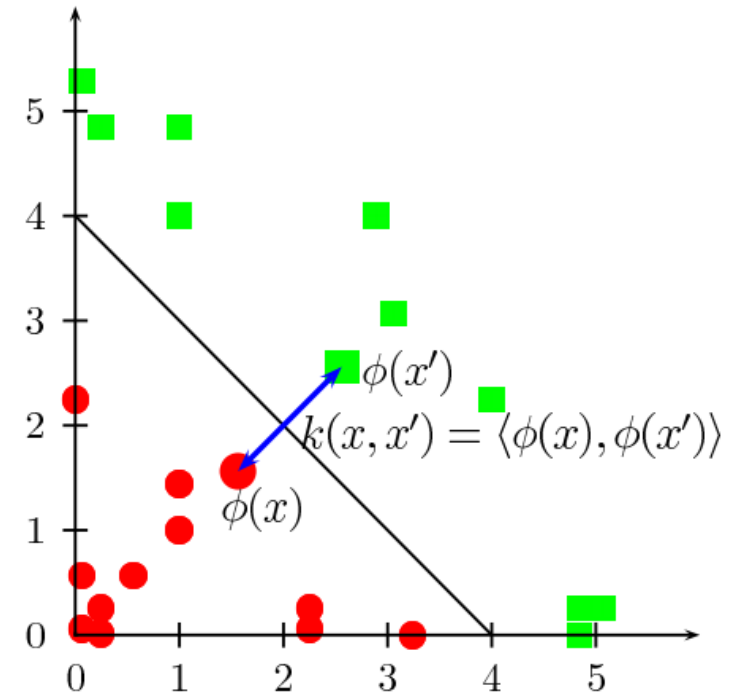


Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C.-H., Angelo, M., Ladd, C., Reich, M., Latulippe, E. & Mesirov, J. P. 2001. Multiclass cancer diagnosis using tumor gene expression signatures. *Proceedings of the National Academy of Sciences*, 98, (26), 15149-15154, doi:10.1073/pnas.211566398.





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

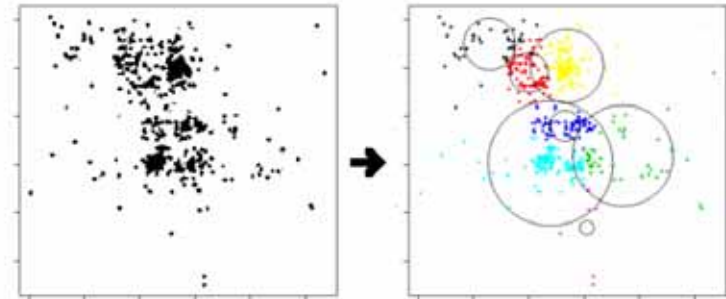


Borgwardt, K., Gretton, A., Rasch, J., Kriegel, H.-P., Schölkopf, B. & Smola, A. 2006. Integrating structured biological data by kernel max. mean discrepancy. Bioinformatics, 22, 14, e49-e57.

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));
- **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));*

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

objective fur

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

---

**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

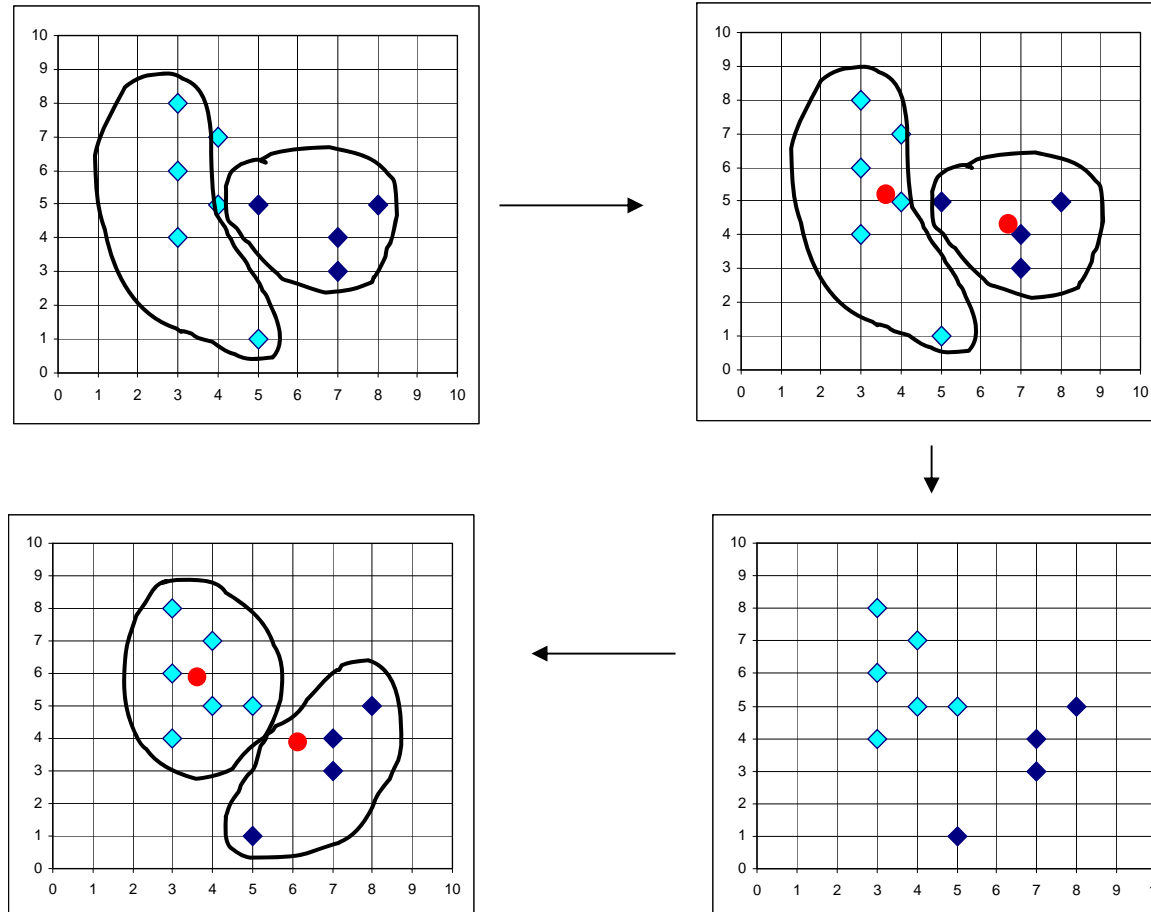
1. Partition  $X$  into a clustering  $C = (C_1, \dots, C_k)$  by assigning  $x_j \in X$  to a cluster  $C_i$  that is closest to site  $s_i \in S$ .
  2. Update each site  $s_i$  as the center of gravity of cluster  $C_i$ ; if  $|C_i| = 0$ , choose  $s_i = x_l$  for a random  $l \leq n$  with  $x_l \neq s_j$  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.

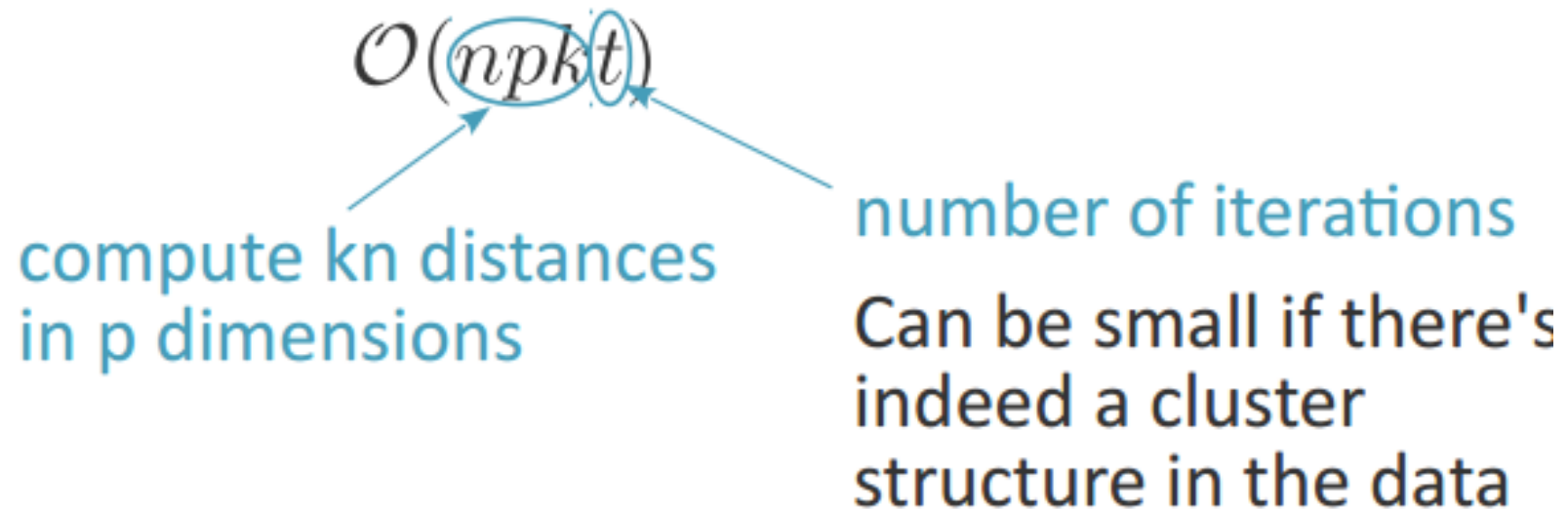
Our results can also be useful when preparing recommendations for antibiotics prescription and to guide the standardized health data record.



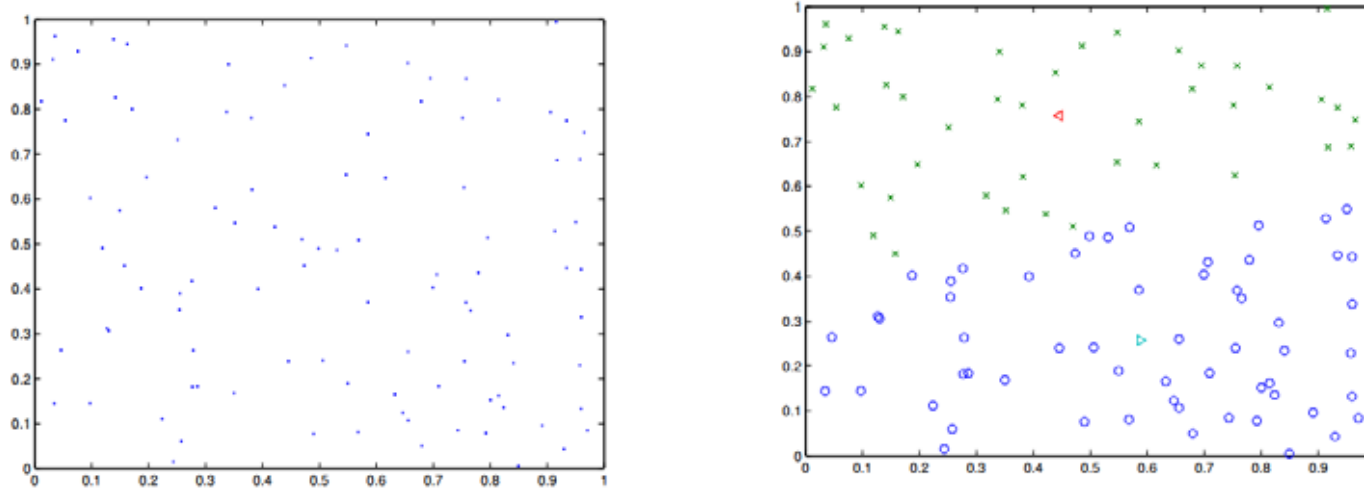
Yildirim, P., Majnarić, L., Ekmekci, O. I. & Holzinger, A. 2013. On the Prediction of Clusters for Adverse Reactions and Allergies on Antibiotics for Children to Improve Biomedical Decision Making. In: Lecture Notes in Computer Science LNCS 8127. 431-445



- What is the computational time of k-means?

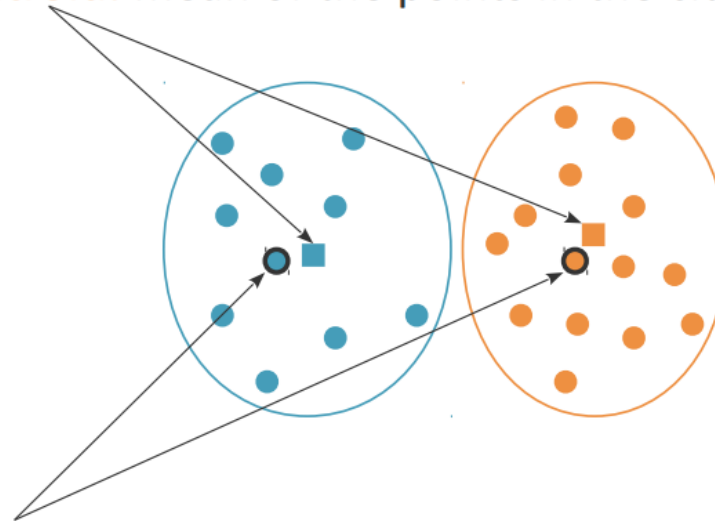


Jain, A. K. 2010. Data clustering: 50 years beyond K-means. *Pattern recognition letters*, 31, (8), 651-666.



- **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.

$$m = \arg \min_{x \in C} d(x, \mu)$$



(“Applied ML is basically feature engineering.  
*Andrew Y. Ng*”).

## 2) Feature Engineering

- 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

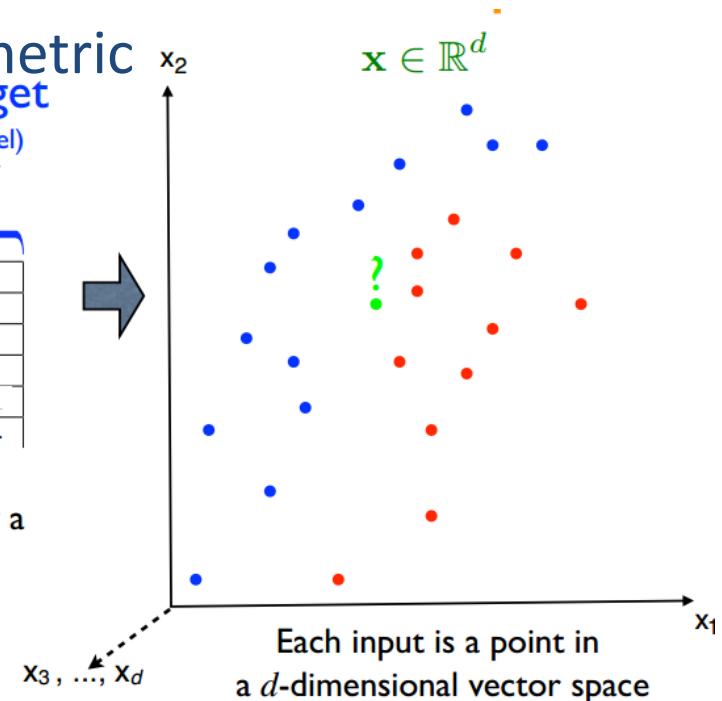
input  $x \in \mathbb{R}^d$  target (label)  $y$

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.

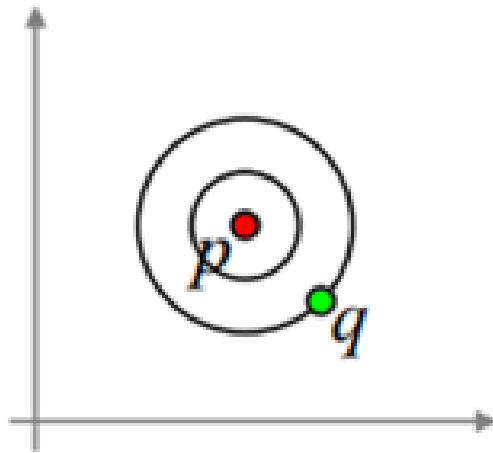
$$d : \mathcal{X} \rightarrow \mathbb{R}$$

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

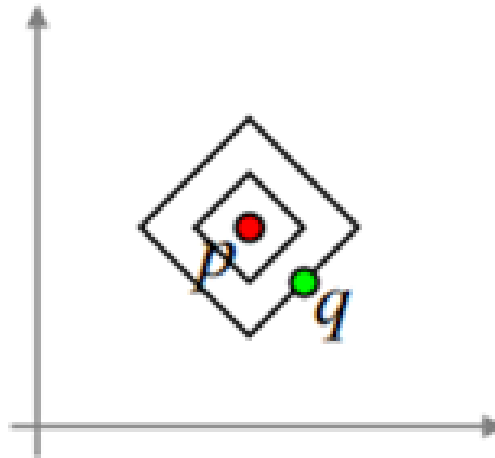
$$d(x^1, x^2) = d(x^2, x^1) \text{ symmetry}$$

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

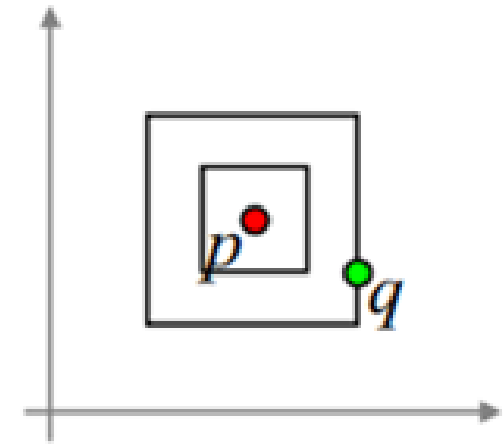
Look at the examples below, which distance measures would you select?



Euclidian norm



Manhattan norm



Maximums norm

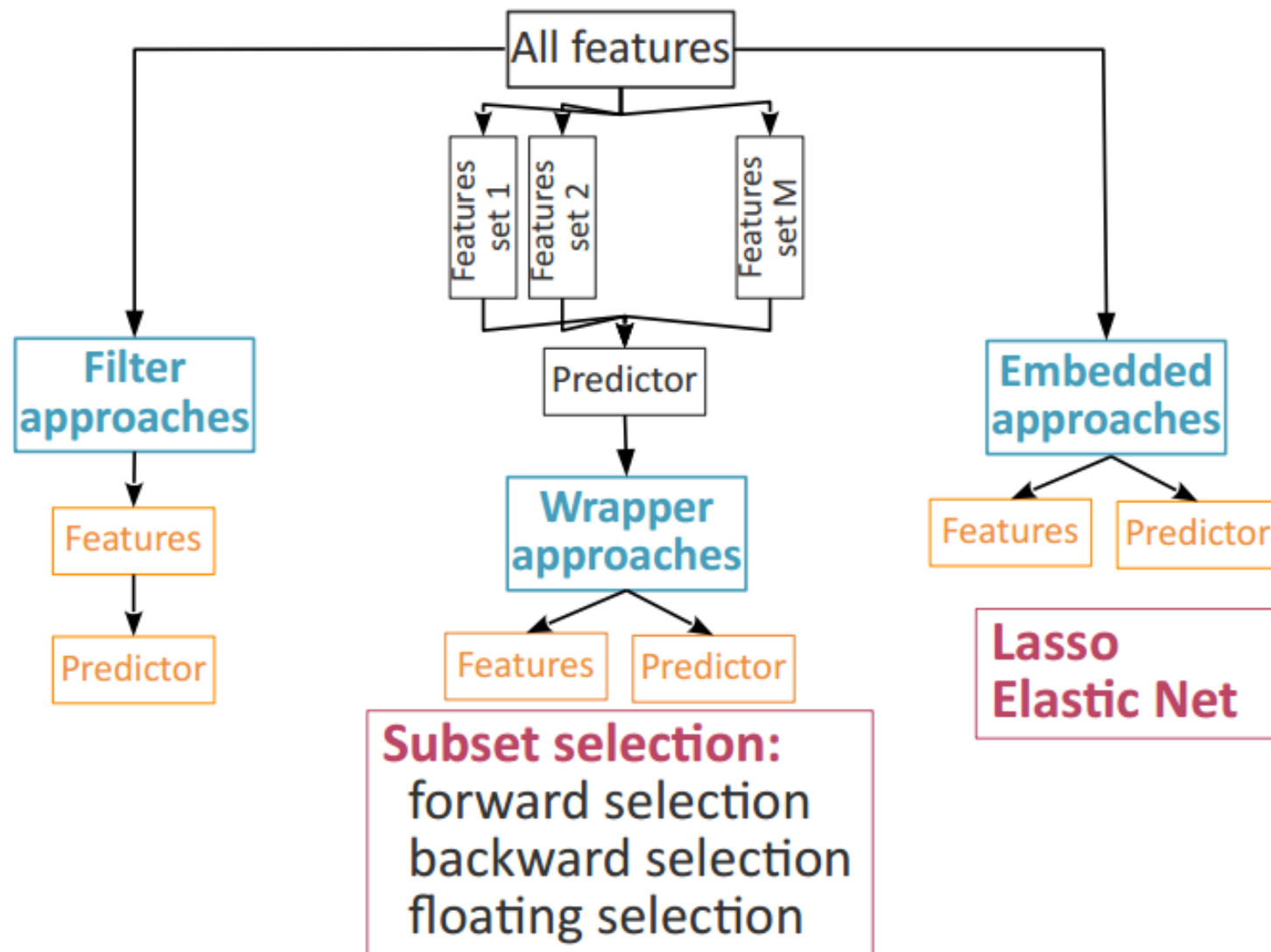
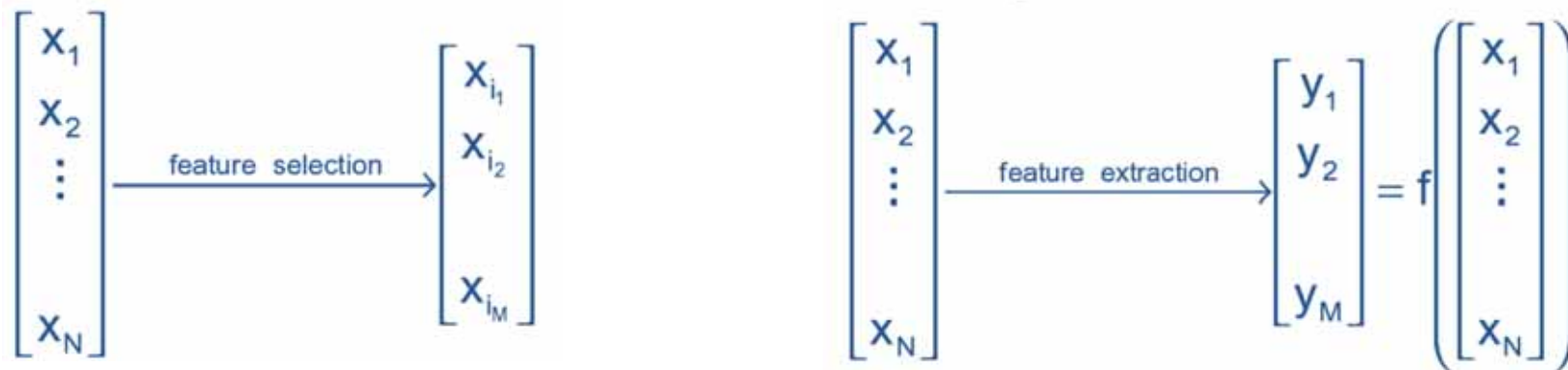


Image credit to Chloe Azencott

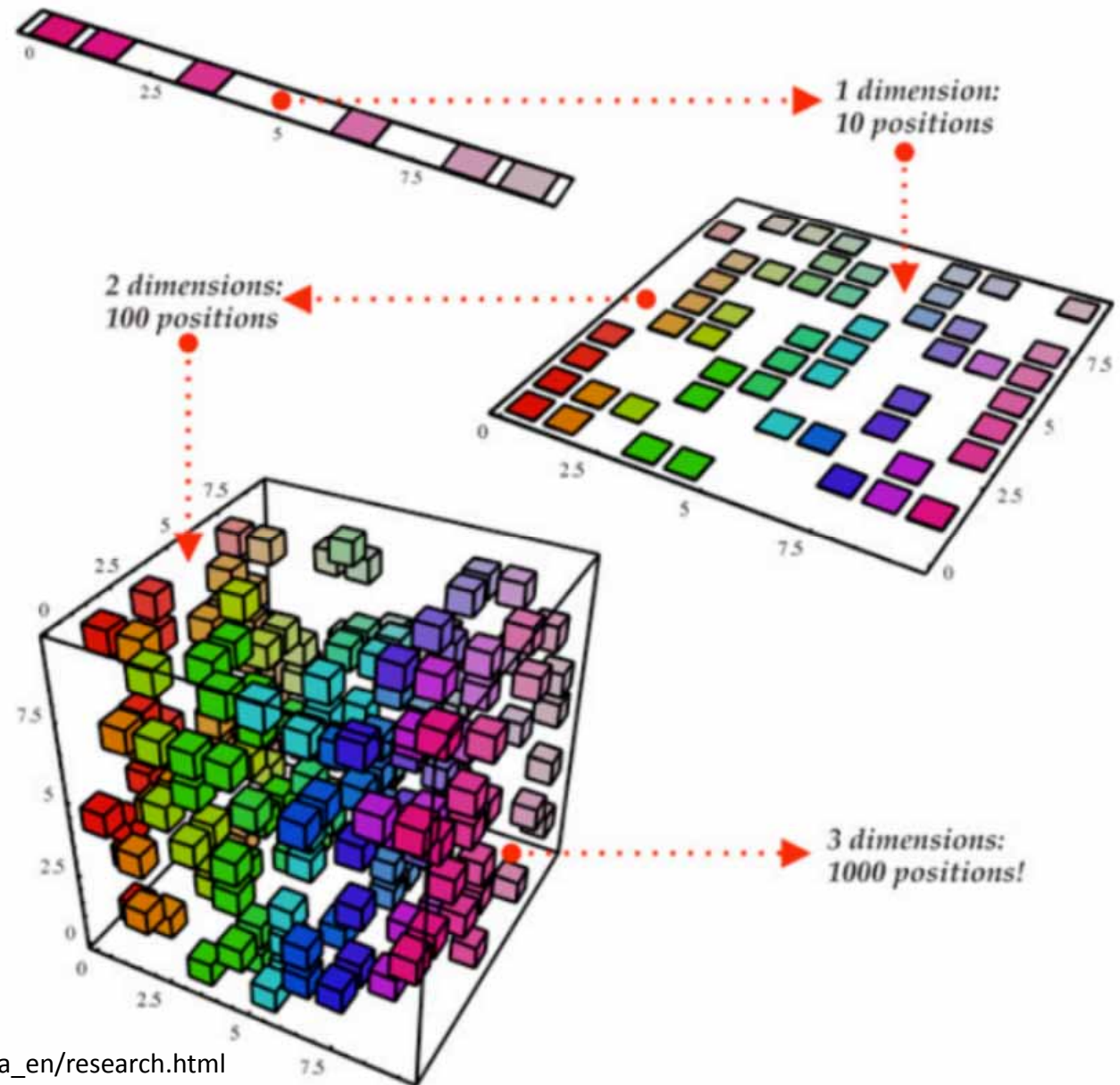
- Feature selection is just selecting a subset of the existing features without any transformation
- Feature extraction is *transforming* existing features into a lower dimensional space



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

# 3) Curse of Dimensionality

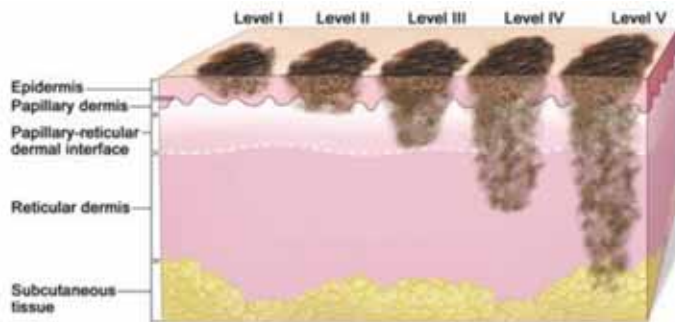




Bengio, S. & Bengio, Y.  
2000. Taking on the curse  
of dimensionality in joint  
distributions using neural  
networks. IEEE Transactions  
on Neural Networks, 11,  
(3), 550-557.

[http://www.iro.umontreal.ca/~bengioy/yoshua\\_en/research.html](http://www.iro.umontreal.ca/~bengioy/yoshua_en/research.html)

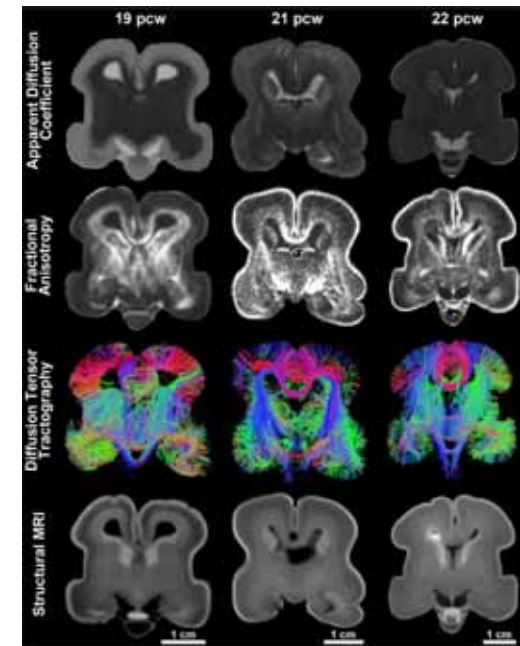
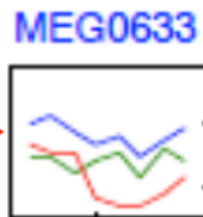
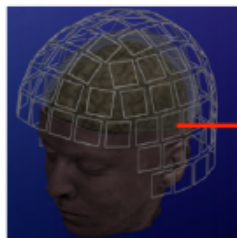
## ■ Medical Image Data (16 - 1000+ features)



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

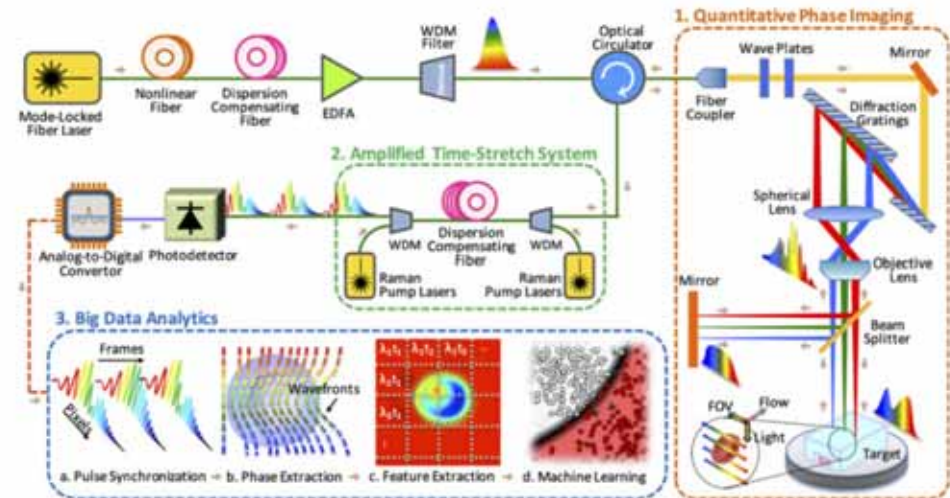
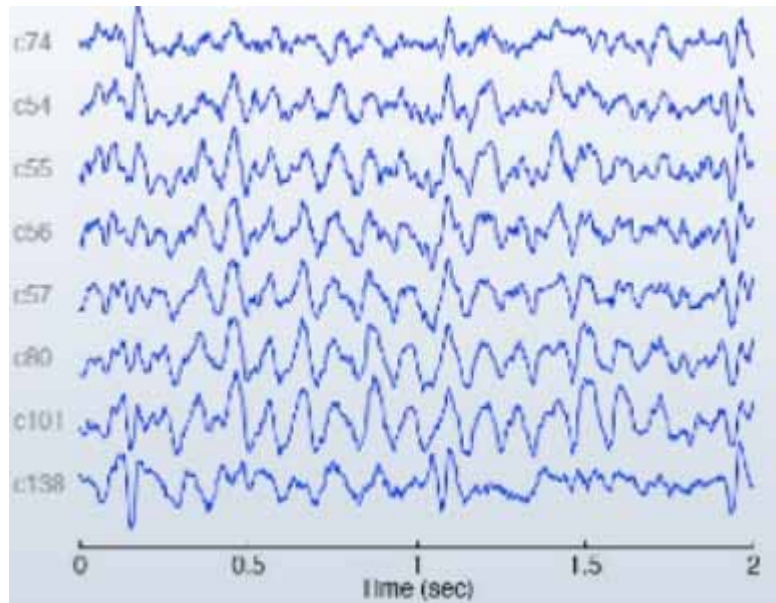
### MEG Brain Imaging

120 locations x 500 time points  
x 20 objects



Nature 508, 199–206  
doi:10.1038/nature13185

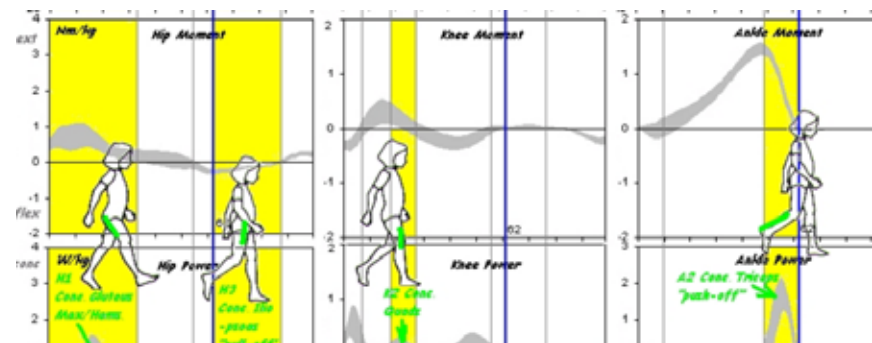
## ■ Biomedical Signal Data (10 - 1000+ features)



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



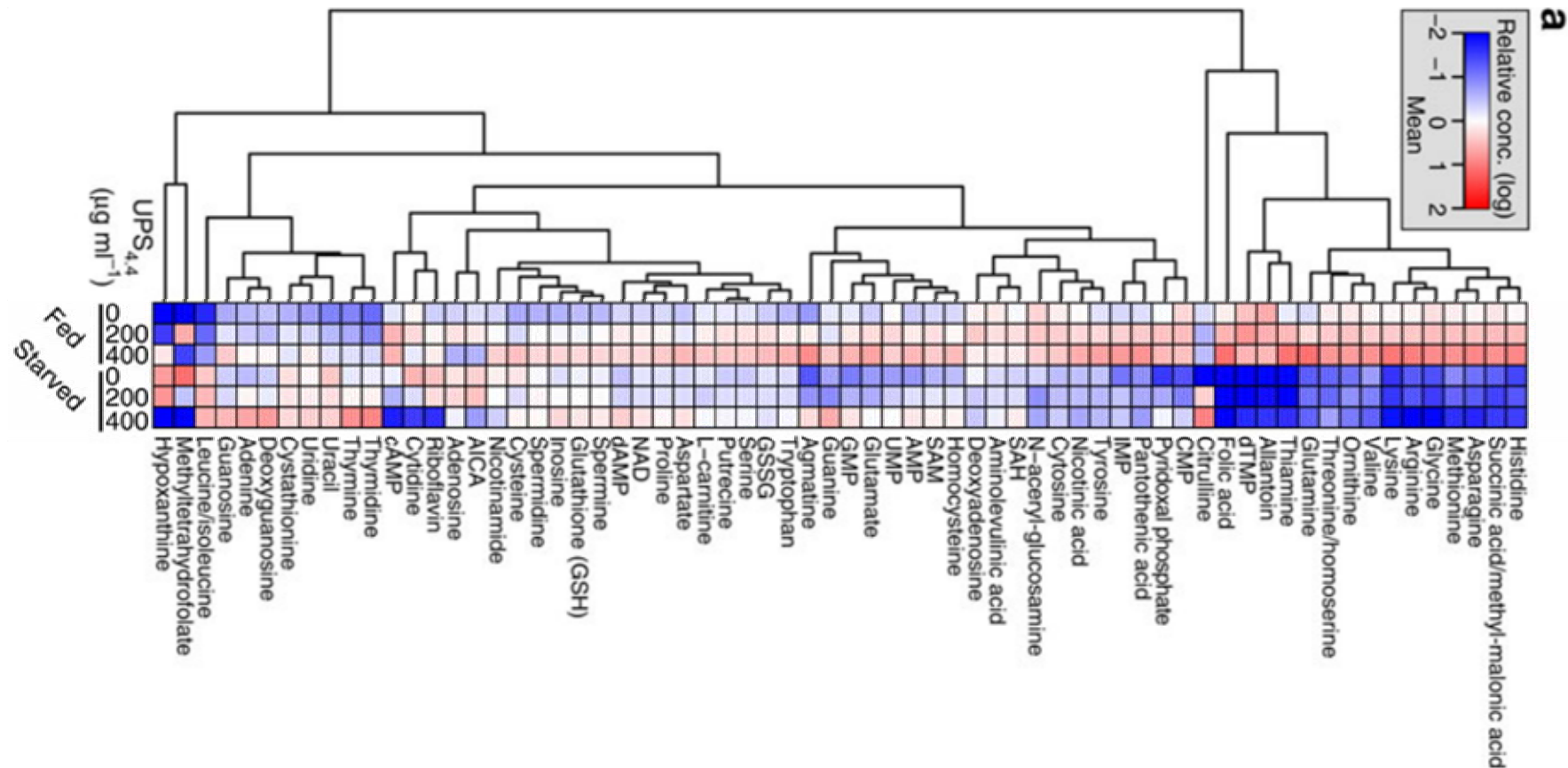
<http://www.mdpi.com/1424-8220/14/4/6124/htm>



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

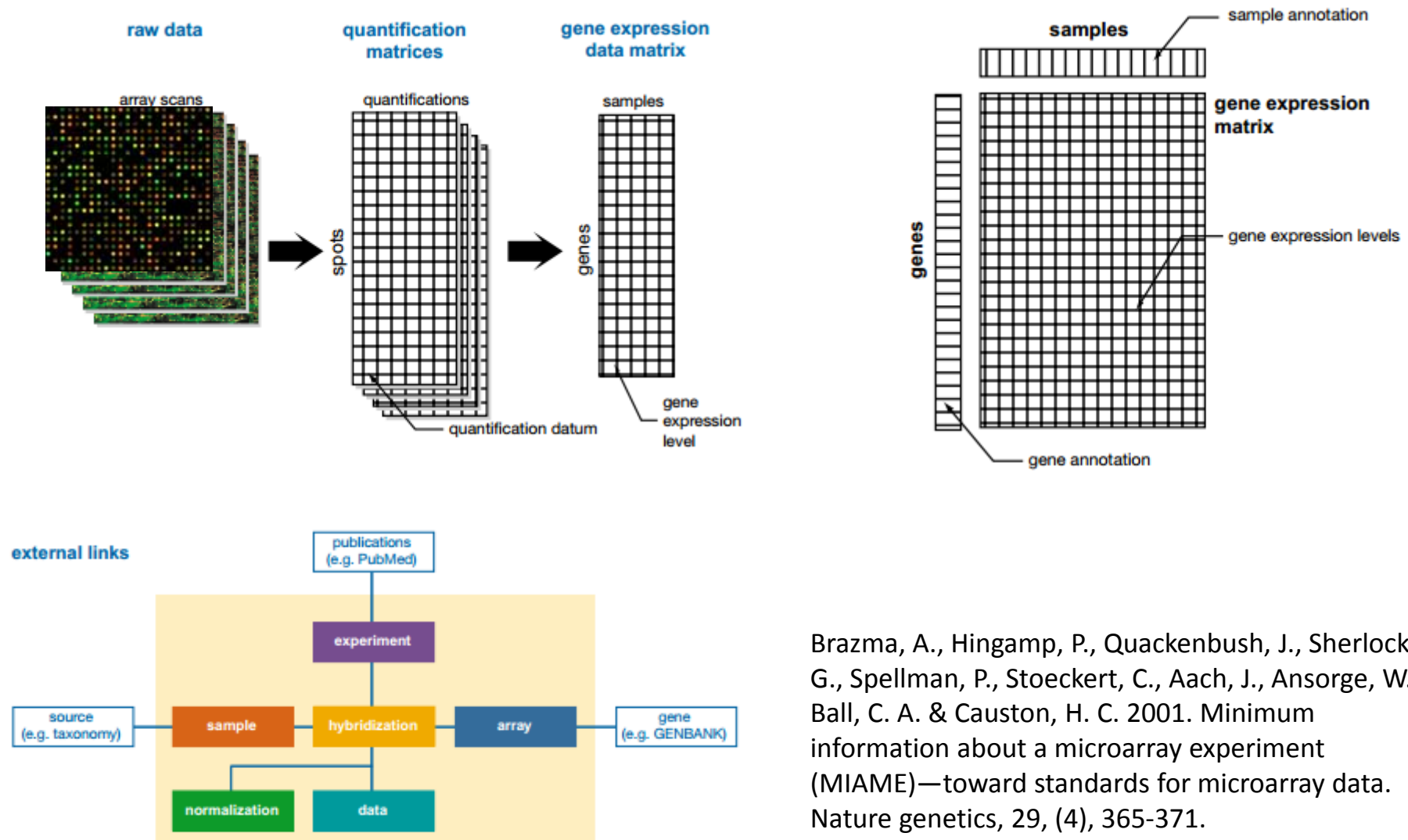


- Metabolome data (feature is the concentration of a specific metabolite; 50 – 2000+ features)



[http://www.nature.com/ncomms/2015/151005/ncomms9524/fig\\_tab/ncomms9524\\_F5.html](http://www.nature.com/ncomms/2015/151005/ncomms9524/fig_tab/ncomms9524_F5.html)

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

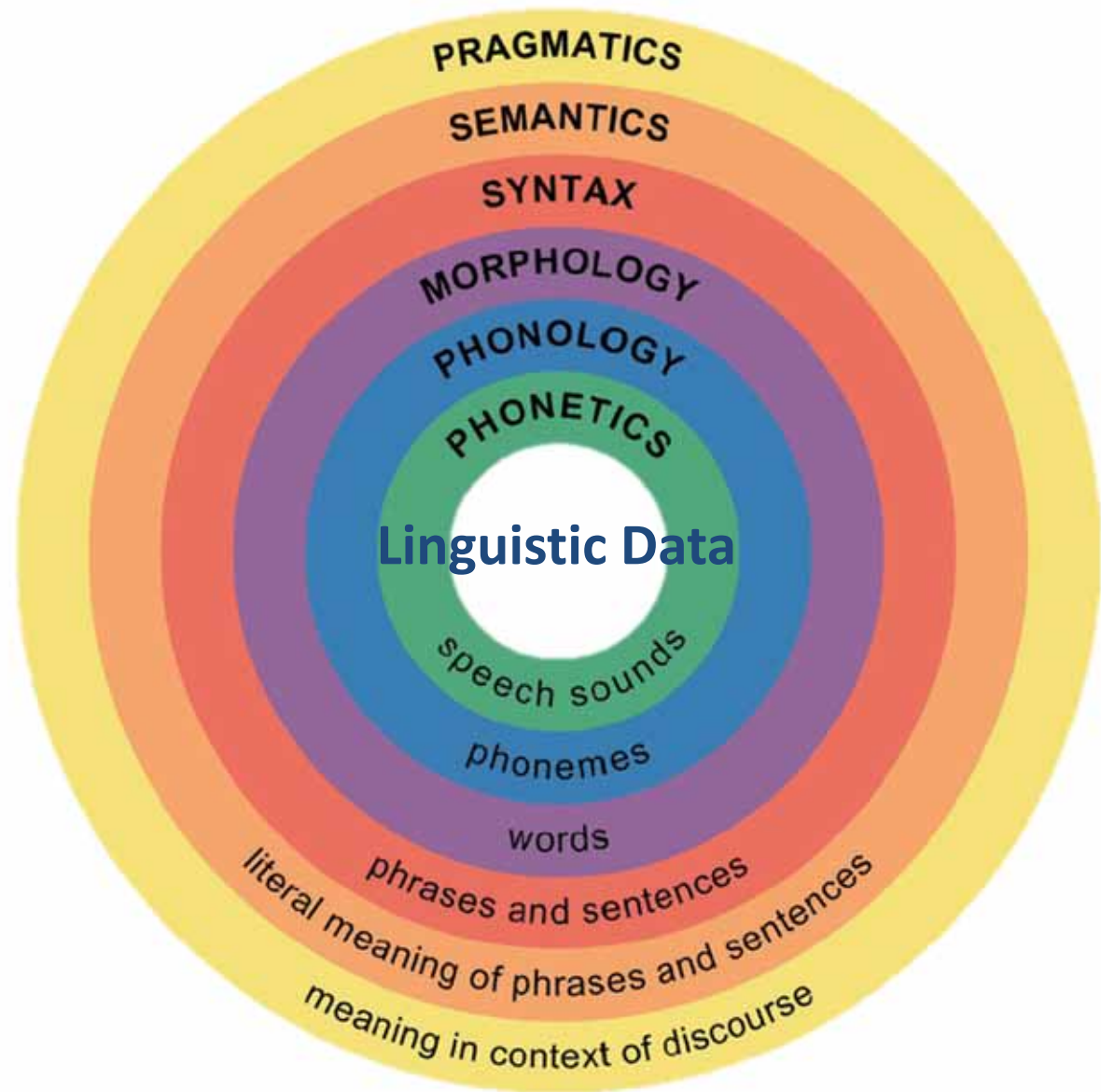


Brazma, A., Hingamp, P., Quackenbush, J., Sherlock, G., Spellman, P., Stoeckert, C., Aach, J., Ansorge, W., Ball, C. A. & Causton, H. C. 2001. Minimum information about a microarray experiment (MIAME)—toward standards for microarray data. *Nature genetics*, 29, (4), 365-371.

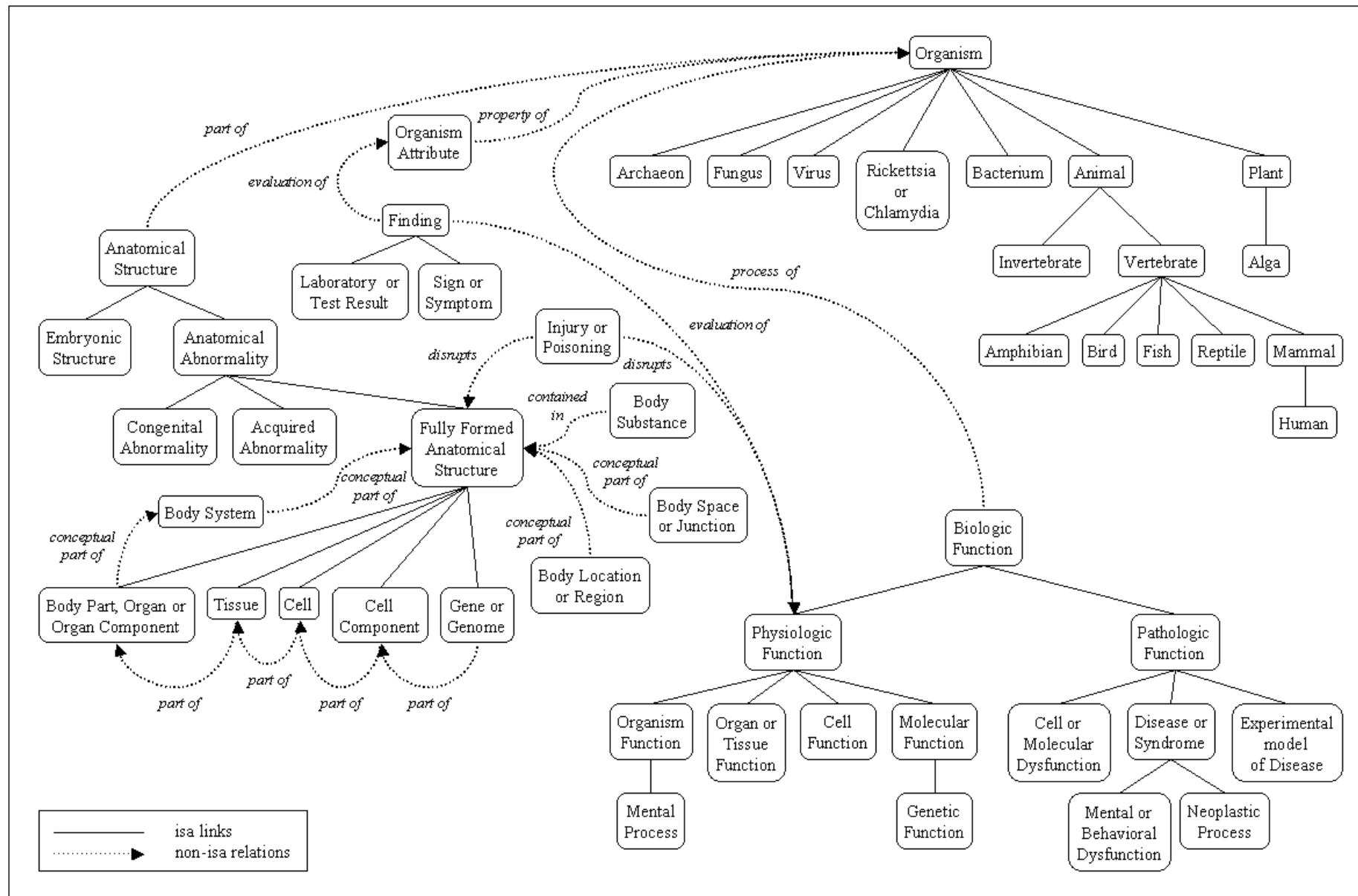
- Text  $> 10^9$  documents  $\times 10^6$  words/n-grams  
features correspond to words or terms, between  
5k to 20k features
- Text (Natural Language) is definitely very  
important for health:
  - Handwritten Notes, Drawings
  - Patient consent forms
  - Patient reports
  - Radiology reports
  - Voice dictations, annotations
  - Literature !!!



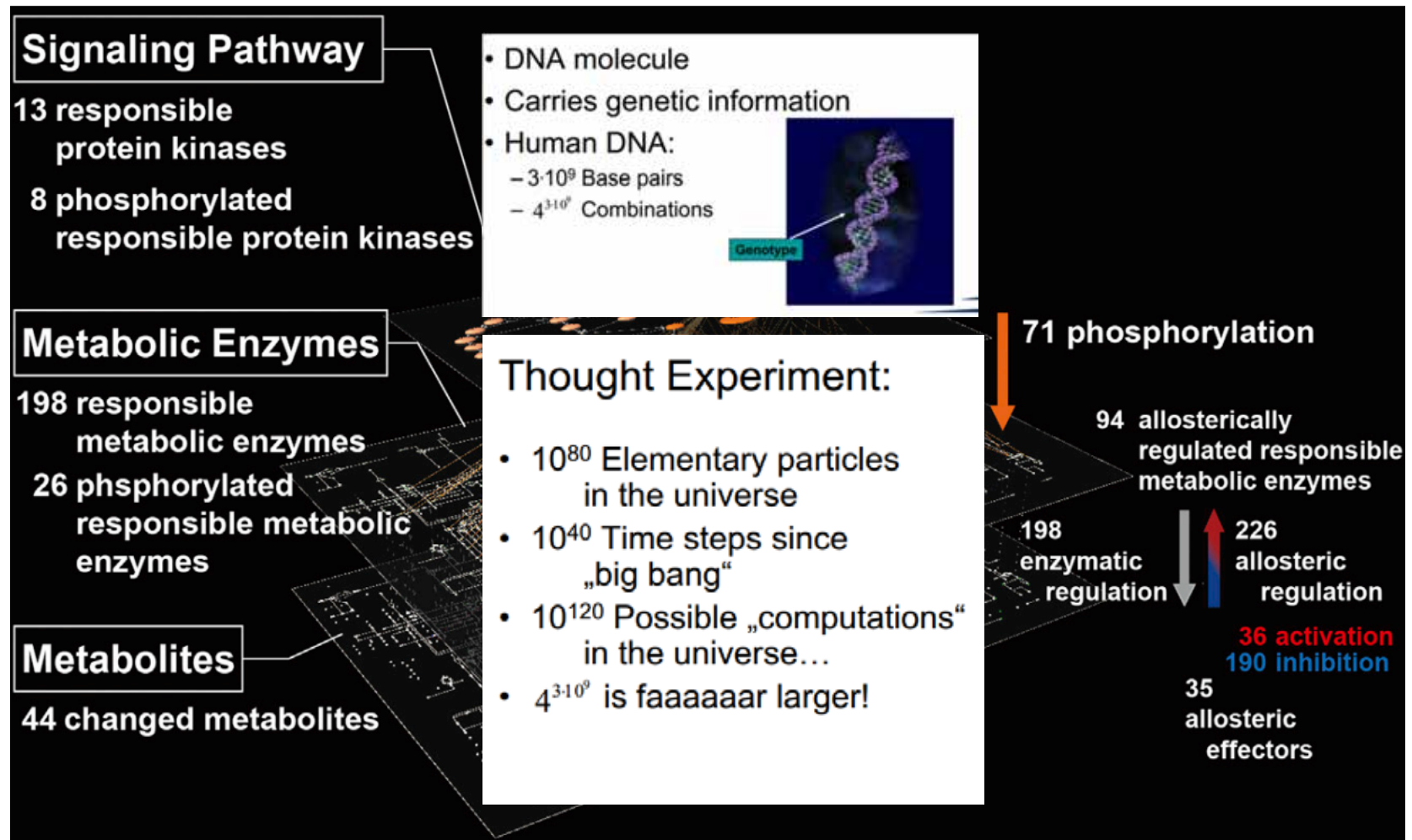
[https://www.researchgate.net/publication/255723699\\_An\\_Answer\\_to\\_Who\\_Needs\\_a\\_Stylus\\_on\\_Handwriting\\_Recognition\\_on\\_Mobile\\_Devices](https://www.researchgate.net/publication/255723699_An_Answer_to_Who_Needs_a_Stylus_on_Handwriting_Recognition_on_Mobile_Devices)



Thomas, J. J. & Cook, K. A.  
2005. *Illuminating the path:  
The research and  
development agenda for  
visual analytics*, New York,  
IEEE Computer Society Press.

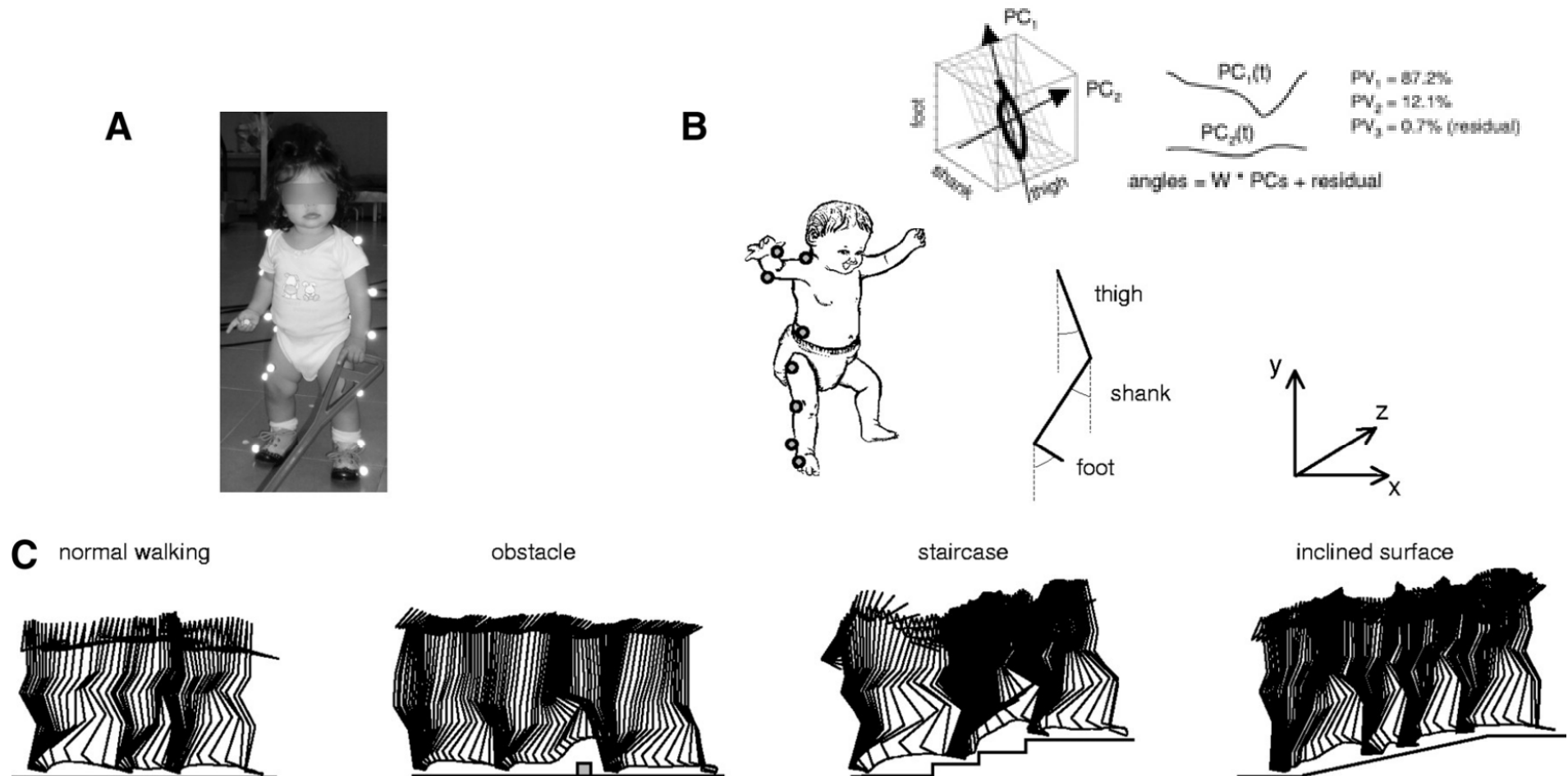






Yugi, K. et al. 2014. Reconstruction of Insulin Signal Flow from Phosphoproteome and Metabolome Data. Cell Reports, 8, (4), 1171-1183, doi:10.1016/j.celrep.2014.07.021.

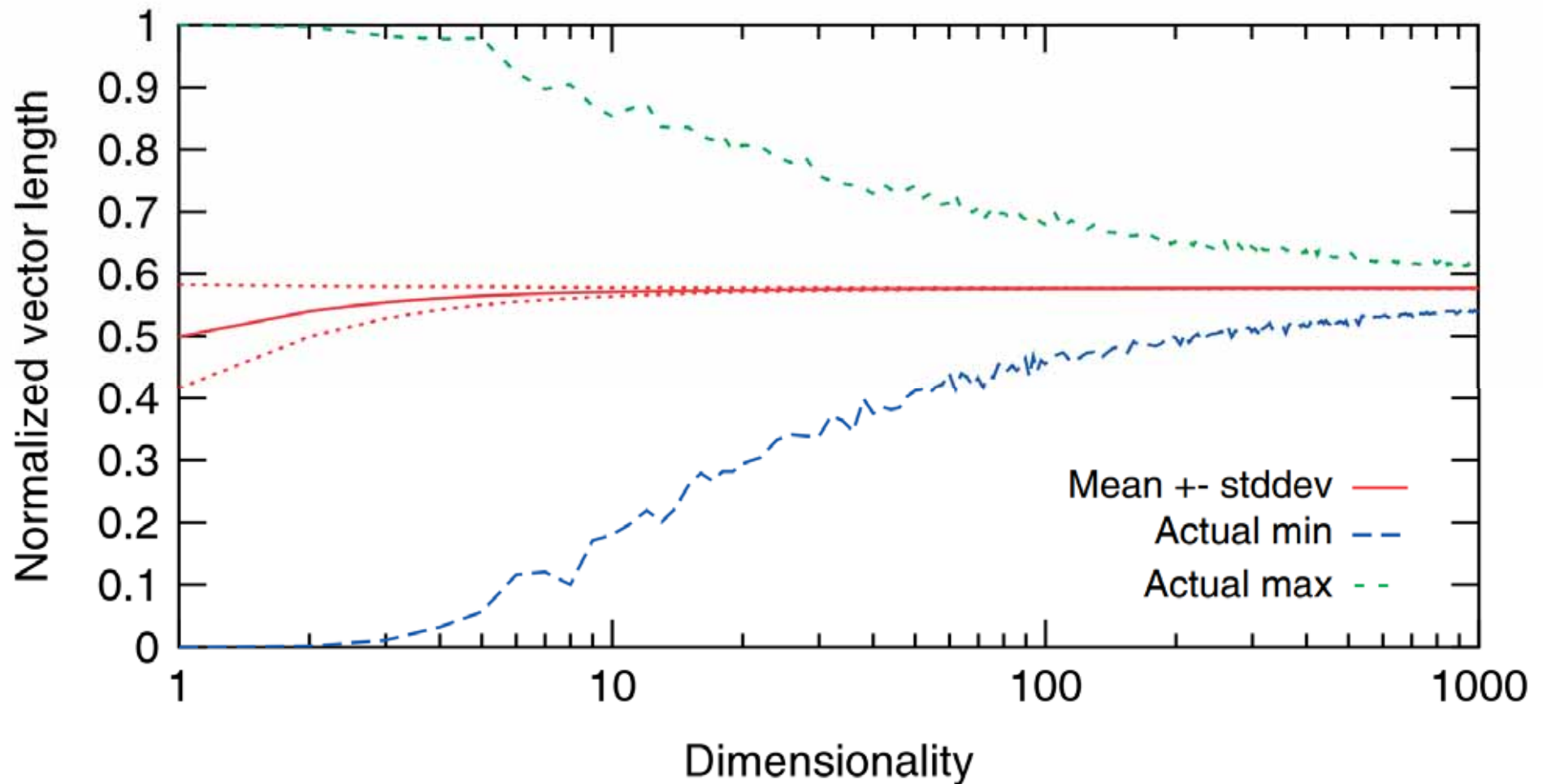
- Hyperspace is large – all points are far apart
- Computationally challenging (in time and space)
- Complexity grows with  $n$  of features
- Complex models less robust – more variance
- Statistically challenging – hard to learn
- Hard to interpret and hard to visualize
- Problem with redundant features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



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.

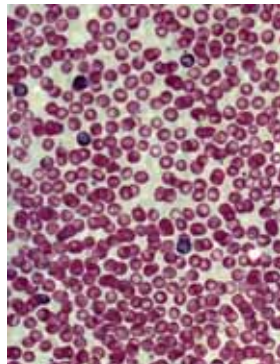
- Aspect 1: Optimization Problem
- Aspect 2: Concentration Effect
- Aspect 3: Irrelevant Attributes
- Aspect 4: Correlated Attributes

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.



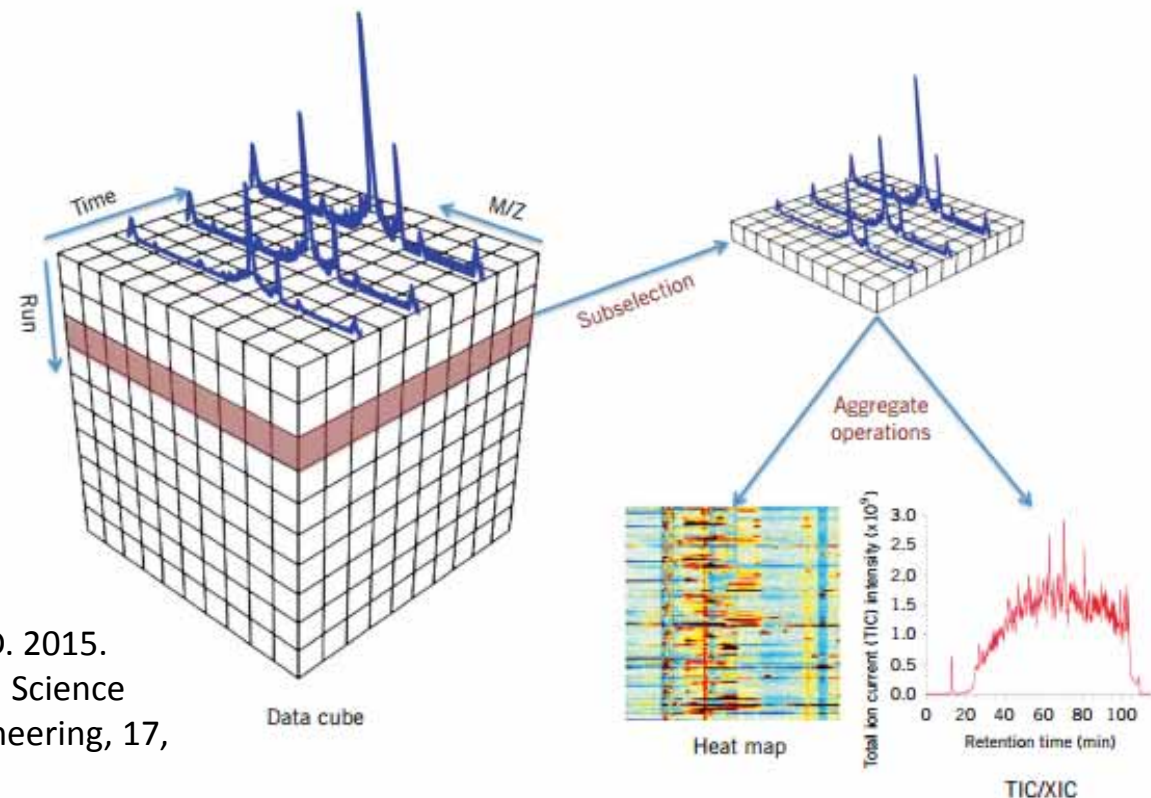
Zimek, A., Schubert, E. & Kriegel, H. P. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5, (5), 363-387, doi:10.1002/sam.11161.





Amino acids (symbols)	Fatty acids (symbols)	Fatty acids (symbols)
Alanine (Ala)	Free carnitine (C0)	Hexadecenoyl-carnitine (C16:1)
Arginine (Arg)	Acetyl-carnitine (C2)	Octadecenoyl-carnitine (C18:1)
Argininosuccinate (Argsuc)	Propionyl-carnitine (C3)	Decenoyl-carnitine (C10:2)
Citrulline (Cit)	Butyryl-carnitine (C4)	Tetradecenoyl-carnitine (C14:2)
Glutamate (Glu)	Isovaleryl-carnitine (C5)	Octadecenoyl-carnitine (C18:2)
Glycine (Gly)	Hexanoyl-carnitine (C6)	Hydroxy-isovaleryl-carnitine (C5-OH)
Methionine (Met)	Octanoyl-carnitine (C8)	Hydroxytetradecenoyl-carnitine (C14-OH)
Ornithine (Orn)	Decanoyl-carnitine (C10)	Hydroxypalmitoyl-carnitine (C16-OH)
Phenylalanine (Phe)	Dodecanoyl-carnitine (C12)	Hydroxypalmitoleyl-carnitine (C16:1-OH)
Pyroglutamate (Pyrglt)	Myristoyl-carnitine (C14)	Hydroxyoleyl-carnitine (C18:1-OH)
Serine (Ser)	Hexadecanoyl-carnitine (C16)	Dicarboxyl-butyl-carnitine (C4-DC)
Tyrosine (Tyr)	Octadecanoyl-carnitine (C18)	Glutaryl-carnitine (C5-DC)
Valine (Val)	Tiglyl-carnitine (C5:1)	Methylglutaryl-carnitine (C6-DC)
Leucine + Isoleucine (Xle)	Decenoyl-carnitine (C10:1)	Methylmalonyl-carnitine (C12-DC)
	Myristoleyl-carnitine (C14:1)	

Fourteen amino acids and 29 fatty acids are analyzed 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.

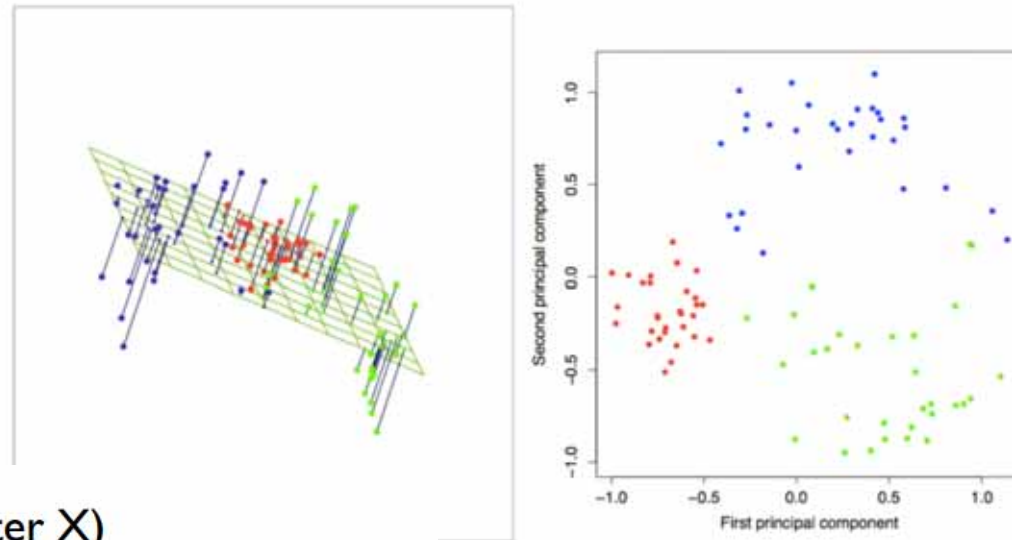
# 3) Dimensionality Reduction



- 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  $r < d$  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”

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



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

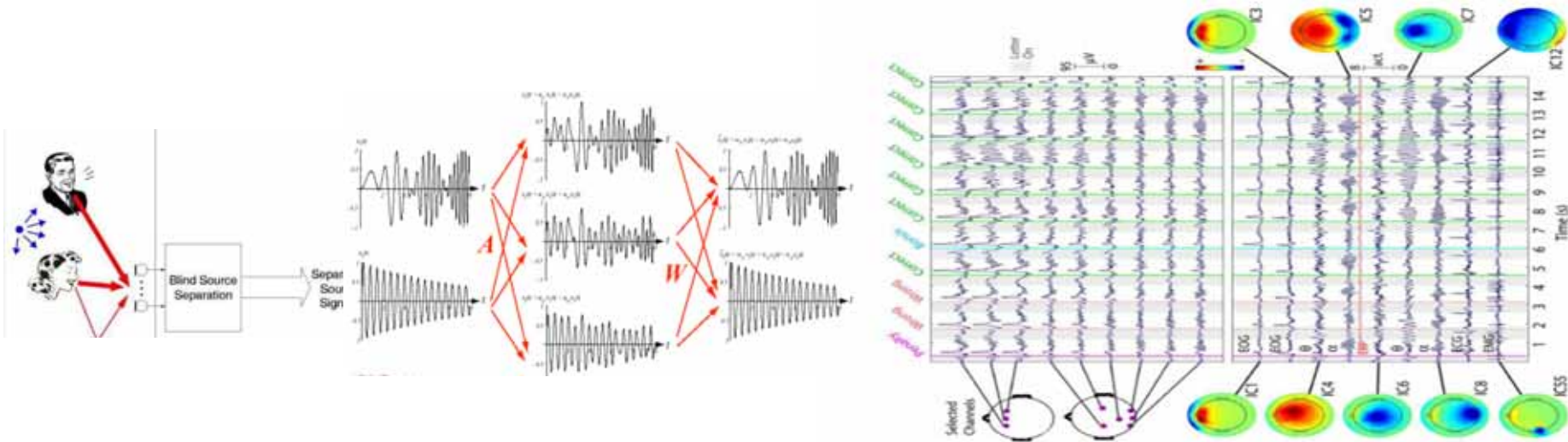
Hastie, T., Tibshirani, R. & Friedman, J. 2009. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second Edition, New York, Springer, doi:10.1007/978-0-387-84858-7.

- 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

**Factor analysis** describes the variability of observations in terms of unobserved latent variables, called **factors**, and noise

- factors explain correlation between the variables
- remaining variance is explained by Gaussian noise

factor analysis is a generative approach and models both the noise of the observations and their correlation

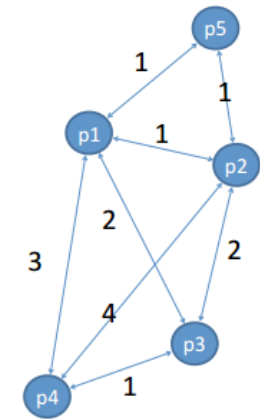
assumptions on the distribution of factors and noise



–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

$\delta_{ij}$  Distance between  $x_i$  and  $x_j$

$d_{ij}$  Distance between  $y_i$  and  $y_j$

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



## Seeking Life's Bare (Genetic) Necessities

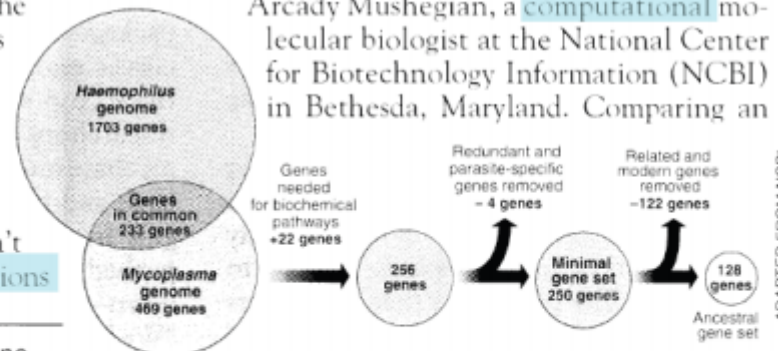
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

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

SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

ADAPTED FROM NCBI

## A Global Geometric Framework for Nonlinear Dimensionality Reduction

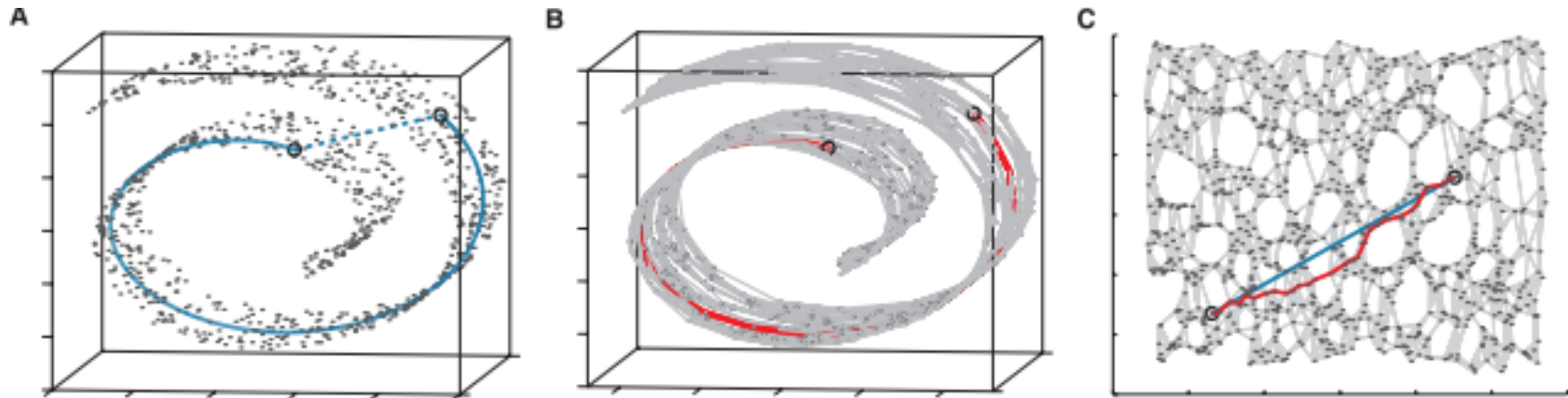
Joshua B. Tenenbaum,<sup>1\*</sup> Vin de Silva,<sup>2</sup> John C. Langford<sup>3</sup>

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.

**Locally linear embedding** (LLE) computes low-dimensional, neighborhood-preserving embeddings / representations.

LLE performs nonlinear mappings. The objective is

$$\varepsilon(\mathbf{W}) = \sum_i \left\| \mathbf{x}_i - \sum_{j=1}^k W_{ij} \mathbf{x}_j \right\|^2 \quad \sum_{j=1}^k W_{ij} = 1$$

Optimized by constrained least squares using neighbors  $\mathbf{x}_j$  of  $\mathbf{x}_i$

The solutions of this problem are invariant to rotations, rescalings, and translations of  $\mathbf{x}_i$

Down-projection optimizes  $\Phi(\mathbf{Y}) = \sum_i \left\| \mathbf{y}_i - \sum_{j=1}^k W_{ij} \mathbf{y}_j \right\|^2$   
where the  $W_{ij}$  are fixed

The representation of  $\mathbf{x}_i$  by its neighbors is transferred to  $\mathbf{y}_i$

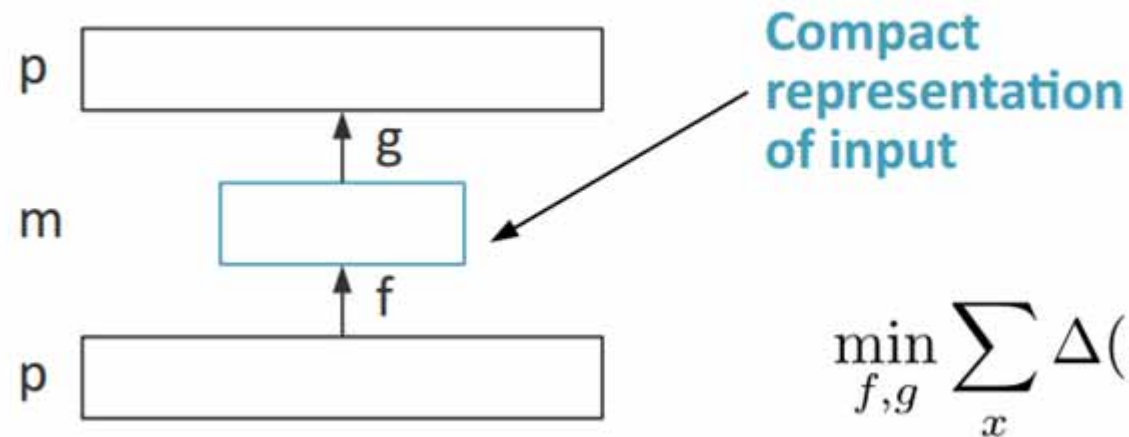
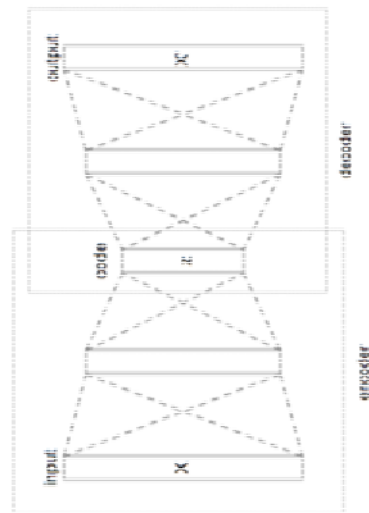
$$\Phi(\mathbf{Y}) = \sum_{ij} M_{ij} \mathbf{y}_i^T \mathbf{y}_j$$

$\delta_{ij} : 1 \text{ for } i=j, 0 \text{ otherwise}$

$$M_{ij} = \delta_{ij} - W_{ij} - W_{ji} + \sum_k W_{ki} W_{kj}$$

optimal embedding: bottom  $d$  eigenvectors of  $\mathbf{M}$  except the last one

Holzinger  $\mathbf{M} = (\mathbf{I} - \mathbf{W})^T (\mathbf{I} - \mathbf{W})$



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

- 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.



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

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

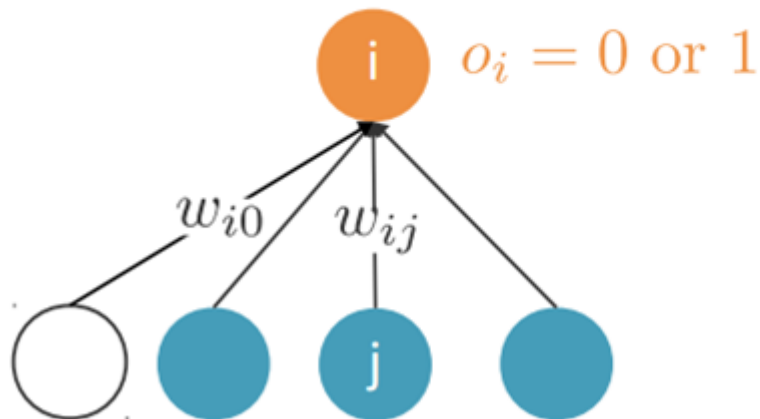
- **Linear autoencoders:** Baldi, P. & Hornik, K. 1989. Neural networks and principal component analysis: Learning from examples without local minima. Neural networks, 2, (1), 53-58.

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

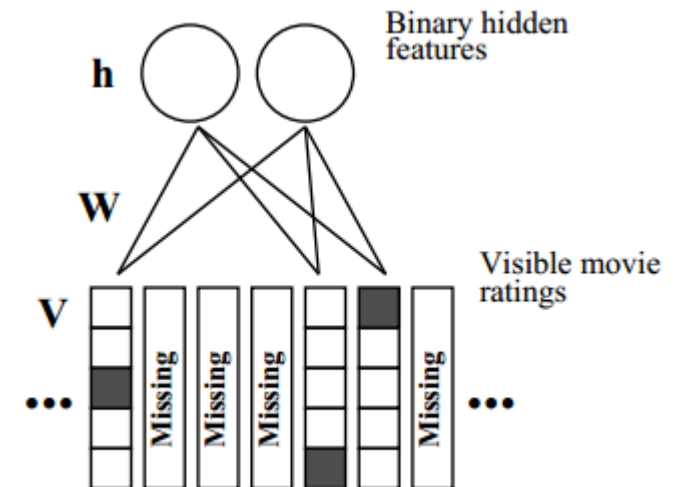
\*) David Rumelhart (1942-2011) was Cognitive Scientist working on math. Psychology

- 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)



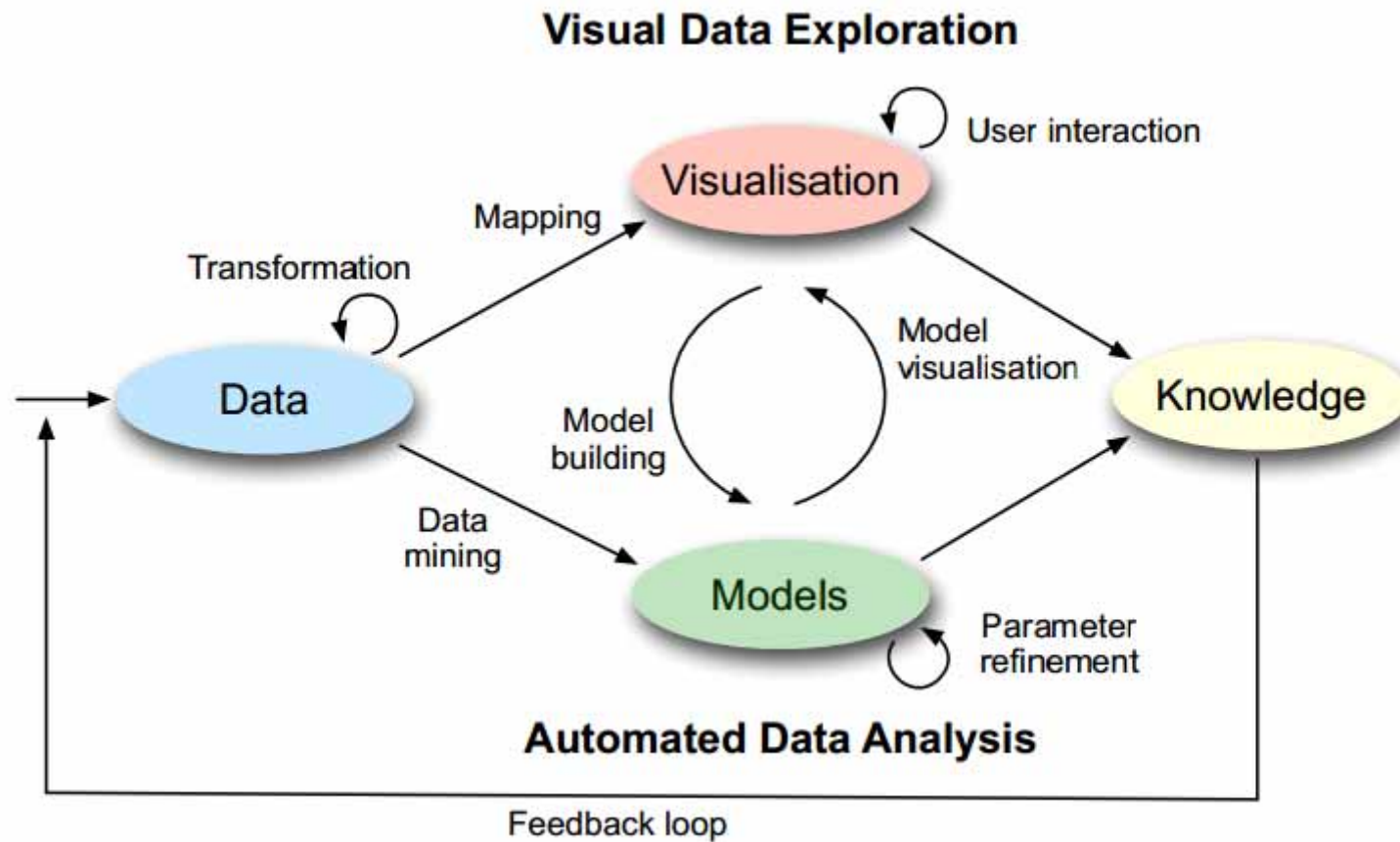
# 4) Subspace Clustering\* and Analysis

\* 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>

## Large Amount of Dimensions

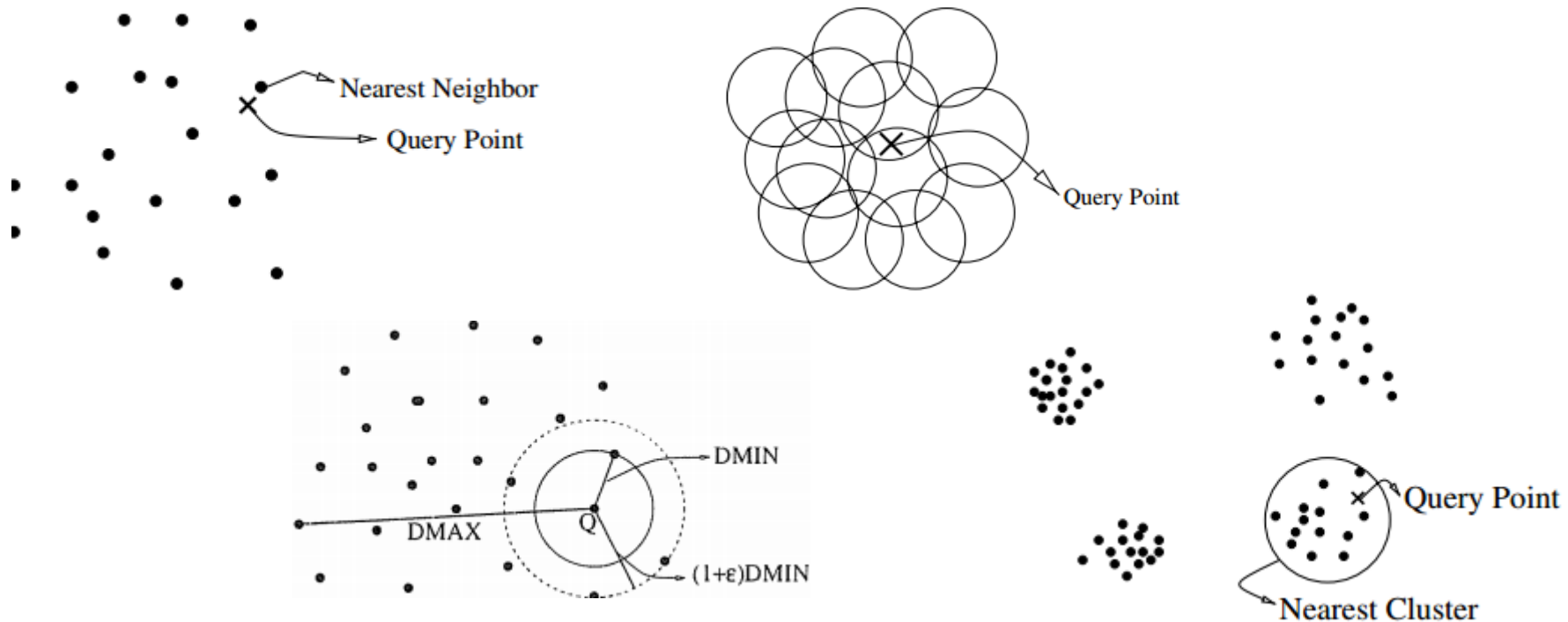
Large Amount of Records

Geography	Product Categories																																										
	Bottles				Bikes and Cages				Hydration				Tires and Tubes				Bottom Br				Derailleurs				Handbars				Pedals				Road Frame Saddles				Touring				Bike R		
Virginia	c2.362	c45	c133	c2.021	c476	c180	c94	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	c316	c561	c894	c237	c80	c15	c5	c5.919	c20	c37	c1.445	c87	c294	c32	c180	c68	c44	c79	c758	c162	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	c2.225	c302	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	c1.13	c987	c109	c23	c6	c2.128	c270	c383	c3.843	c179	c.406	c1.029	c315	c1764	c47	c54	c101	c72	c21	c1	c4	c703	c533									
Illinois	c576	c27	c33	c489	c237	c45				c1522	c159	c1.401	c1.069	c51	c5.915	c110	c14	c1.436	c25	c2.06	c700	c53	c2.564	c1.036	c728	c1.037	c21	c60	c1.046	c4	c238	c65	c0	c1.260	c390								
Indiana	c1.250	c33	c32	c1.330	c474	c45	c24	c14	c334	c48	c458	c649	c136	c23	c44	c22	c2.78	c36	c167	c153	c291	c82	c1.076	c62	c78	c9	c86	c295	c207	c35	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	c406	c545	c463	c119	c170	c21	c162	c40	c245	c40	c86	c11	c5.000	c7.500									
Michigan	c2.421	c88	c140	c2.842	c691	c60	c94	c22	c478	c24	c918	c351	c723	c33	c57	c30	c1.589	c164	c6	c2.512	c669	c831	c252	c326	c194	c0	c478	c1.186	c1.509	c62	c75	c17	c14.59	c59.000									
Missouri	c1.368	c63	c91	c1.140	c660	c75	c60			c483	c193	c406	c309	c250	c25	c23	c13	c111	c6	c58	c2.170	c1.06	c560	c297	c238	c197	c26	c255	c63	c651	c2	c24	c50	c340									
Nevada	c1.656	c122	c149	c1.621	c738					c12	c1.389	c195	c187	c672	c549	c581	c309	c18	c392	c220	c130	c3.032	c1.131	c2.410	c1.239	c1.188	c1.958	c26	c68	c990	c1.766	c4.598	c2.714	c194	c8.000								
New Mexico	c1.531	c56	c133	c1.996	c594	c105	c48	c14	c337	c129	c742	c136	c323	c64	c48	c6	c2.91	c3	c212	c1.904	c108	c571	c368	c159	c240	c3	c348	c1.83	c232	c525	c1105	c79	c9.911	c29									
New York	c3.217	c185	c312	c4.317	c2.070	c165	c108	c43	c1.571	c829	c429	c3.962	c1.461	c101	c163	c14	c3.28	c201	c127	c2.265	c3.265	c731	c80	c509	c97	c46	c128	c4.048	c4.599	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	c91	c146	c361	c0	c0	c0	c650	c4	c30	c1.020	c466	c0	c0	c11	c2.450								
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	c20	c7	c8	c63	c287	c249	c464	c42	c421	c3.532	c2.898									
Arizona	c1.327	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	c398	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	c16	c2.382	c1.166	c1.730	c3	c64	c496	c764	c841	c966	c775	c5.352	c68										
Illinois	c1.376	c150	c145	c1.721	c219	c738	c81	c51	c1.006	c3	c682	c96	c35	c43	c46	c23	c1.405	c1	c36	c168	c72	c742	c57	c55	c75	c0	c246	c1.166	c77	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	c22	c211	c235	c275	c4.286	c1									
Maine	c430	c9	c22	c558	c742	c79	c5	c10	c24	c43	c49	c43	c130	c20	c3	c1	c100	c3	c28	c2.003	c448	c109	c2	c9	c133	c7	c9	c33	c37	c34	c2	c1	c181	c727									
Michigan	c1.615	c356	c0	c2.498	c533	c48	c165	c32	c1.473	c65	c3	c133	c149	c2	c26	c6	c595	c16	c0	c2.21	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	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	c229	c2	c2.275	c997									
New Mexico	c2.209	c61	c39	c1.881	c806	c75	c72	c19	c71	c316	c561	c894	c237	c80	c87	c5	c5.919	c20	c37	c1.445	c87	c294	c32	c180	c68	c44	c79	c758	c162	c59	c5	c1	c25	c727									
New York	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	c2.225	c302	c1.194	c53	c5	c3	c1.477	c727									
Ohio	c4.422	c182	c206	c3.848	c1.068	c180	c144	c33	c1.941	c889	c1.208	c1.13	c987	c109	c23	c6	c2.128	c270	c383	c3.843	c179	c.406	c1.029	c315	c1764	c47	c54	c101	c72	c21	c1	c4	c703	c533									
Virginia	c576	c27	c33	c489	c237	c45				c1522	c159	c1.401	c1.069	c51	c5.915	c110	c14	c1.436	c25	c2.06	c700	c53	c2.564	c1.036	c728	c1.037	c21	c60	c1.046	c4	c238	c65	c0	c1.260	c390								
Arizona	c1.250	c33	c32	c1.330	c474	c45	c24	c14	c334	c48	c458	c649	c136	c23	c44	c22	c2.78	c36	c167	c153	c291	c82	c1.076	c62	c78	c9	c86	c295	c207	c35	c6	c1	c87	c578									
Indiana	c2.069	c69	c137	c1.948	c507	c60	c132	c12	c372	c259	c701	c476	c324	c49	c242	c8	c1.375	c42	c1	c2.326	c406	c545	c463	c119	c170	c21	c162	c40	c245	c40	c86	c11	c5.000	c7.500									
Maine	c2.421	c88	c140	c2.842	c691	c60	c94	c22	c478	c24	c918	c351	c723	c33	c57	c30	c1.589	c164	c6	c2.512	c669	c831	c252	c326	c194	c0	c478	c1.186	c1.509	c62	c75	c17	c14.59	c59.000									
Michigan	c1.368	c63	c91	c1.140	c660	c75	c60			c483	c193	c406	c309	c250	c25	c23	c13	c111	c6	c58	c2.170	c1.06	c560	c297	c238	c197	c26	c255	c63	c651	c2	c24	c50	c340									
Missouri	c1.656	c122	c149	c1.621	c738					c12	c1.389	c195	c187	c672	c549	c581	c309	c18	c392	c220	c130	c3.032	c1.131	c2.410	c1.239	c1.188	c1.958	c26	c68	c990	c1.766	c4.598	c2.714	c194	c8.000								
Nevada	c1.531	c56	c133	c1.996	c594	c105	c48	c14	c337	c129	c742	c136	c323	c64	c48	c6	c2.91	c3	c212	c1.904	c108	c571	c368	c159	c240	c3	c348	c1.83	c232	c525	c1105	c79	c9.911	c29									
New Mexico	c3.217	c185	c312	c4.317	c2.070	c165	c108	c43	c1.571	c829	c429	c3.962	c1.461	c101	c163	c14	c3.28	c201	c127	c2.265	c3.265	c731	c80	c509	c97	c46	c128	c4.048	c4.599	c1.416	c1.500	c253	c4.003	c2.219									
New York	c1.656	c51	c67	c1.091	c462					c1	c1.249	c194	c286	c487	c266	c0	c0	c1	c454	c101	c91	c146	c361	c0	c0	c0	c650	c4	c30	c1.020	c466	c0	c0	c11	c2.450								
Ohio	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	c20	c7	c8	c63	c287	c249	c464	c42	c421	c3.532	c2.898									
Arizona	c1.327	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	c398	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	c16	c2.382	c1.166	c1.730	c3	c64	c496	c764	c841	c966	c775	c5.352	c68										
Illinois	c1.376	c150	c145	c1.721	c219	c738	c81	c51	c1.006	c3	c682	c96	c35	c43	c46	c23	c1.405	c1	c36	c168	c72	c742	c57	c55	c75	c0	c246	c1.166	c77	c211	c235	c275	c4.286	c1									
Indiana	c38	c8	c20	c334	c1.075	c18	c4	c15	c19	c45	c82																																

- 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: Beerli, C. & Buneman, P. (eds.) Database Theory ICDT 99, LNCS 1540.* Berlin: Springer, pp. 217-235.

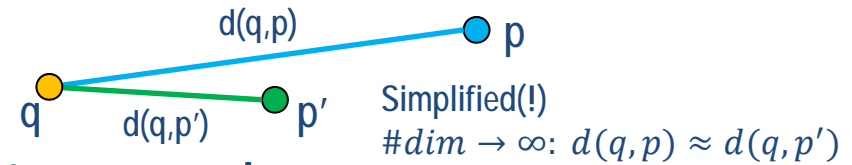
- NN problem: Given  $n$  data points and a query point in an  $m$  –dimensional metric space
- find the data point closest to the query point.



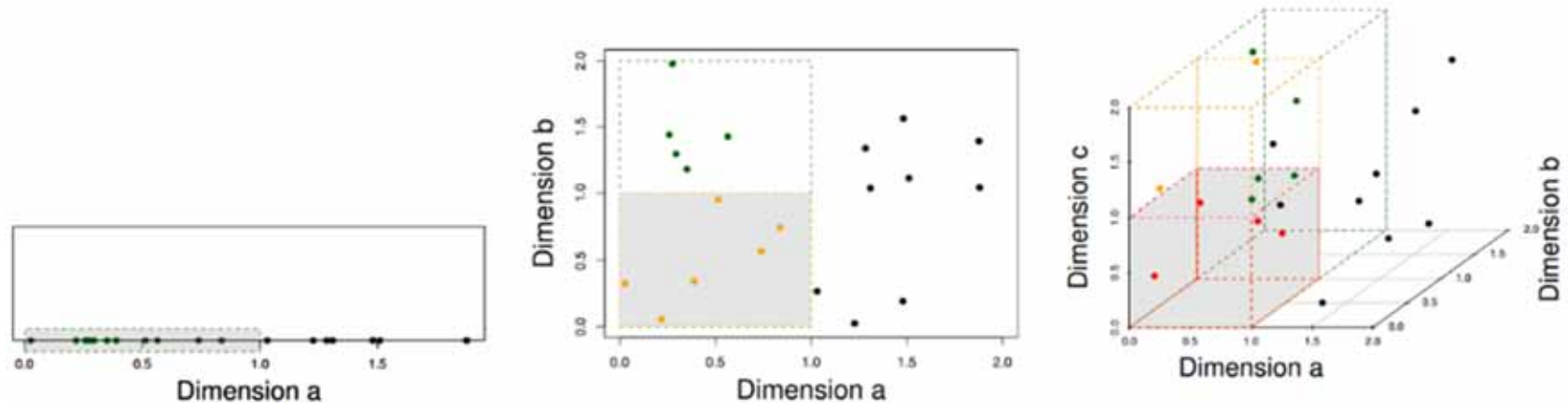
Beyer, K., Goldstein, J., Ramakrishnan, R. & Shaft, U. 1999. When is "nearest neighbor" meaningful? *In: Beer, C. & Buneman, P. (eds.) Database Theory ICDT 99, LNCS 1540.* Berlin: Springer, pp. 217-235.

- Concentration Effect

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



- High-Dimensional Data is Sparse

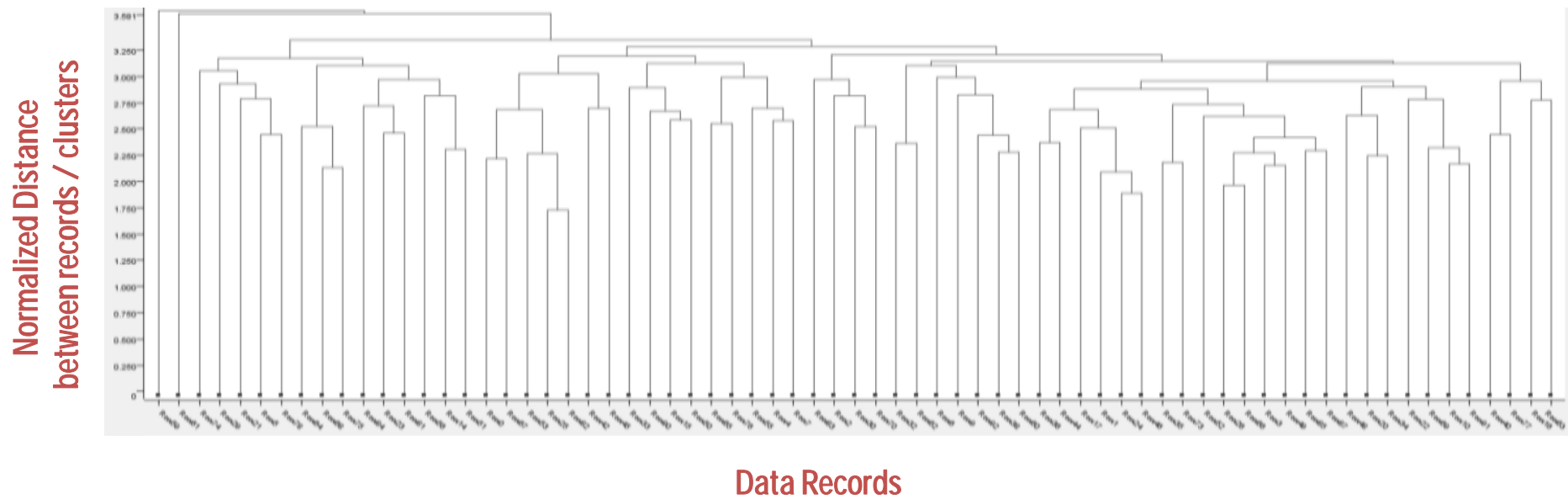


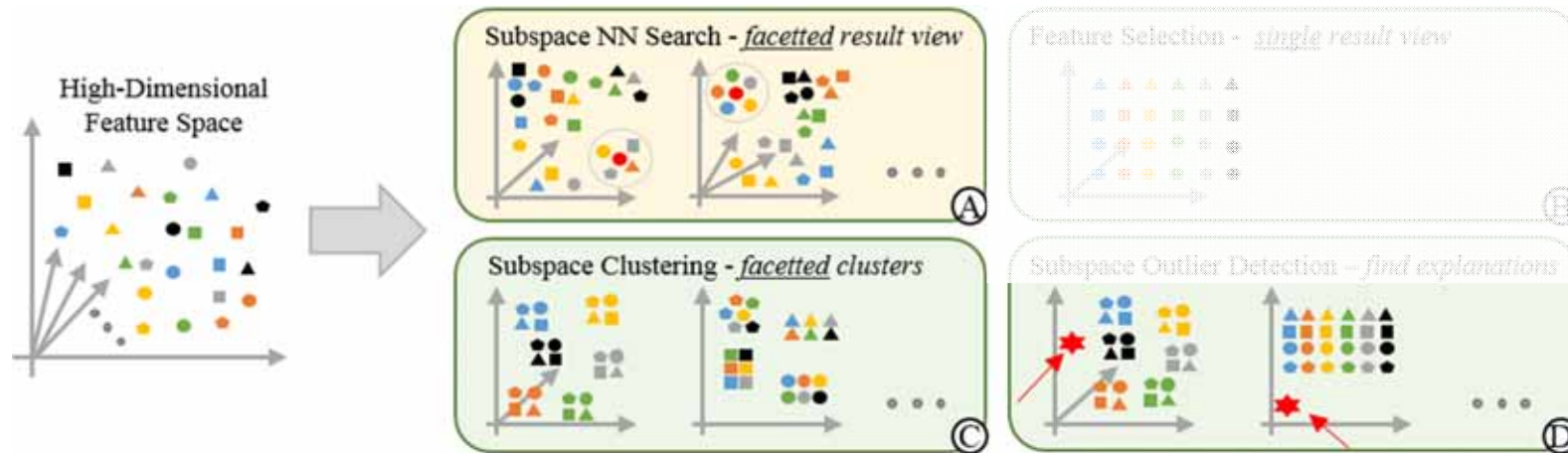
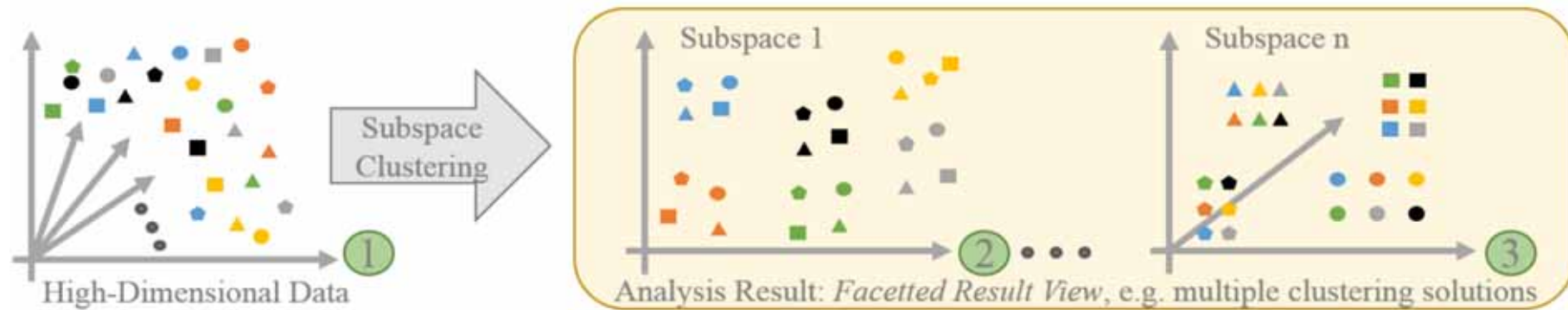
Optimization Problem and Combinatorial Issues

Feature selection and dimension reduction

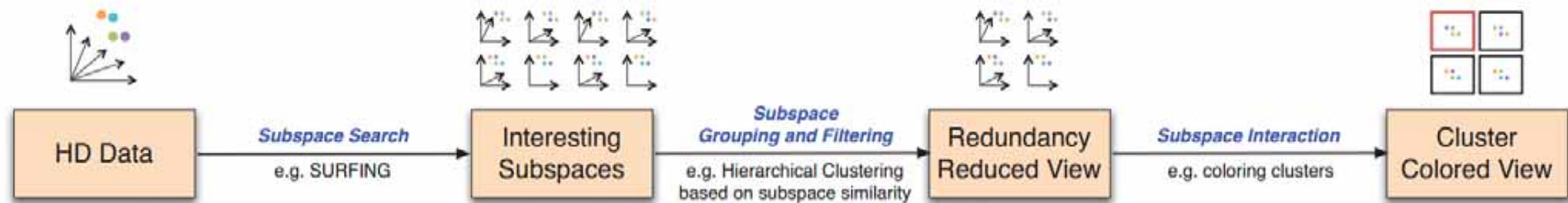
$2^d - 1$  possible subsets of dimensions (  $\rightarrow$  subspaces)



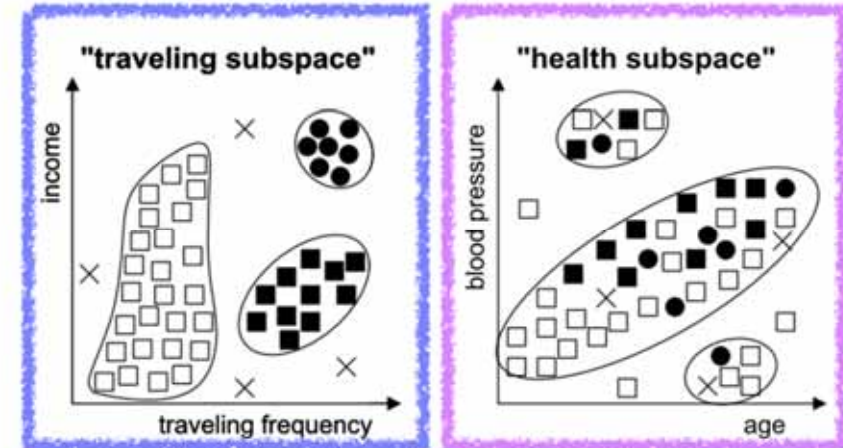




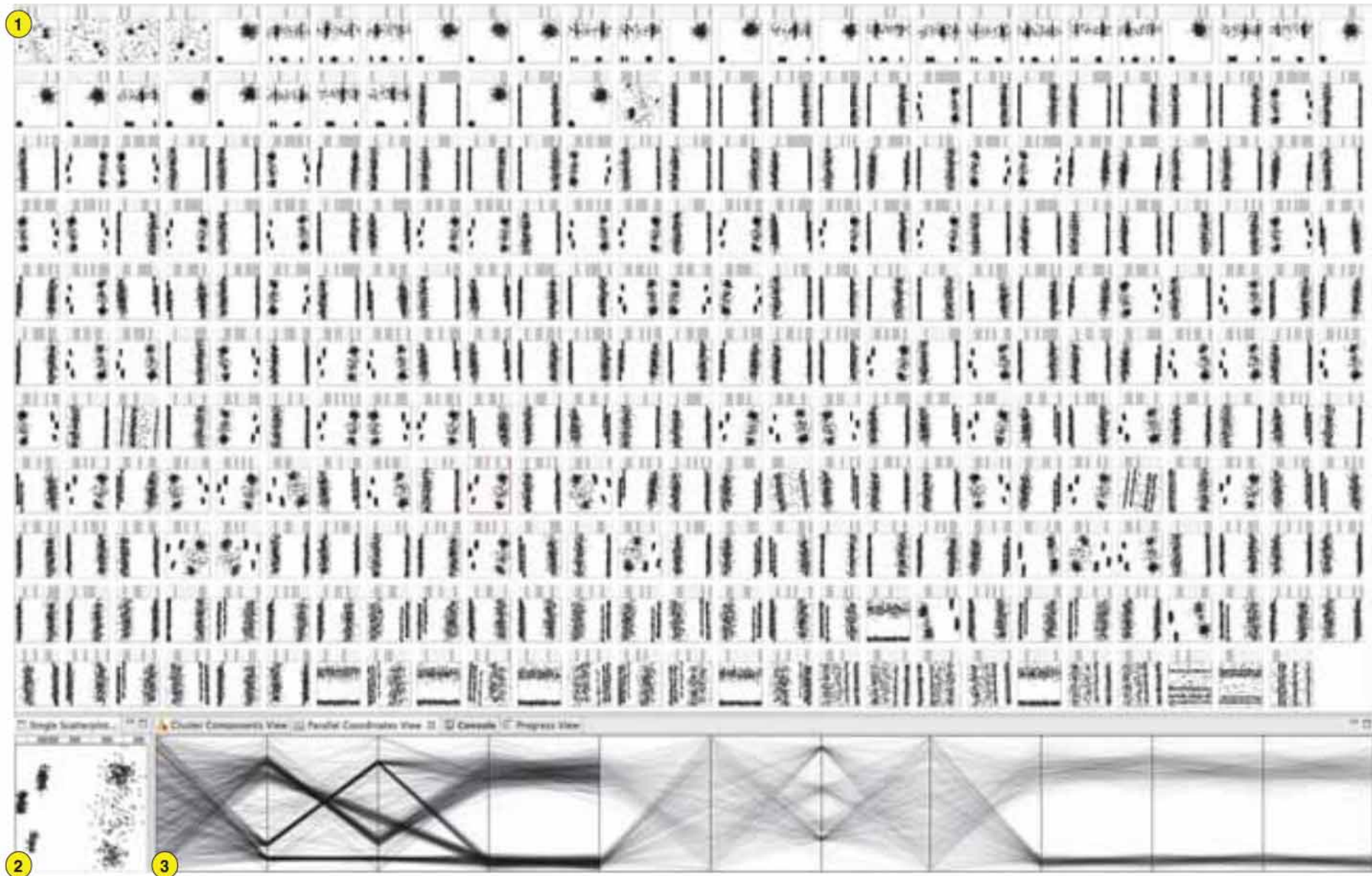
- Patterns may be found in subspaces (dimension combinations)
- Patterns may be complementary or redundant to each other



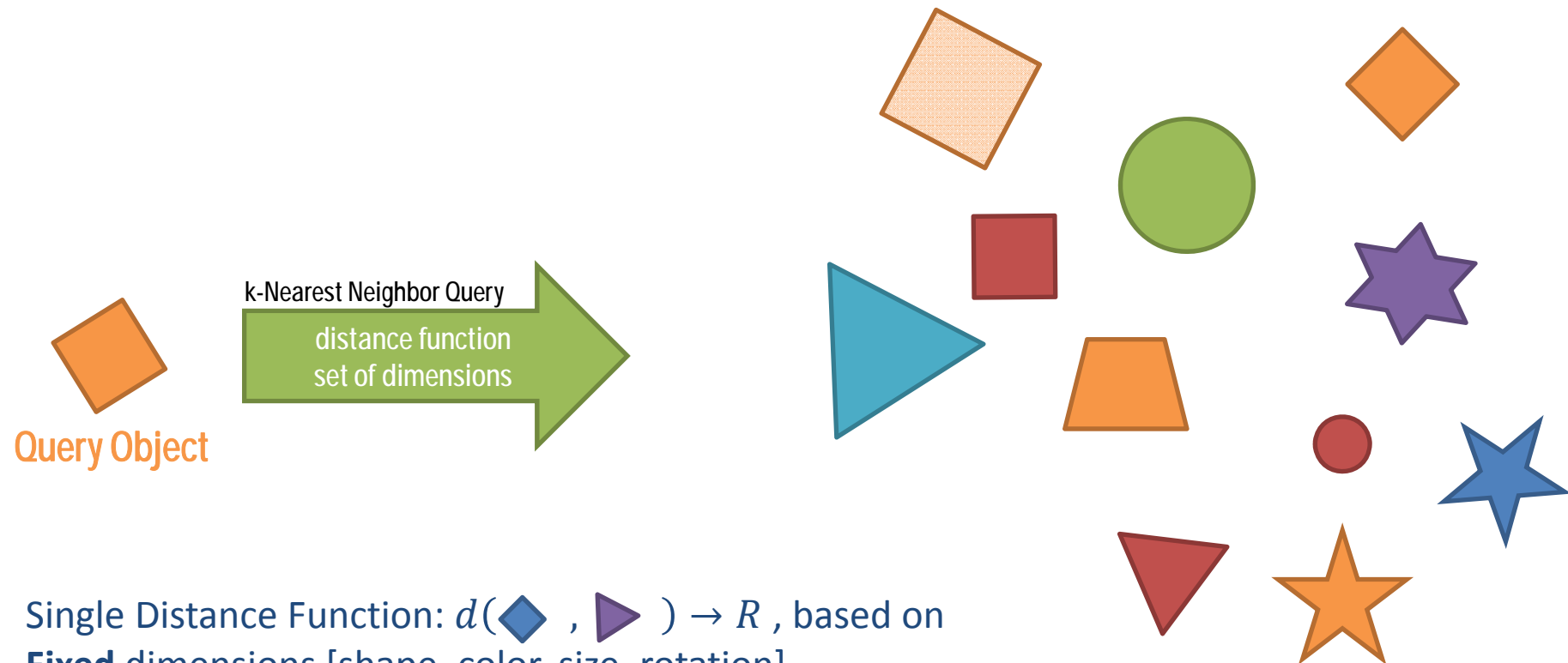
objectID	age	blood pres.	sportactiv	income	trav. freq.
1	ABC	ABC	ABC	ABC	ABC
2	ABC	ABC	ABC	ABC	ABC
3	ABC	ABC	ABC	ABC	ABC
4	ABC	ABC	ABC	ABC	ABC
5	ABC	ABC	ABC	ABC	ABC
6	ABC	ABC	ABC	ABC	ABC
7	ABC	ABC	ABC	ABC	ABC
8	ABC	ABC	ABC	ABC	ABC
9	ABC	ABC	ABC	ABC	ABC



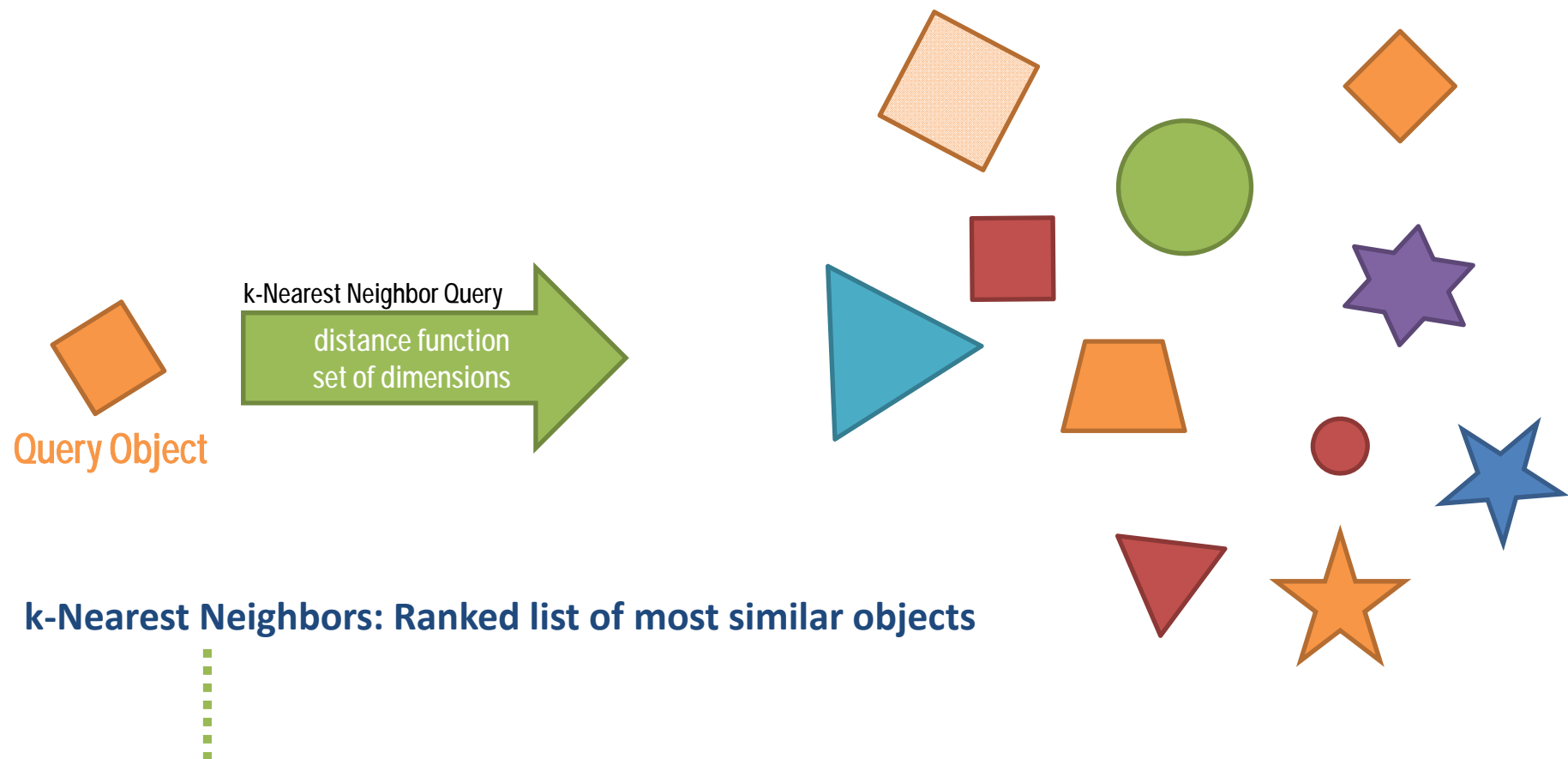
Tatu, A., Maass, F., Faerber, I., Bertini, E., Schreck, T., Seidl, T. & Keim, D. Subspace search and visualization to make sense of alternative clusterings in high-dimensional data. IEEE Symposium on Visual Analytics Science and Technology (VAST), 2012 Seattle. IEEE, 63-72, doi:10.1109/VAST.2012.6400488.





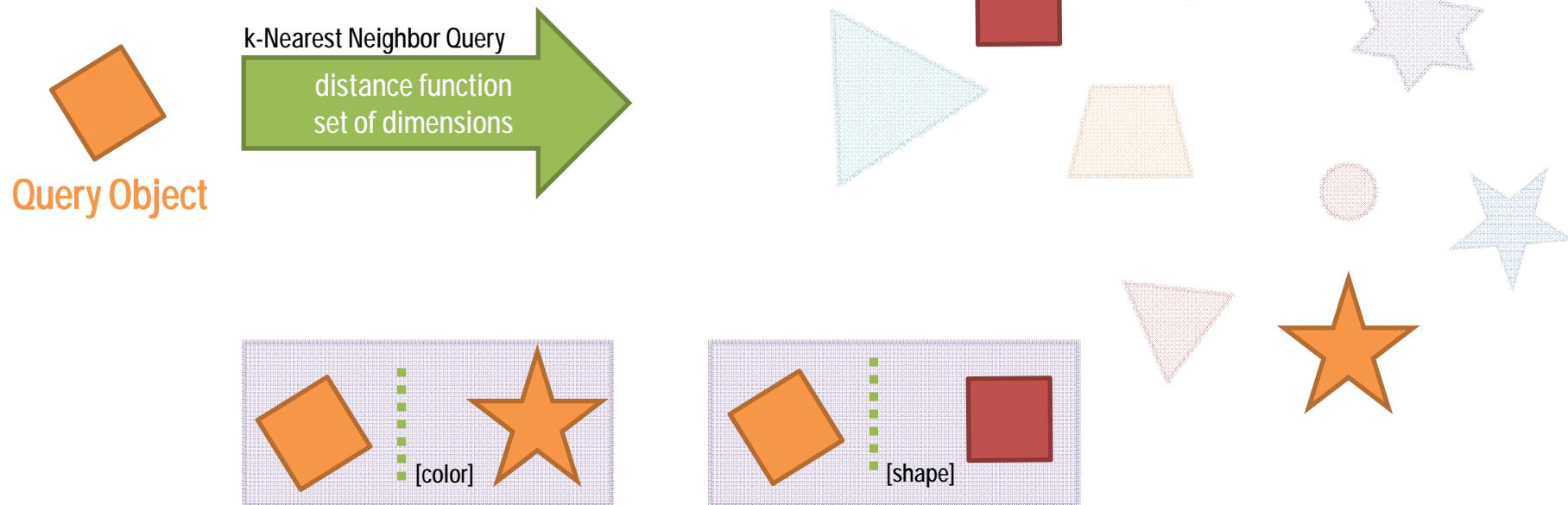


Hund, M., Behrisch, M., Färber, I., Sedlmair, M., Schreck, T., Seidl, T. & Keim, D. 2015. Subspace Nearest Neighbor Search-Problem Statement, Approaches, and Discussion. *Similarity Search and Applications*. Springer, pp. 307-313.





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





Nearest Neighbor  
Search

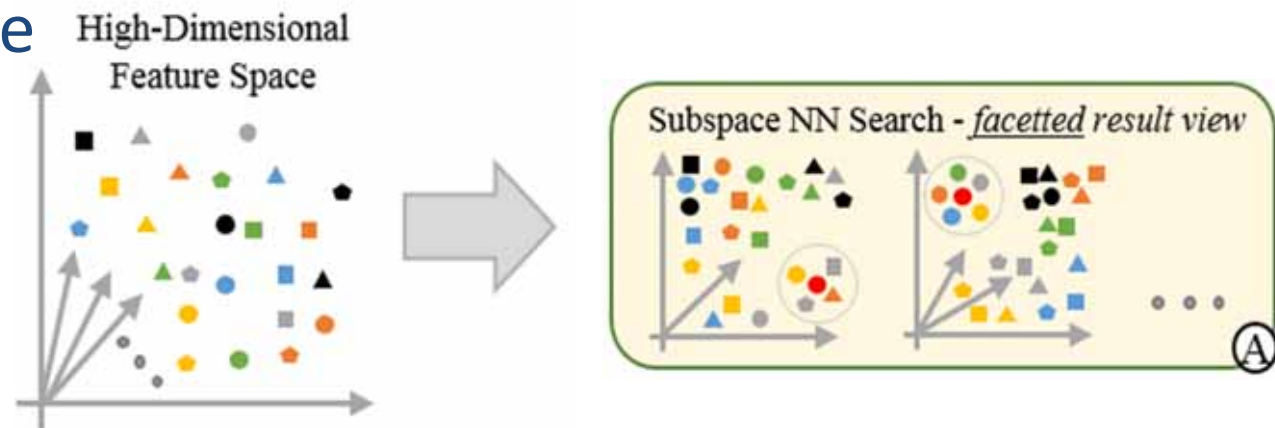


Sex, Age, Blood Type,  
Blood Pressure,  
Former Diseases,  
Medication, ...

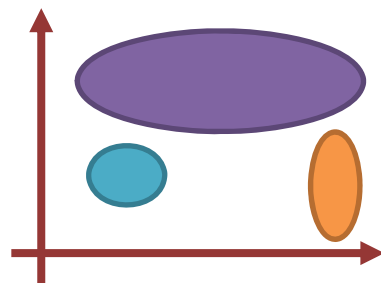
- (1) Relevant subspaces *depend on patient* and are *unknown* beforehand
- (2) *Multiple* subspaces might be relevant
- (3) Subspaces helps to *interpret* the nearest neighbors (*semantic* meaning)

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

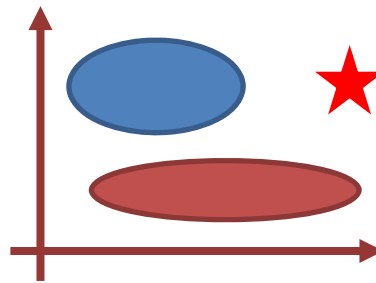
Characteristics:



- 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



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.

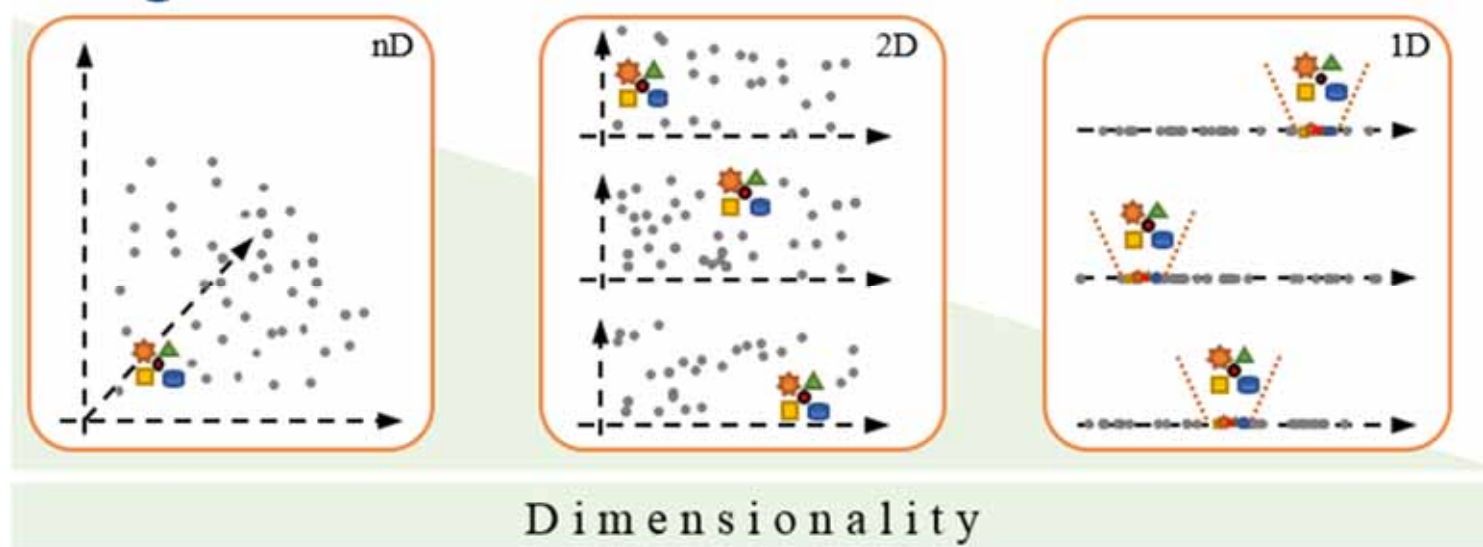
[2] Zimek, A., Schubert, E. & Kriegel, H. P. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining*, 5, (5), 363-387.

## 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

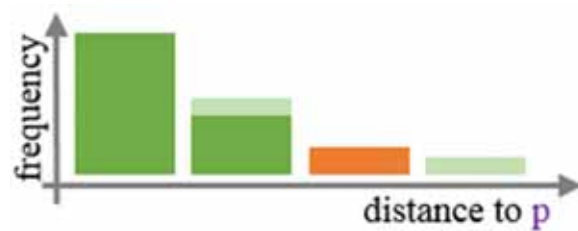
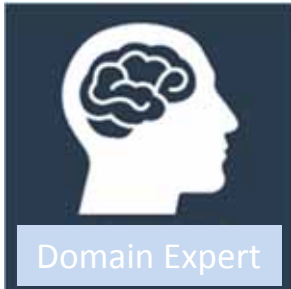
A subspace is considered **relevant**, iff it contains relevant nearest neighbors



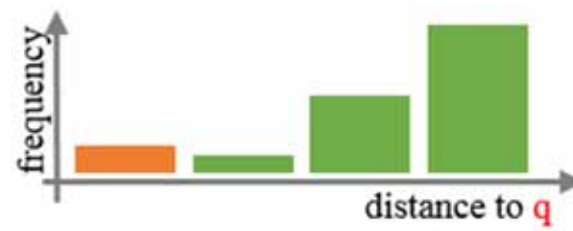
Hund, M., Behrisch, M., Färber, I., Sedlmair, M., Schreck, T., Seidl, T. & Keim, D. 2015. Subspace Nearest Neighbor Search-Problem Statement, Approaches, and Discussion. Similarity Search and Applications. Springer, pp. 307-313.

- **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

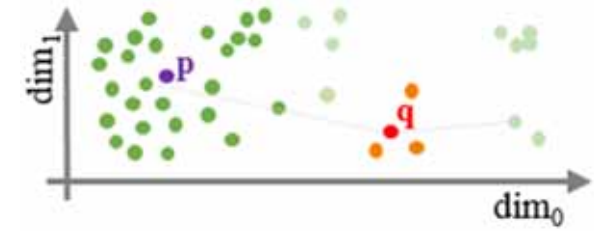




Non-Characteristic  
Dimension



Characteristic  
Dimension



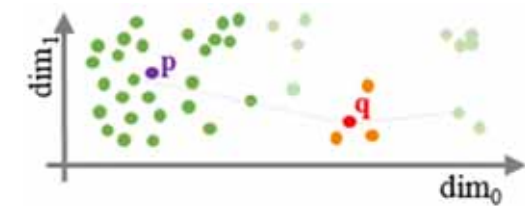
Data Distribution



query = butter



query = gauda cheese



## Supplementary Material

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

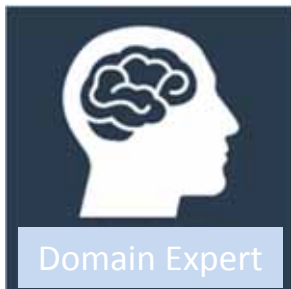
## Dataset

- 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, anhydrous	butter, without salt
butter oil, anhydrous	butter, without salt	salad drsng, mayo
kellogg's, fruit bars	lard	margarine
margarine	salad drsng, mayo	chicken, broilers
pancakes	oil, soybn	pork, backfat
waffle	oil, cocnt	candies, butterscotch
cream	oil, olive	candies, hard
cheese, cream	oil, safflower	candies, jellybeans
pie crust	vegetable oil, palm kernel	candies, mars snackfood
cheese, mozzarella	oil, canola	chewing gum
kellogg's cereals	oil, sunflower	puddings, vanilla
soup	margarine	jellies
cheese, limburger	shortening	sweeteners, tabletop
peppers	chicken, broilers	syrups, corn
sauce tabasco	oil, corn, peanut, and olive	syrups, maple

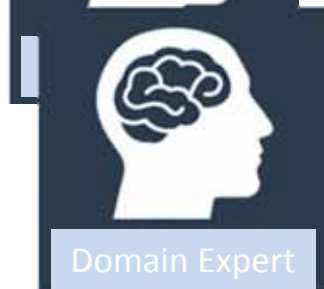


(1) Determine Nearest Neighbors per Dimension

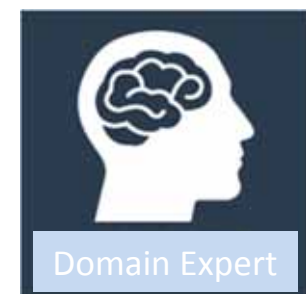
(2) Efficient Search Strategy

(3) Query-Based Interestingness for Dimensions

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

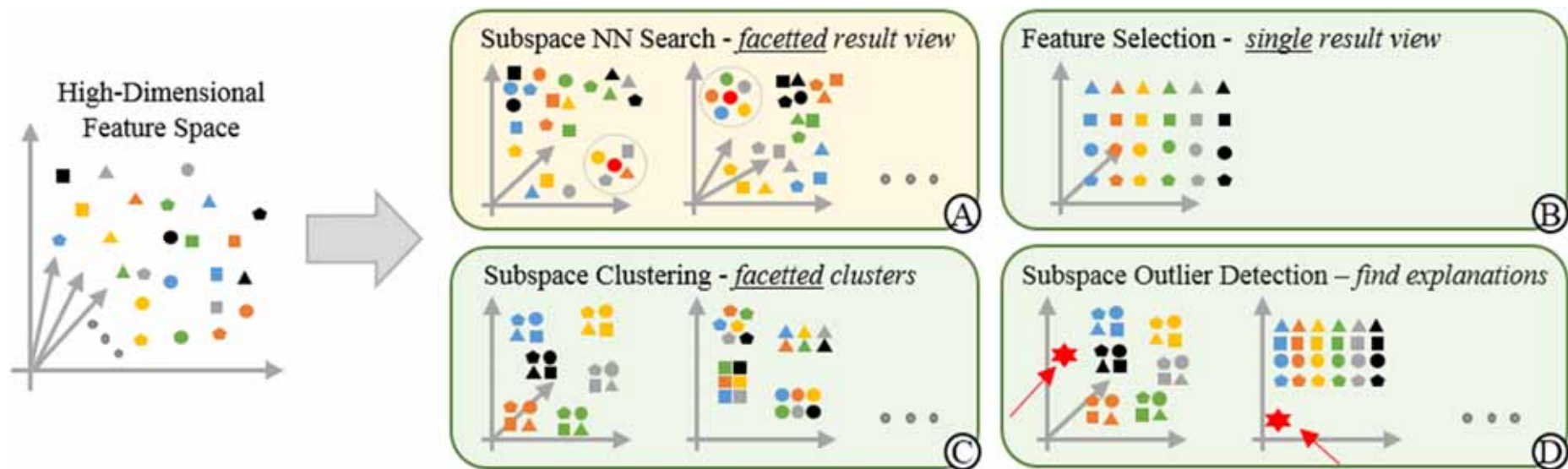


(5) Evaluation Methods and Development of Benchmark Datasets



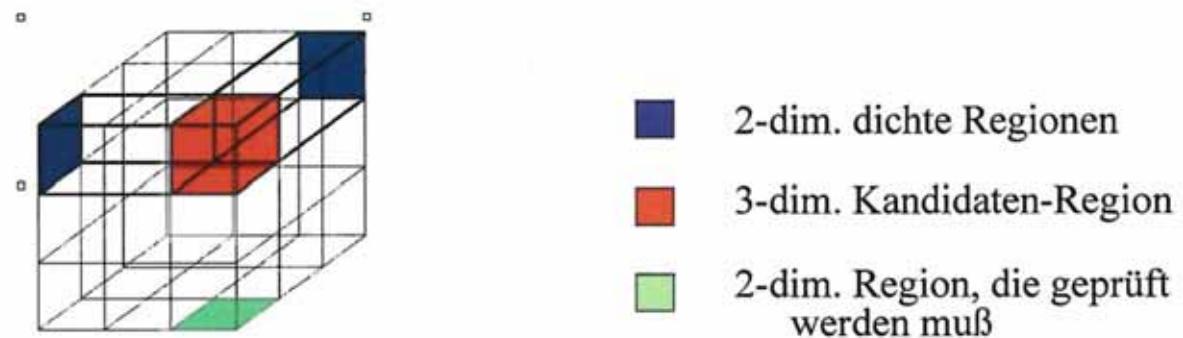
(6) Multi-input Subspace Nearest Neighbor Search

(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.

- 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)



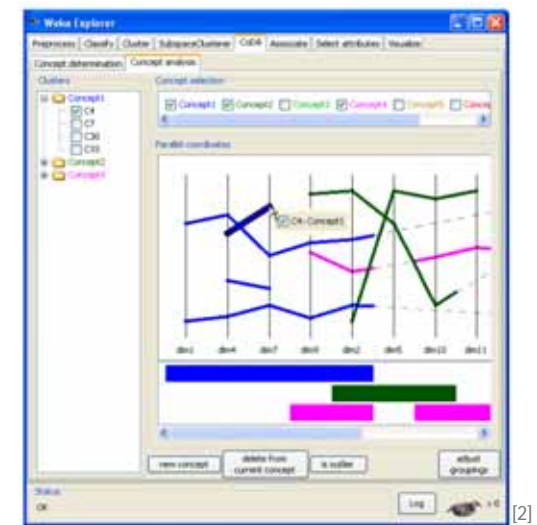
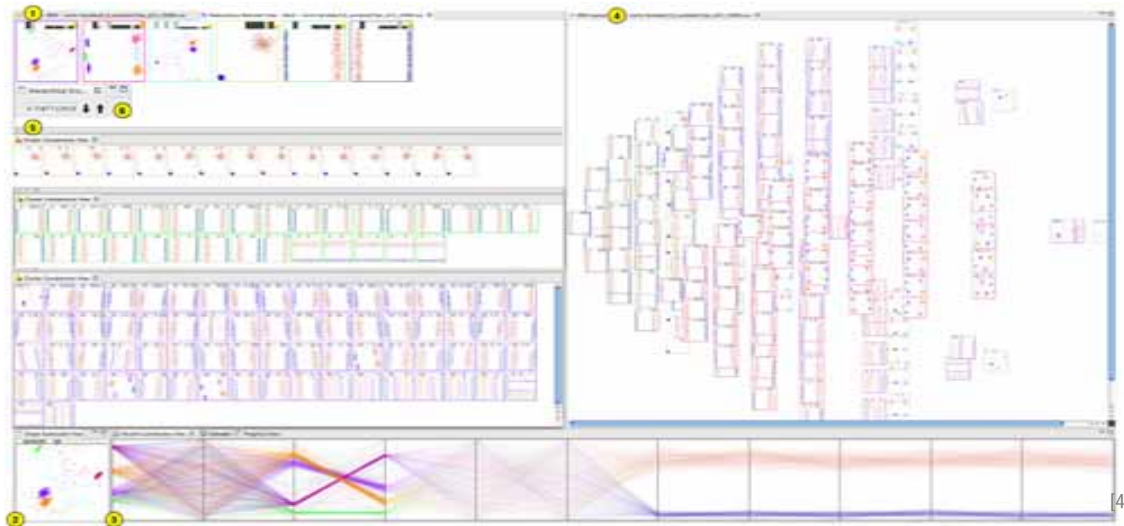
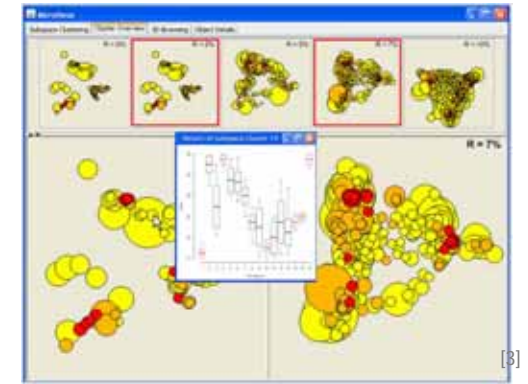
Which dimensions occur more often in clusters?  
Which occur often together?  
Which values do records in a specific cluster have?



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.

- 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



- 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

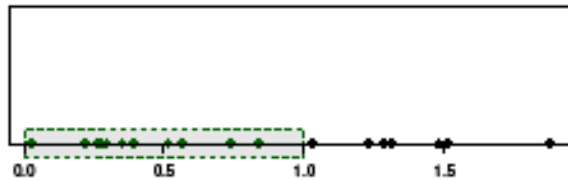
Domain Expert

- Pruning of intermediate results
- Adjust parameters to domain knowledge
- ...

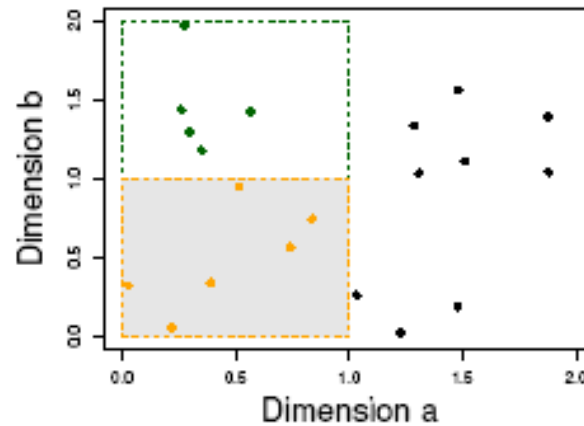


**Fig. 3** A screenshot of our visual analytics tool SubVIS. It enables the user to interactively explore a large number of subspace clusters. A general overview of the similarities between the subspaces is given by an MDS projection (A). Small multiples (B) allows to preview projections of different 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 lense (F) provides details on demand.

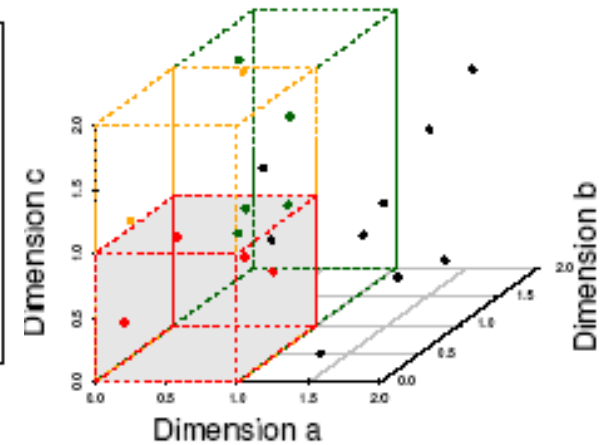
Hund, M., Sturm, W., Schreck, T., Ullrich, T., Keim, D., Majnaric, L. & Holzinger, A. 2015. Analysis of Patient Groups and Immunization Results Based on Subspace Clustering. In: Guo, Y., Friston, K., Aldo, F., Hill, S. & Peng, H. (eds.) Brain Informatics and Health, Lecture Notes in Artificial Intelligence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



(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

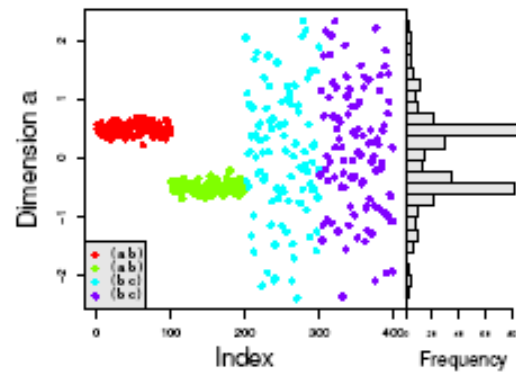
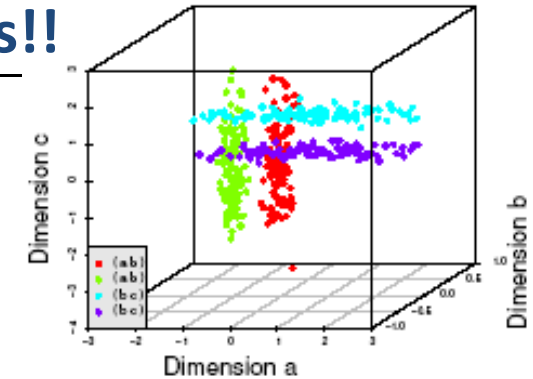
Distance measure becomes meaningless—due to equidistance

- 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.

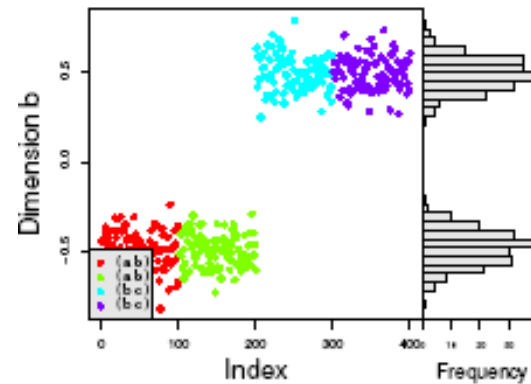


# Interesting Clusters may ONLY exist in subspaces!!

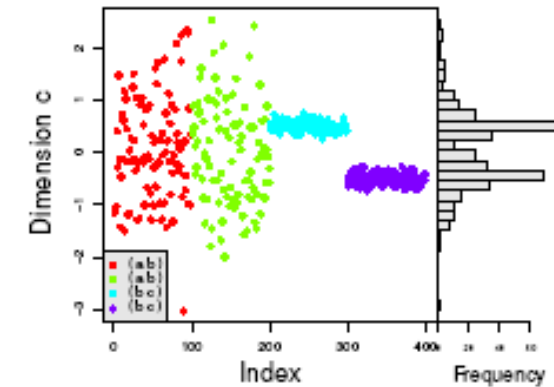
Parsons, L., Haque, E. & Liu, H. 2004. Subspace clustering for high dimensional data: a review. SIGKDD Explorations 6, (1), 90-105.



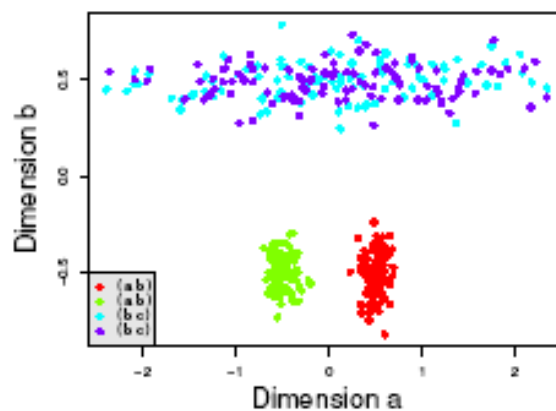
(a) Dimension *a*



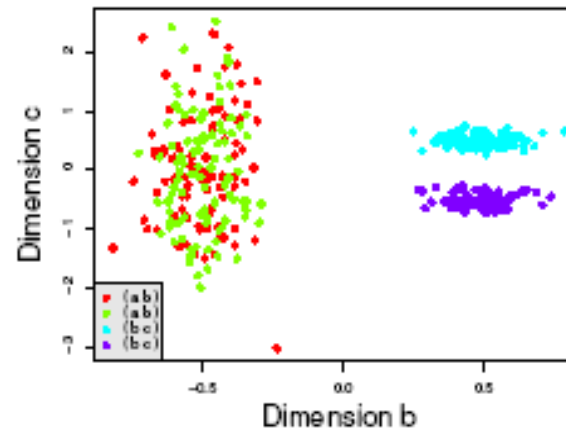
(b) Dimension *b*



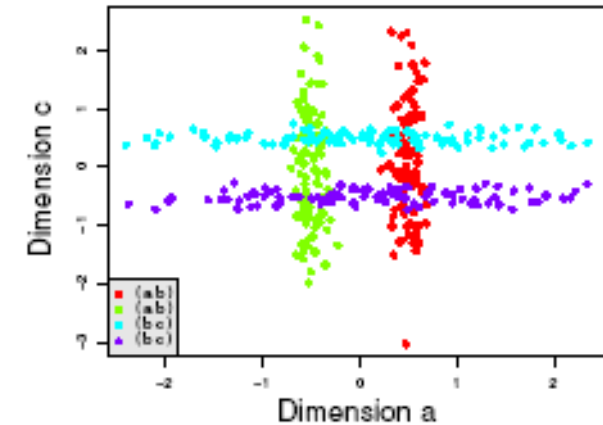
(c) Dimension *c*



(a) Dims *a* & *b*

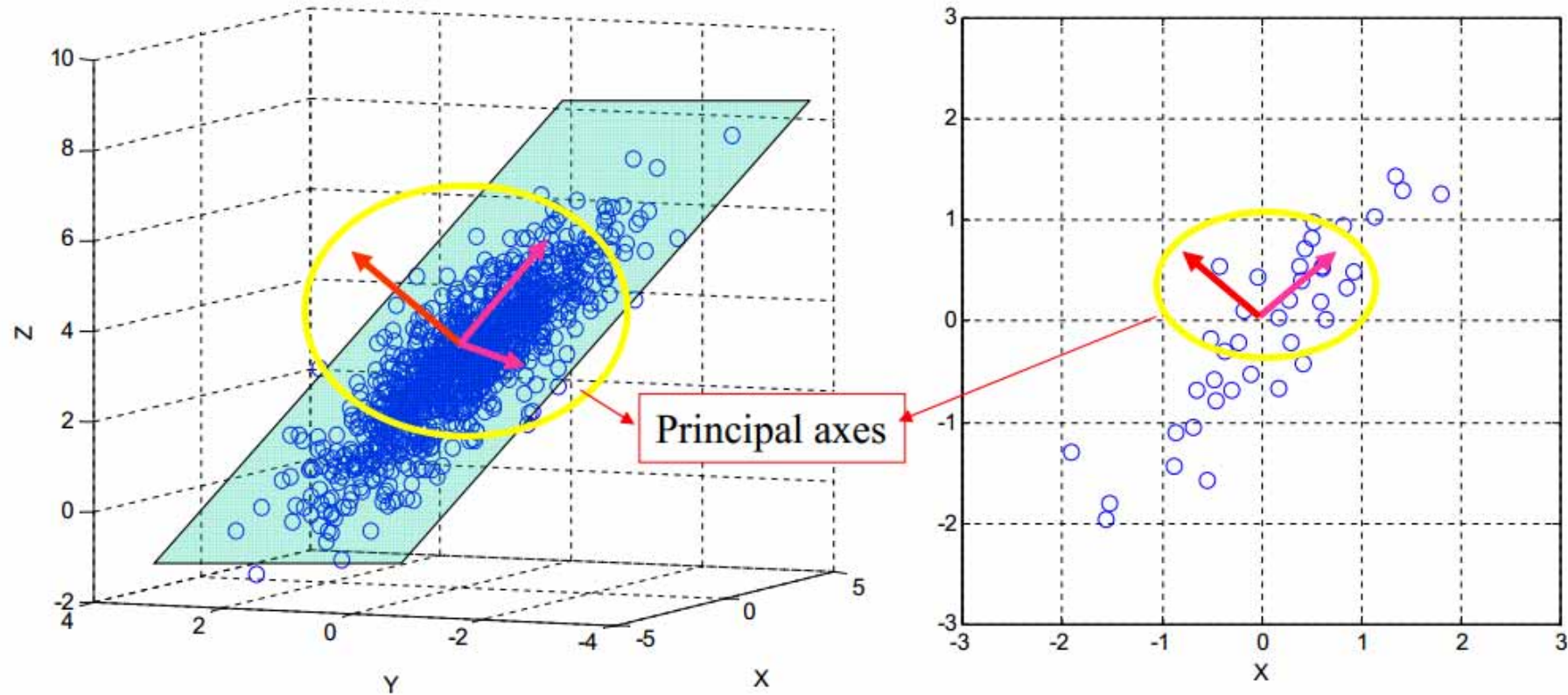


(b) Dims *b* & *c*



(c) Dims *a* & *c*





# 6) “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.

## Projection pursuit:

least Gaussian (“interesting”) projections of the data

how to define non-Gaussianity?

covariance and mean given: Gaussian distribution maximizes the entropy

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

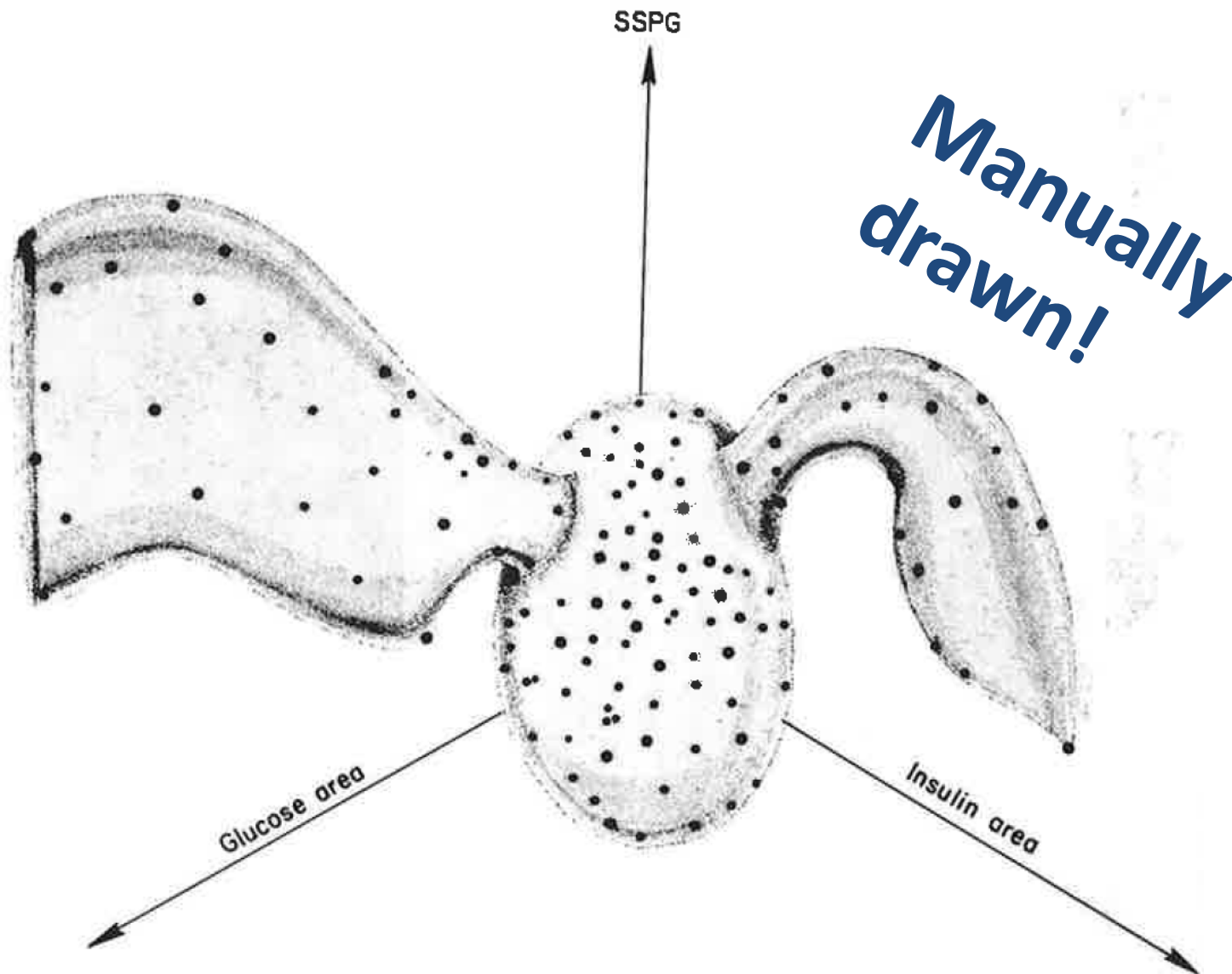
This is difficult to optimize

- finding unimodal super-Gaussians
- finding multimodal distributions

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

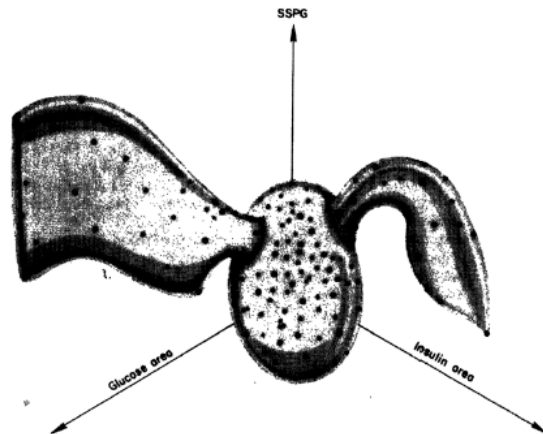
- 145 diabetes patients
- 6 dimensional data set:
  - 1) age,
  - 2) relative weight,
  - 3) fasting plasma glucose,
  - 4) area under the plasma glucose curve for the three hour glucose tolerance test (OGTT),
  - 5) area under the plasma insulin curve for the OGTT,
  - 6) steady state plasma glucose 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.



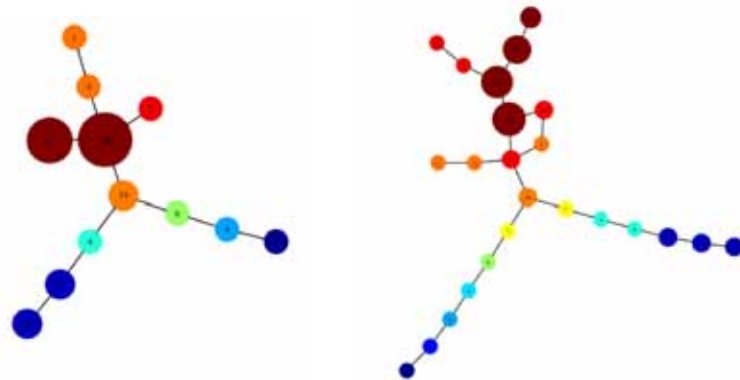


Given a point cloud data set  $X$  and a covering  $U$   
 $\Rightarrow$  *simplicial complex*

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

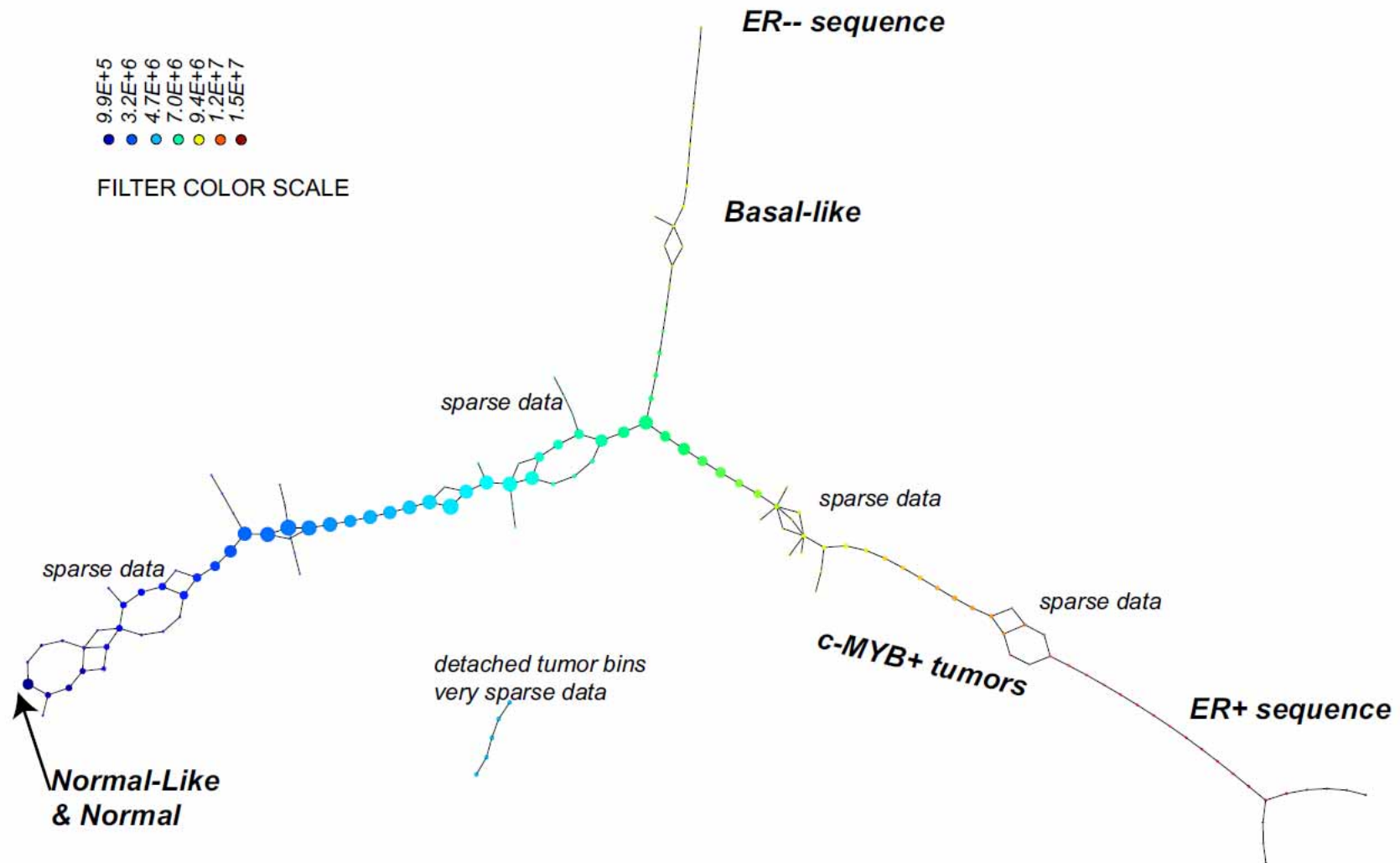
$$f: X \rightarrow Z$$

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



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

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.

- 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?



# Thank you!

- 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?