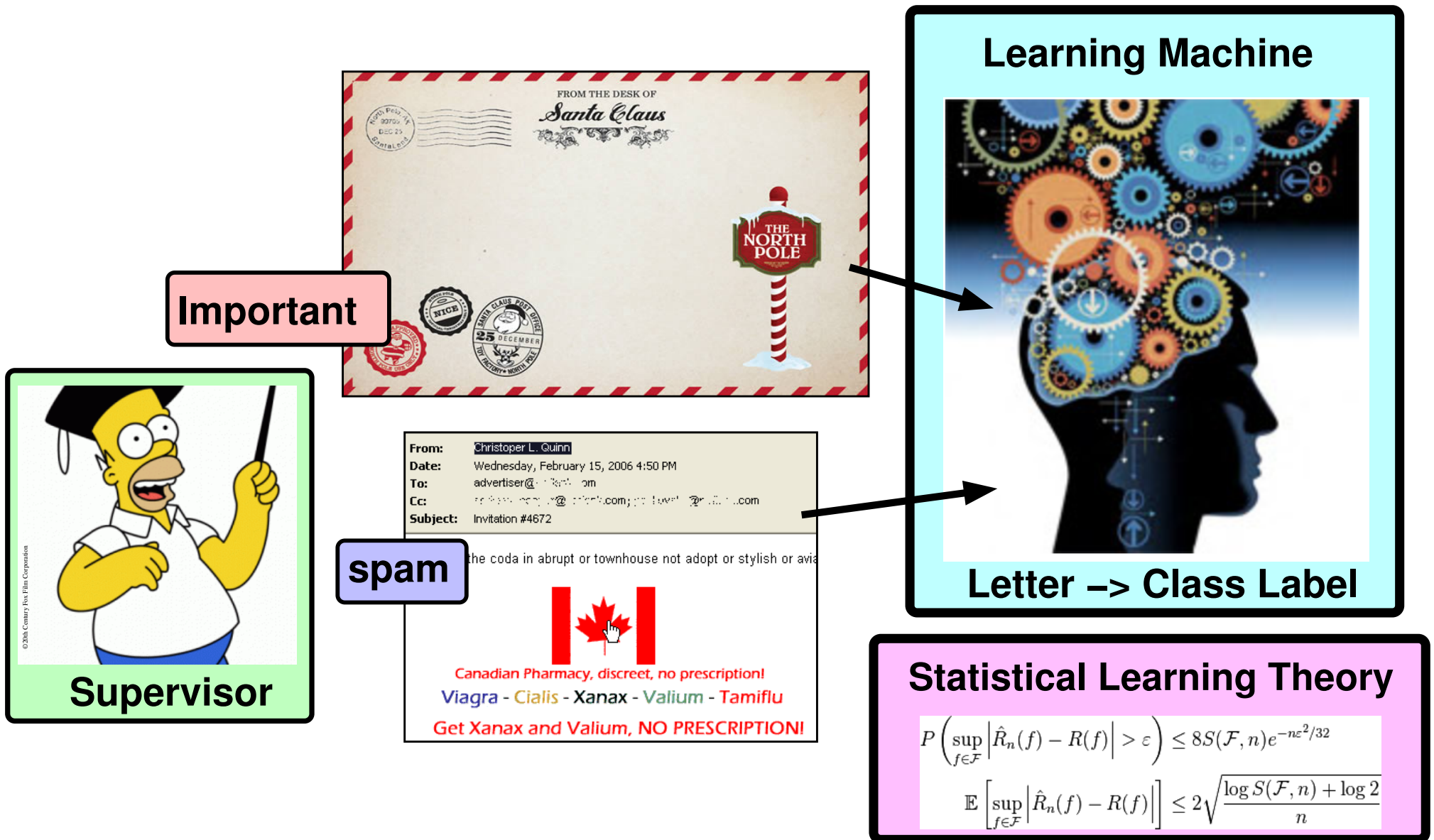


# Archetype Analysis: a framework for selecting representative objects

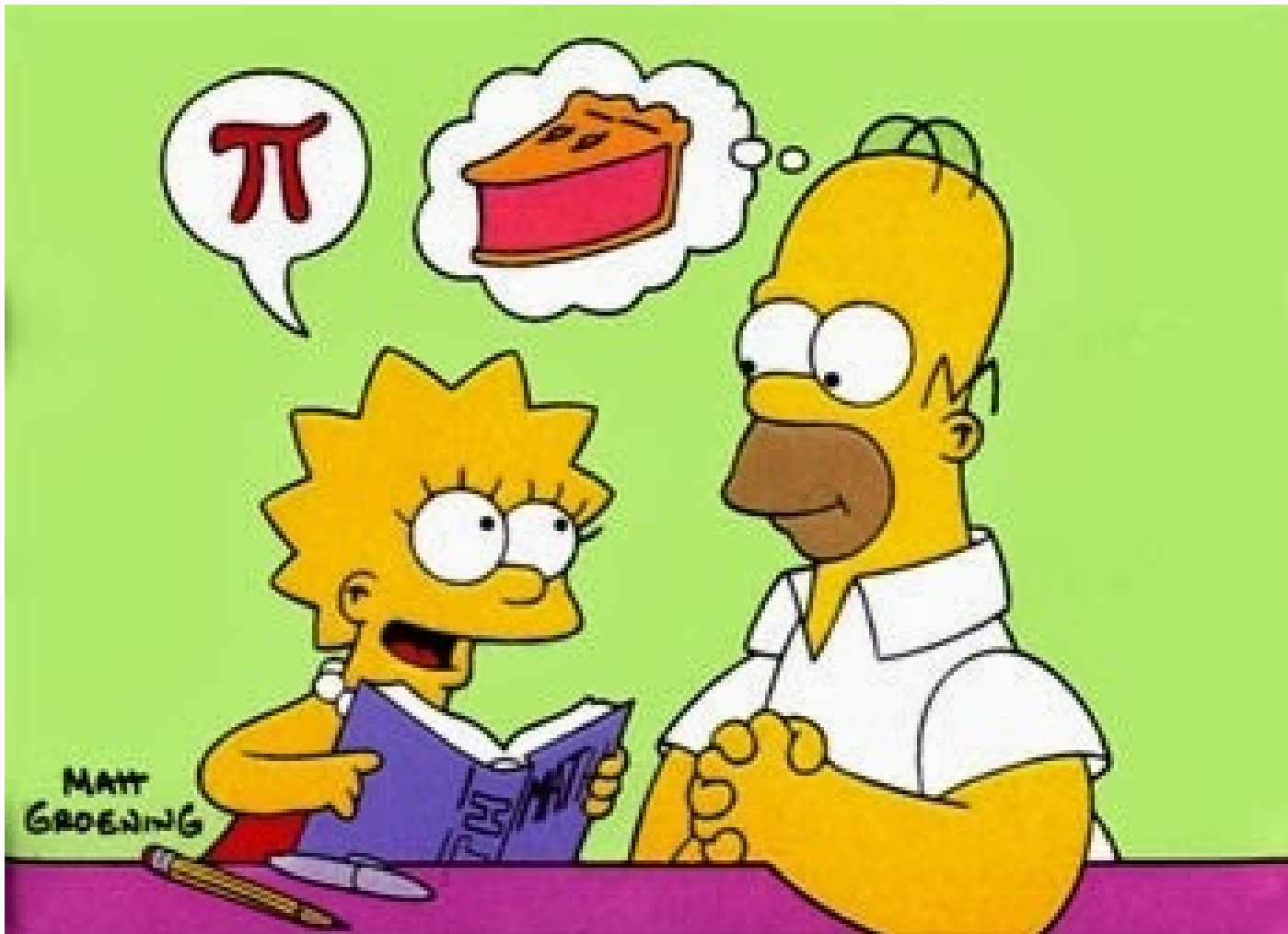
Volker Roth, Department of Mathematics and Computer Science,  
University of Basel

# Machine Learning: Supervised Setting

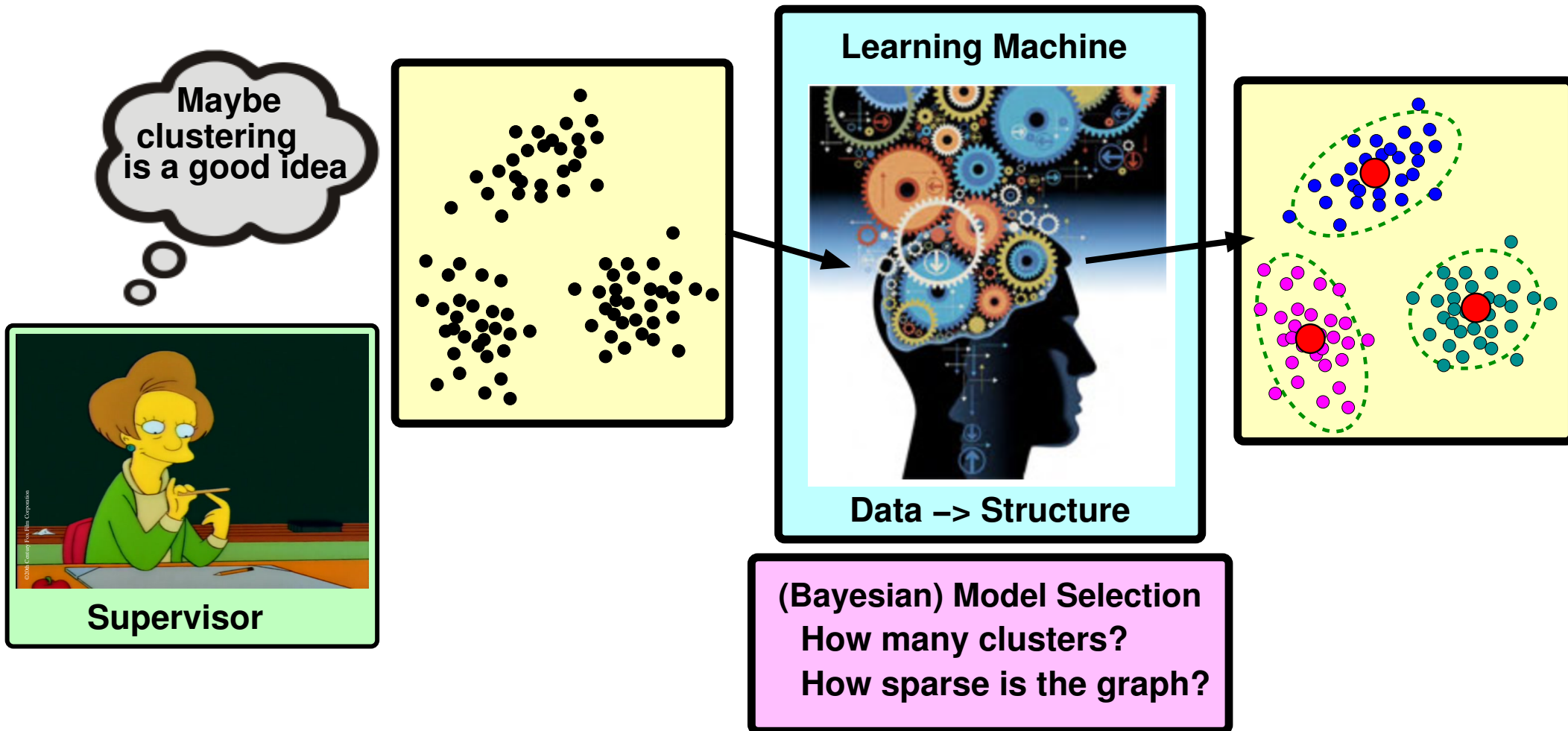


# Machine Learning: Uncertain Labels

Sometimes even the best experts make errors.  
Labels might be uncertain or missing...

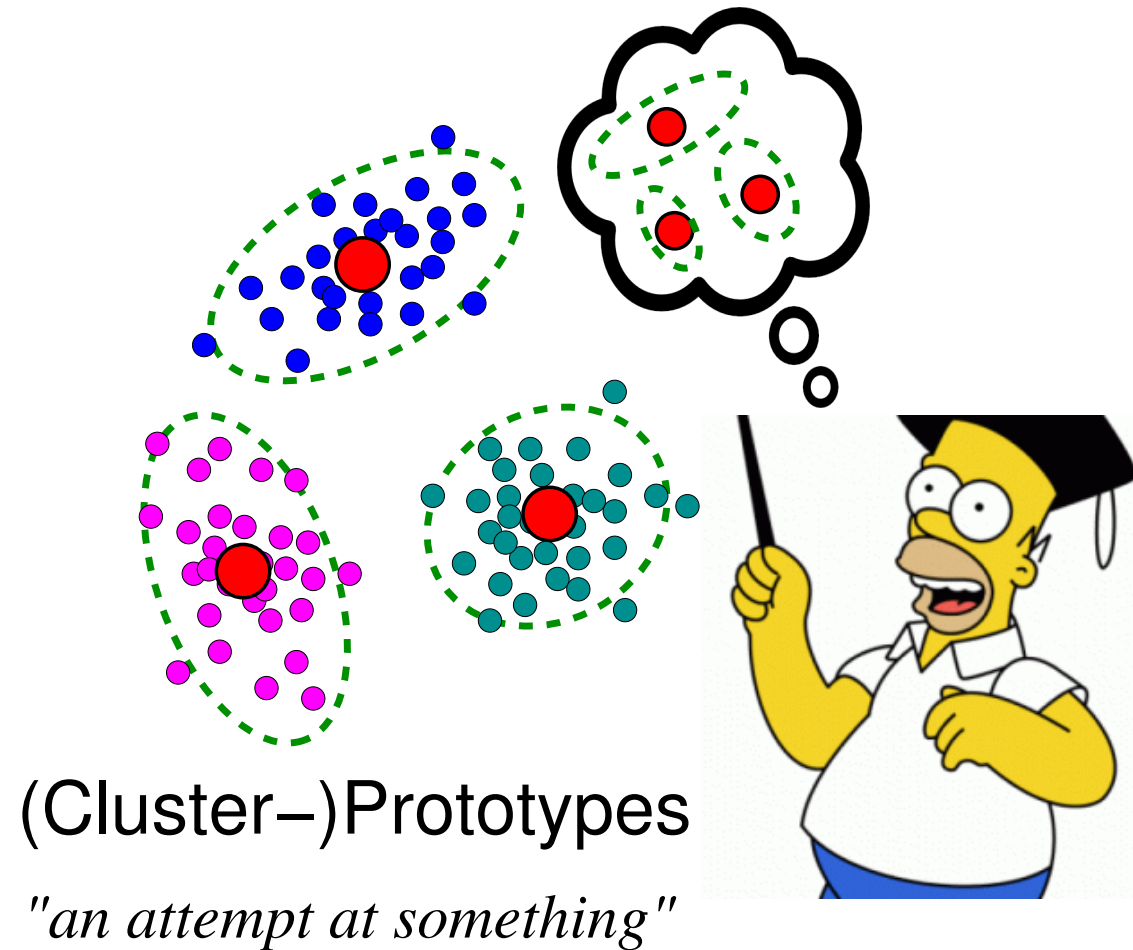
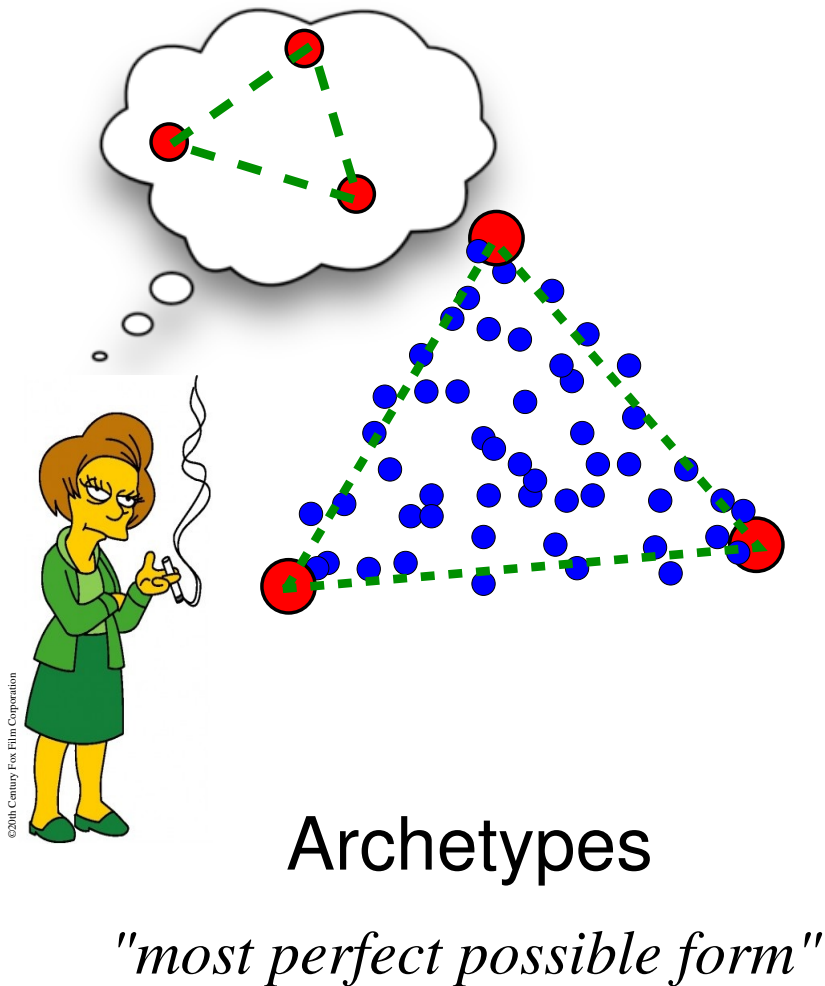


# Machine Learning: Unsupervised Setting





# A Repeated Pattern: Search for Representative Observations

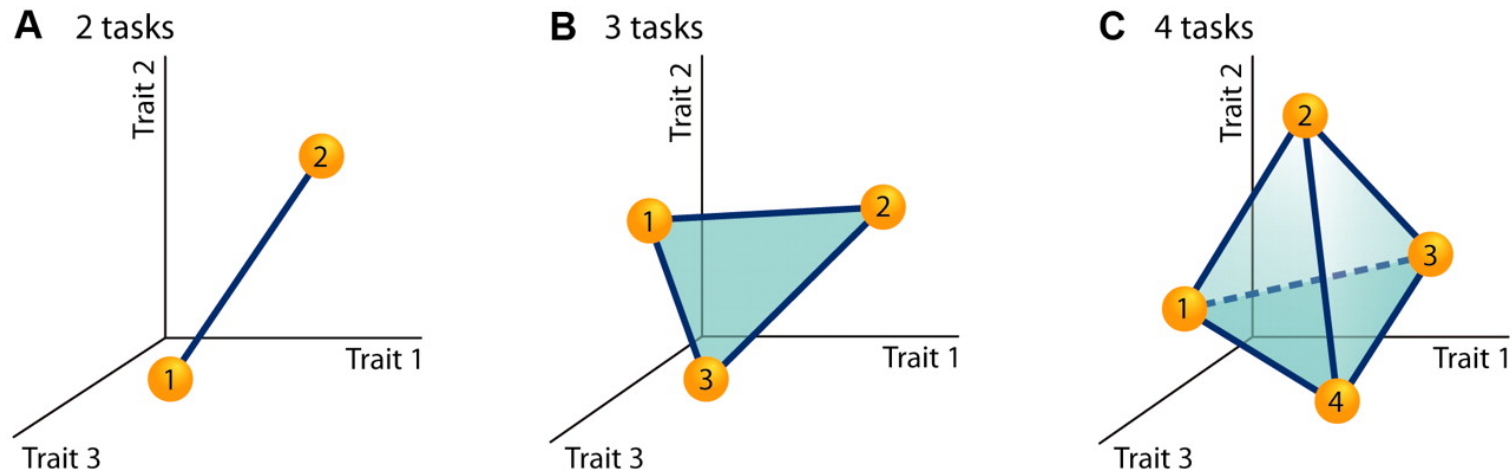
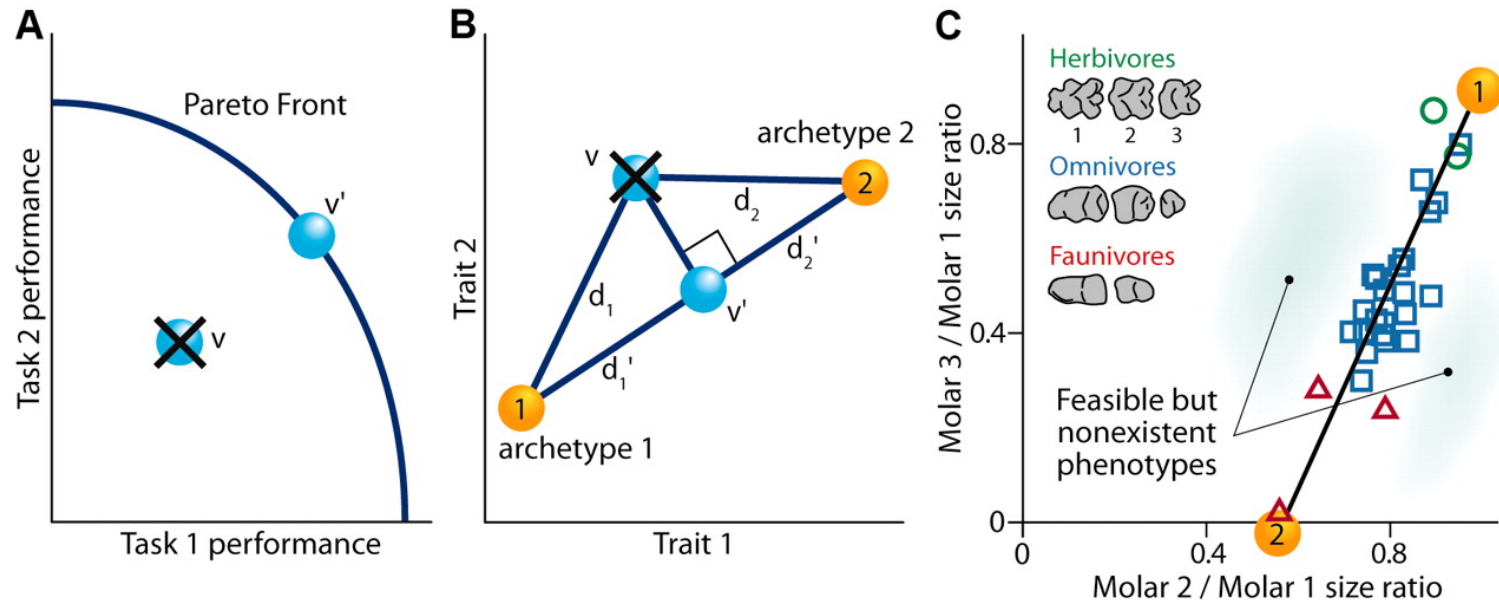


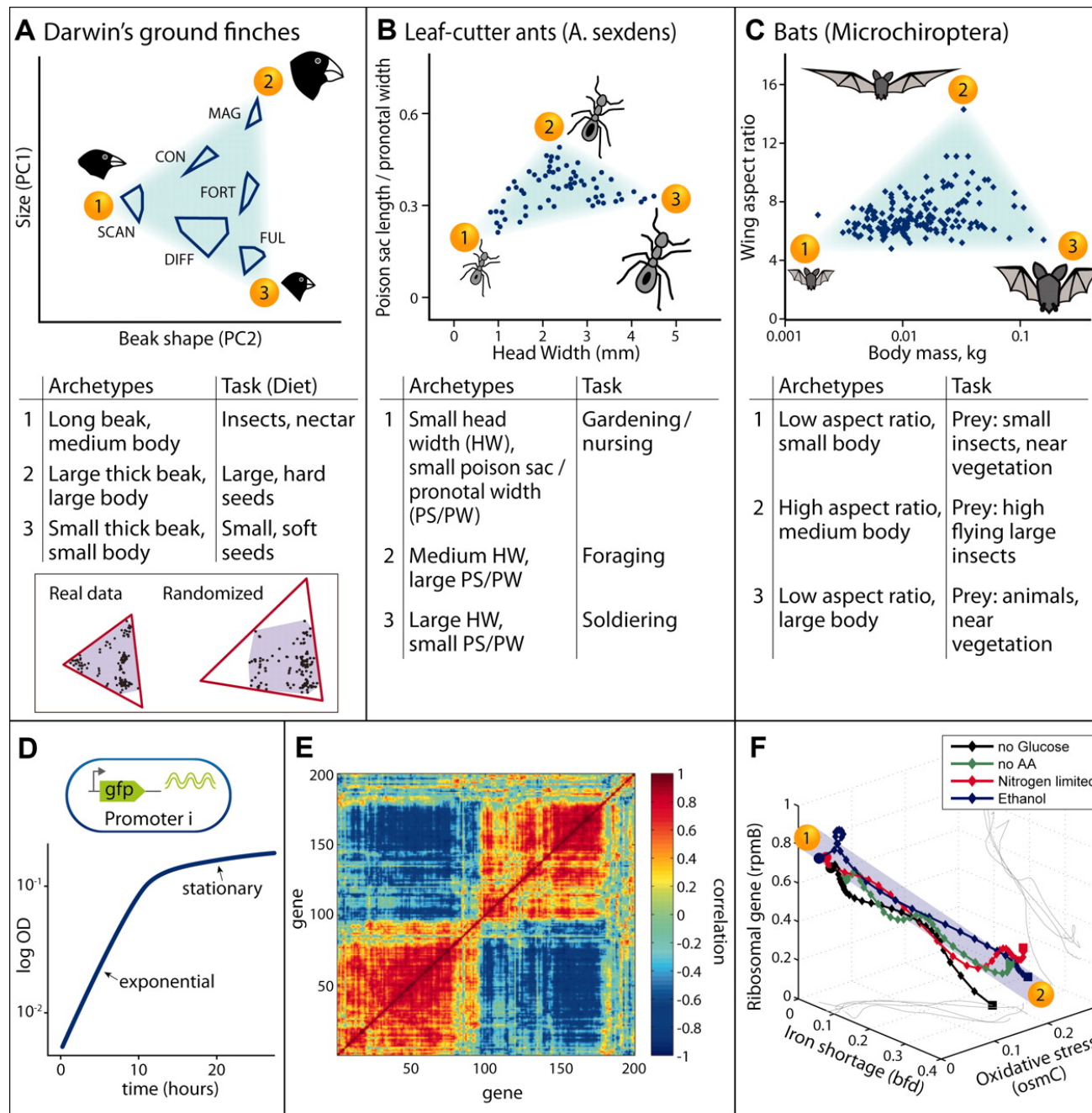
# Archetype Analysis: Biological Motivation

Is there a theoretical foundation of the “archetype concept”?

~> O. Shoval, H. Sheftel, G. Shinar, Y. Hart, O. Ramote, A. Mayo, E. Dekel, K. Kavanagh, U. Alon: **Evolutionary Trade-Offs, Pareto Optimality, and the Geometry of Phenotype Space**. Science, 2012

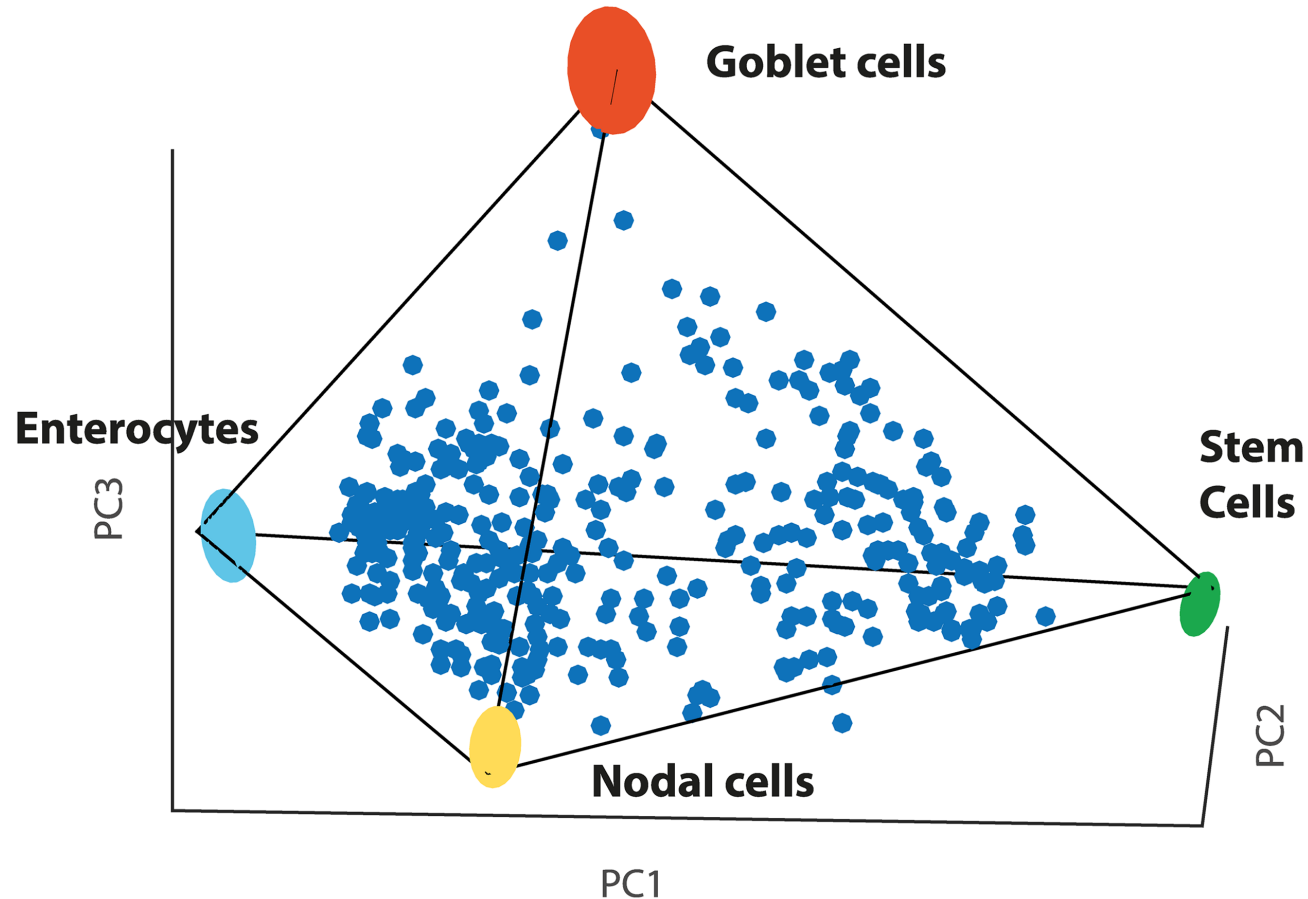
# Archetypes and Evolutionary Trade-offs





# Gene Expression Space

Human colon crypt cells fall in a tetrahedron in gene expression space.



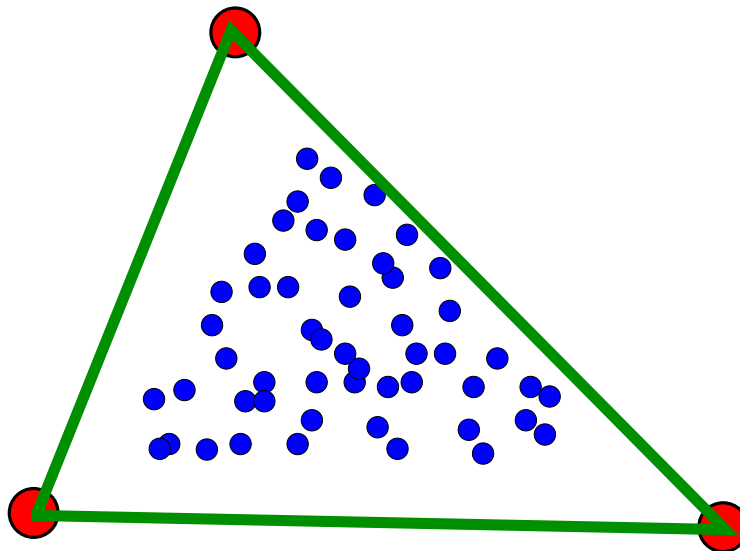
The four vertices of this tetrahedron are each enriched with genes for a specific task related to stemness and early differentiation.

# Computational Archetype Selection

Cutler & Breiman, *Archetypal Analysis*, Technometrics 1994.

- $n$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^p$ , as rows of data matrix  $X \in \mathbb{R}^{n \times p}$
- Aim: find  $K$  archetypes  $\Rightarrow Z \in \mathbb{R}^{K \times p}$ ;  $K \ll n$  fixed.
- Observations are convex mixtures of archetypes

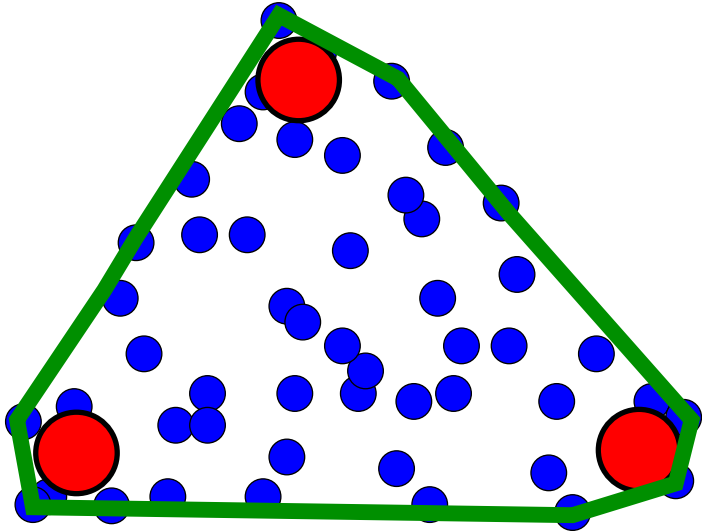
$$\mathbf{x}_i = Z^t \mathbf{a}_i + \epsilon_i, \quad a_{ij} \geq 0 \quad \text{and} \quad \sum_{j=1}^K a_{ij} = 1.$$



# Computational Archetype Selection

- Archetypes themselves are convex mixtures of observations:

$$\mathbf{z}_i = \sum_{j=1}^n b_{ij} \mathbf{x}_j, \quad \text{where } b_{ij} \geq 0 \quad \text{and} \quad \sum_{j=1}^n b_{ij} = 1$$



**Archetypes approximate the convex hull!**

- Constrained optimization problem involving two sets of coefficients  $\{a_{ij}\}$  and  $\{b_{ij}\}$  : iteratively minimize sum of squares

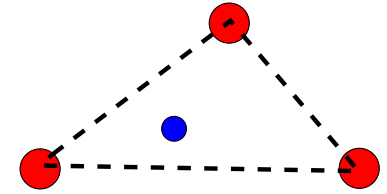
$$(\hat{A}, \hat{B}) = \operatorname{argmin}_{A, B} \|X - AZ\|^2 = \|X - ABX\|^2.$$



# Basic Algorithm

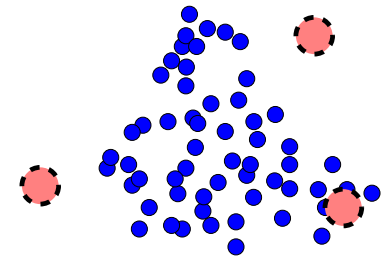
1. Given  $K$  archetypes  $Z_{K \times p}$ , update the compositions  $\mathbf{a}_i$  of the  $i$ -th object (a QP):

$$\mathbf{a}_i = \arg \min_{\mathbf{a} \in \mathbb{R}_+^p : \mathbf{a}^t \mathbf{1} = 1} \|\mathbf{x}_i - Z^t \mathbf{a}\|^2.$$



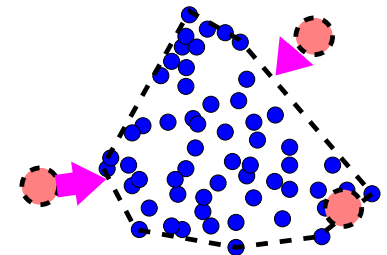
2. Given compositions  $A_{n \times K}$ , update archetypes by solving the least-squares problem

$$Z = \arg \min_{Z \in \mathbb{R}^{K \times p}} \|X - AZ\|^2.$$



3. Move ATs back to the convex hull (also a QP):

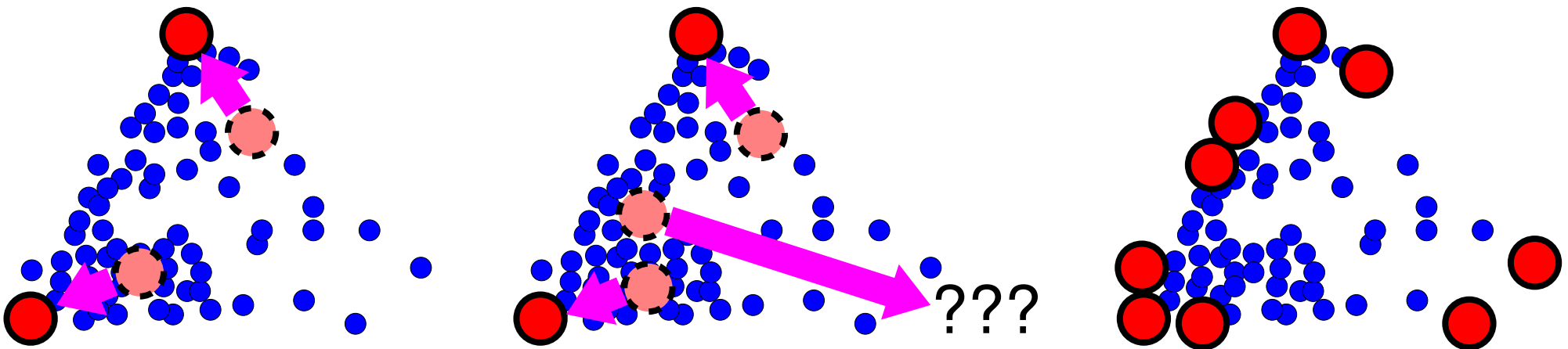
$$\mathbf{b}_k = \arg \min_{\mathbf{b} \in \mathbb{R}_+^n : \mathbf{b}^t \mathbf{1}_n = 1} \|\mathbf{z}_k - X^t \mathbf{b}\|^2.$$



# Computational Archetype Selection

Problems:

- High computational complexity,
- Solution heavily depends on initialization of archetypes,
- Have to fix the number of archetypes  $K$  a priori.
- Assume that we have a suitable **representation** such that it makes sense to search for “triangles” ...

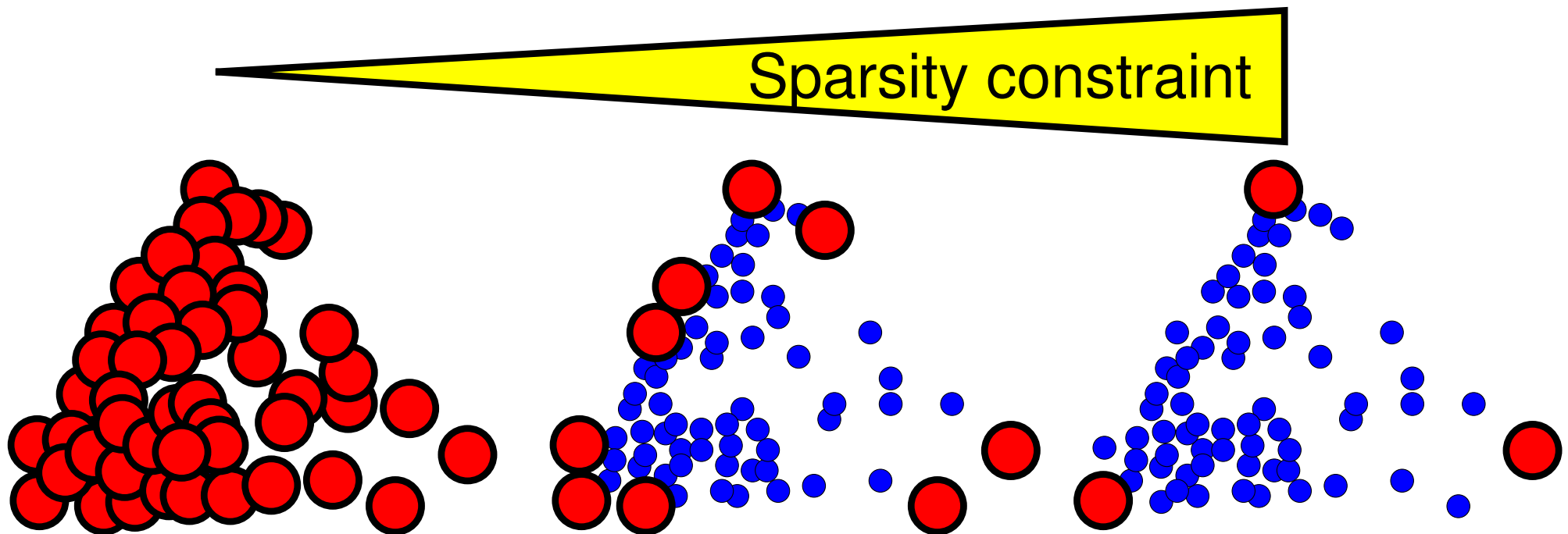


# Archetype Analysis: Model Selection

**How many archetypes?**

# Model selection: how many archetypes?

- Solve for all numbers and choose “best”  
     $\leadsto$  **Bayesian model comparison.**
- A clever way for looking at all numbers?
- Initialize all observations as archetypes, apply **sparse regression** methods to shrink most archetypes to zero.

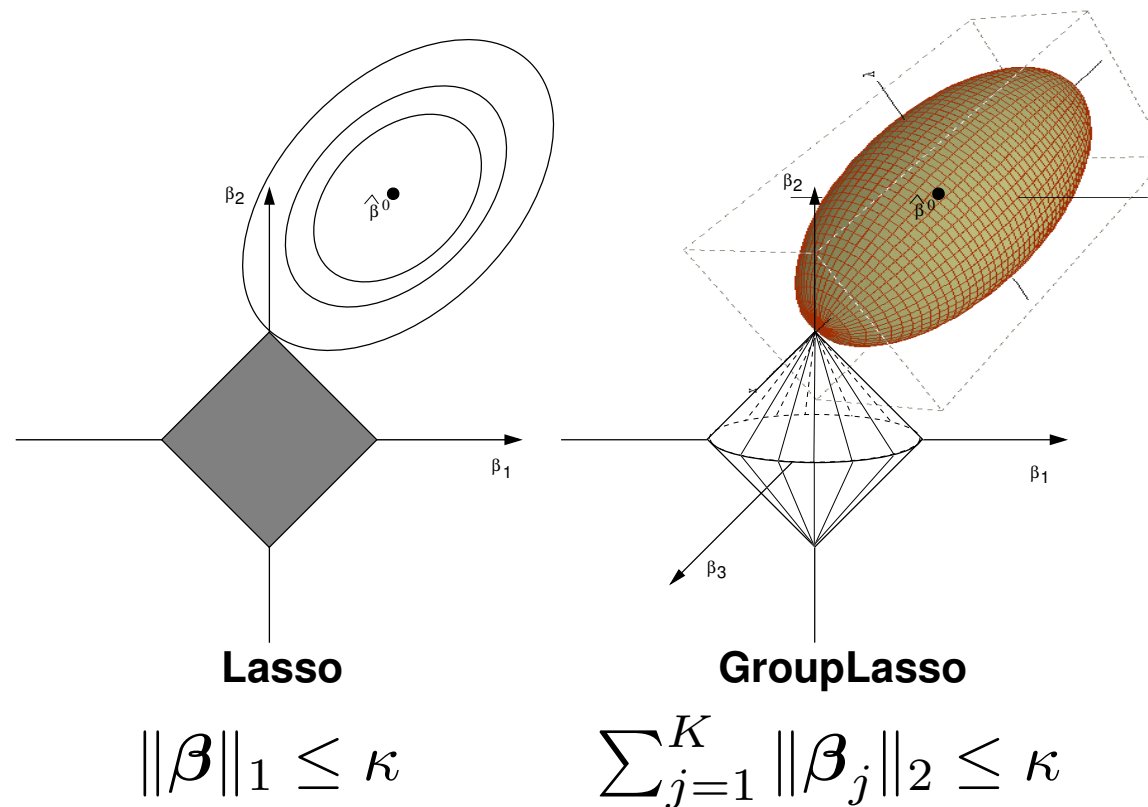


# Sparse Regression via the Group-Lasso

Least-squares:  $\min_{\beta} \|X\beta - \mathbf{y}\|^2$ .

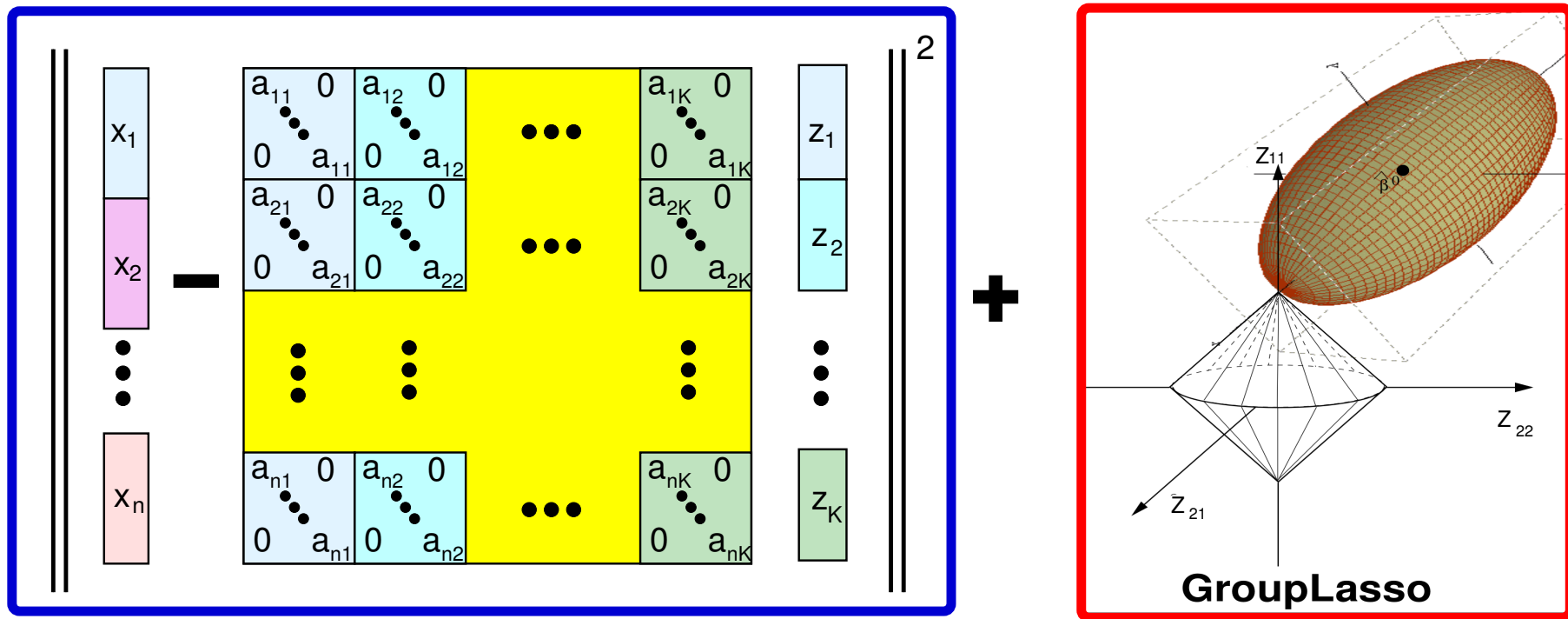
Treat regression coefficients  $\beta_i$  as **resources** needed for optimization.

Idea: **limit resources**  $\leadsto$  model must concentrate on **important features**.



# Automatic detection of archetypes

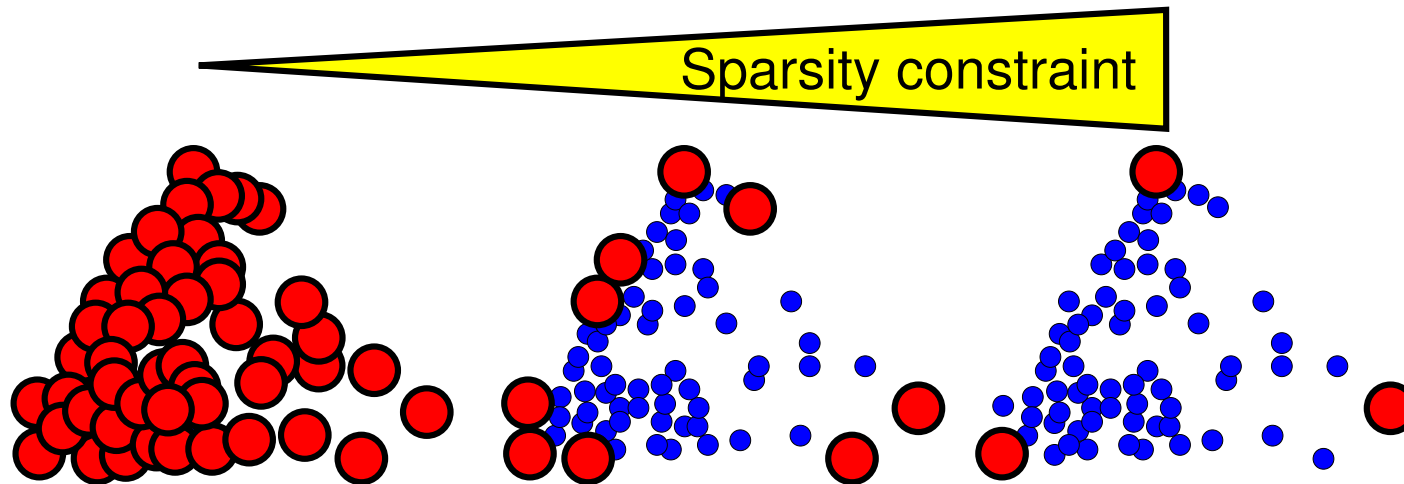
Rewrite  $\|X - AZ\|$  as  $\|\vec{x} - (A \otimes I_p)\vec{z}\|$ , where  $\vec{x}$  means column-wise vectorization of  $X$  and use  $\ell_{1,2}$  **block-norm constraint**:  $\sum_{j=1}^K \|\mathbf{z}_j\| \leq \kappa$ .



GroupLasso constraint will shrink some archetypes  $z_j$  to zero, depending on constraint value  $\kappa$ .

# A solution for the model selection problem

- Use GL algorithm that approximates the whole solution path on a fine grid of  $\kappa$  values



- For every  $\kappa$ -value, use BIC score

$$\text{BIC}(\hat{\mu}) = \frac{\|\vec{x} - \hat{\mu}\|^2}{n\sigma^2} + \frac{\log(n)}{n} \cdot \hat{df}(\hat{\mu}),$$

and select the best scoring model.



# Archetype Analysis: Algorithms

**Will it work for huge datasets?**

# Efficient algorithms

(Bauckhage et al., 2015):

Solve step 1 and combined (2,3)-step with Frank-Wolfe algorithm

Idea: use **linear optimization oracle** over constraint set.

- Linear minimization oracle

$$\begin{aligned}\Delta(\boldsymbol{x}) &= \operatorname{argmin}_{\boldsymbol{z}} \langle \boldsymbol{x}, \boldsymbol{z} \rangle, \\ \text{s.t. } &g(\boldsymbol{z}) \leq \kappa\end{aligned}$$

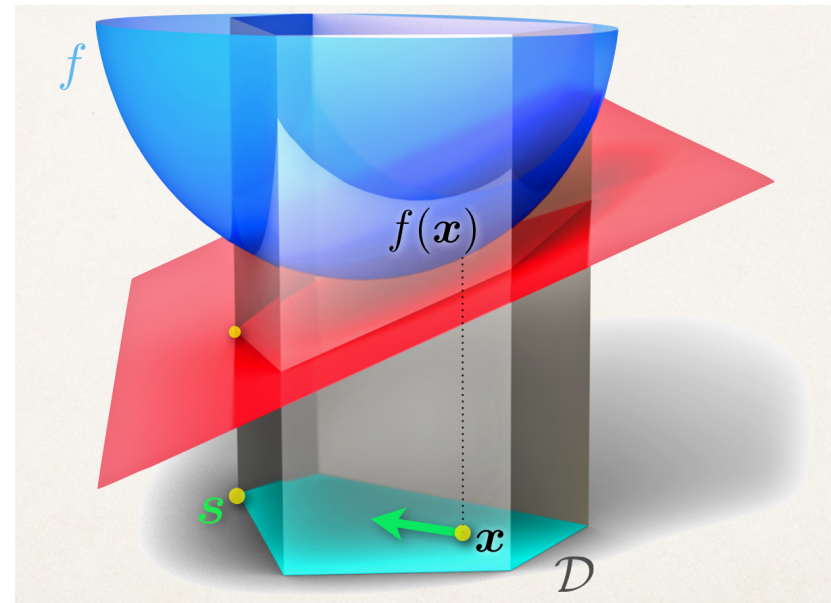
- Update  $\boldsymbol{x} \leftarrow (1 - \gamma)\boldsymbol{x} + \gamma\Delta(\nabla f(\boldsymbol{x}))$

- Decrease  $\gamma$

**Pros:** Highly efficient for one fixed value of  $\kappa$

**Cons:** Not efficient for the whole solution path

$\leadsto$  **model selction trick does not work well.**



M. Jaggi, 2015

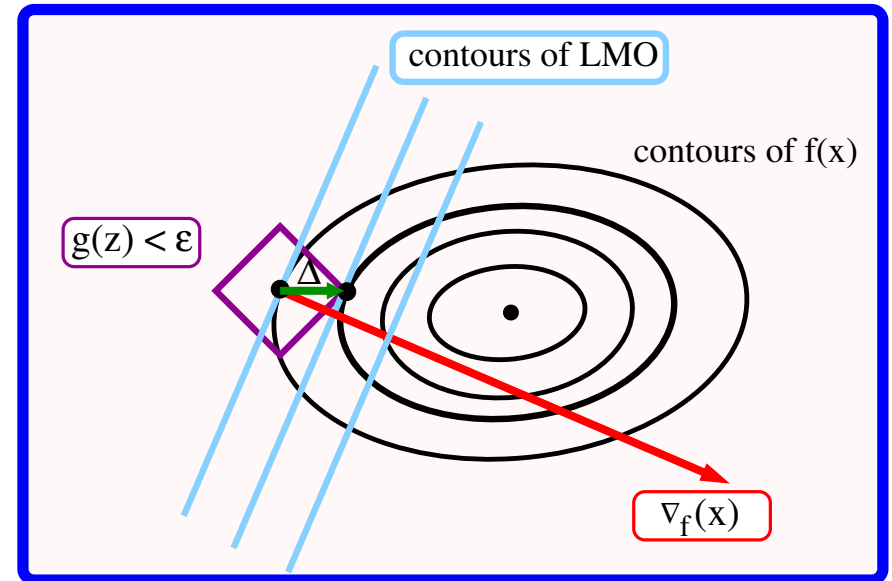
# Alternative: Forward stagewise

**Technically similar to Frank-Wolfe:**

- Linear minimization oracle (LMO)

$$\begin{aligned}\Delta(\mathbf{x}) &= \operatorname{argmin}_{\mathbf{z}} \langle \mathbf{x}, \mathbf{z} \rangle, \\ \text{s.t. } g(\mathbf{z}) &\leq \epsilon \ll \kappa\end{aligned}$$

- Update  $\mathbf{x} \leftarrow \mathbf{x} + \Delta(\nabla_f(\mathbf{x}))$



**...but very different behaviour:**

- incremental path following behaviour  
 $\leadsto$  efficient for computing the whole solution path
- built-in monotonicity “regularization”  $\leadsto$  very stable,  
can be extended to non-convex “norms” for increased sparsity.

# Forward stagewise

Consider step 1 in AT analysis:

Given  $K$  archetypes  $Z_{K \times p}$ , update the compositions  $\mathbf{a}_i$  of the  $i$ -th object  $\mathbf{x}_i$  under convexity constraints:

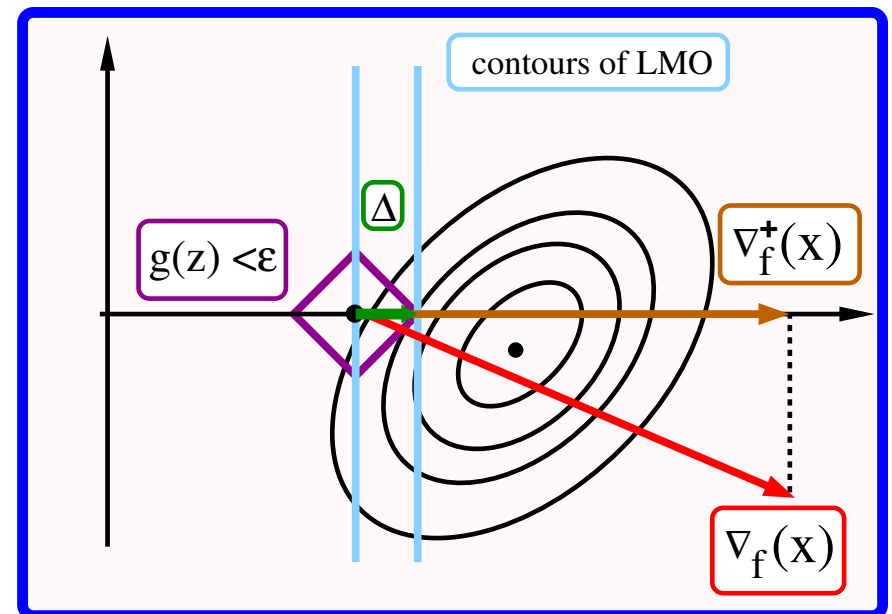
$$\mathbf{a}_i = \arg \min_{\mathbf{a} \in \mathbb{R}_+^p : \mathbf{a}^t \mathbf{1} = 1} \|\mathbf{x}_i - Z^t \mathbf{a}\|^2.$$

This is a non-negative lasso estimate

$$\begin{aligned} \min_{\mathbf{a}} \quad & \|\mathbf{x}_i - Z^t \mathbf{a}\|^2 \\ \text{s.t.} \quad & \|\mathbf{a}\|_1 = 1, \quad a_j \geq 0. \end{aligned}$$

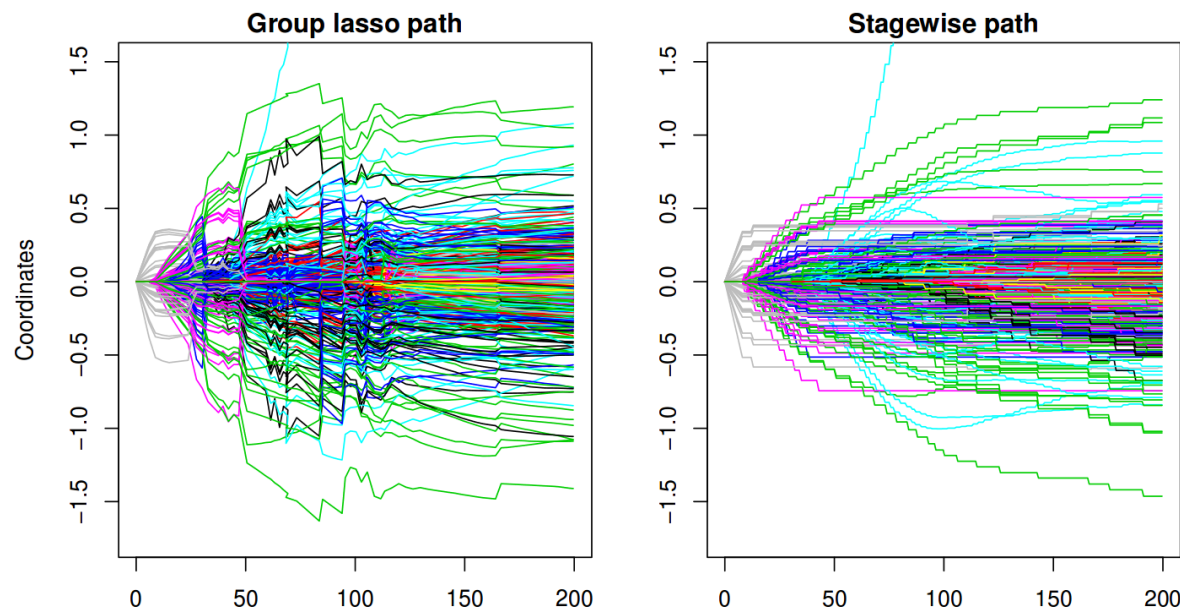
Use projected gradient:

$$\mathbf{a} \leftarrow \mathbf{a} + \Delta(\nabla_f^+(\mathbf{a}))$$



# Forward stagewise

- Useful if LMO can be computed easily: all norms, block-norms etc.
- NN lasso: LMO is intersection of simplex and linear function  
 $\leadsto$  find best vertex  $l$ , update  $l$ -th component of  $\mathbf{a}$  as  $a_l \leftarrow a_l + \epsilon$   
 $\leadsto \mathbf{a}$  is monotone increasing in every component.
- Simple and efficient: iterate until  $\|\mathbf{a}\|_1 = 1 \leadsto 1/\epsilon$  iterations needed.
- Conceptually the same behaviour for group-lasso estimate in step 2.



R. Tibshirani, 2015

## Further algorithmic tricks

- Pre-select candidate points on convex hull:
  - Points on **convex hull in any linear projection** are also on the “full” convex hull.
  - Convex hull computation very efficient in 2D:  $O(n \log n)$
  - Randomly project data to planes, compute points on convex hull, aggregate.
- Alternative to random projections: use **pairwise PCA projections**.
  - PCA is a good preprocessing step anyway, since convex polygon with  $K$  vertices can be embedded in a space with  $< K$  dimensions.
- We can easily solve AT problems with  $> 100$  millions of objects.

# Archetype Analysis: Data representation

**How to encode the data such that we can see “triangles”?**



# Representation issues

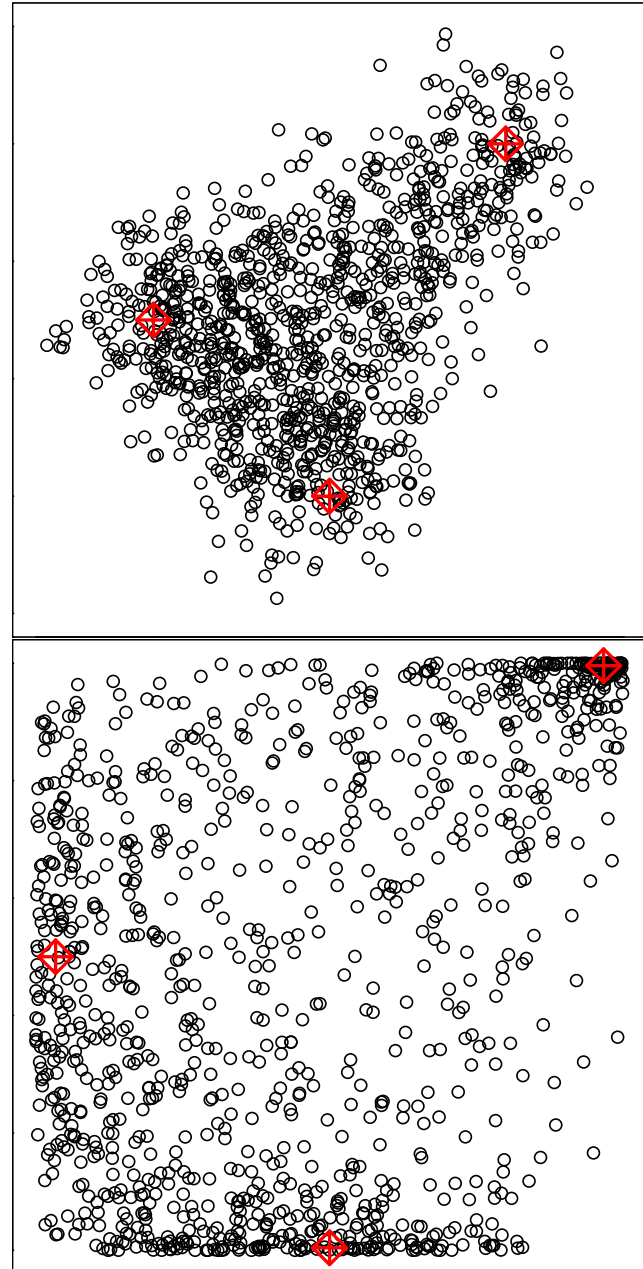
## Representations

- different domains  
( $p$ -values, body size, weight)
- monotone transformations (e.g. log)

**Problem:** archetypal analysis is sensitive to choice of representation

**Solution:** Copula extension

- representation independence
- robust against outliers
- mixed data & missing values



# Copula

## Copula density

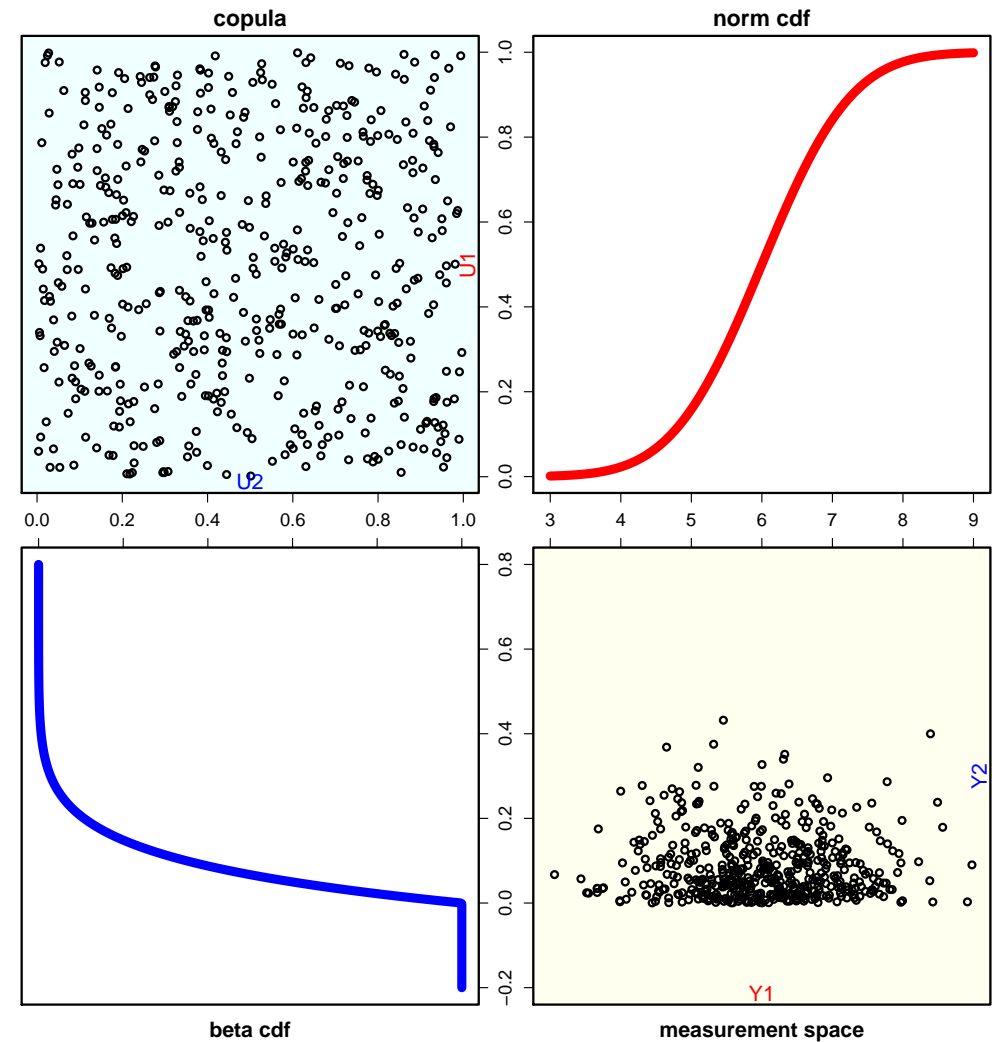
- $p$ -dim pdf on  $[0, 1]^p$
- uniform marginals
- defines dependency structure

## Property

- construct arbitrary multivariate distribution

$$y_j = F_j^{-1}(u_j)$$

(Sklar 1959)



Independence Copula

# Copula

## Copula density

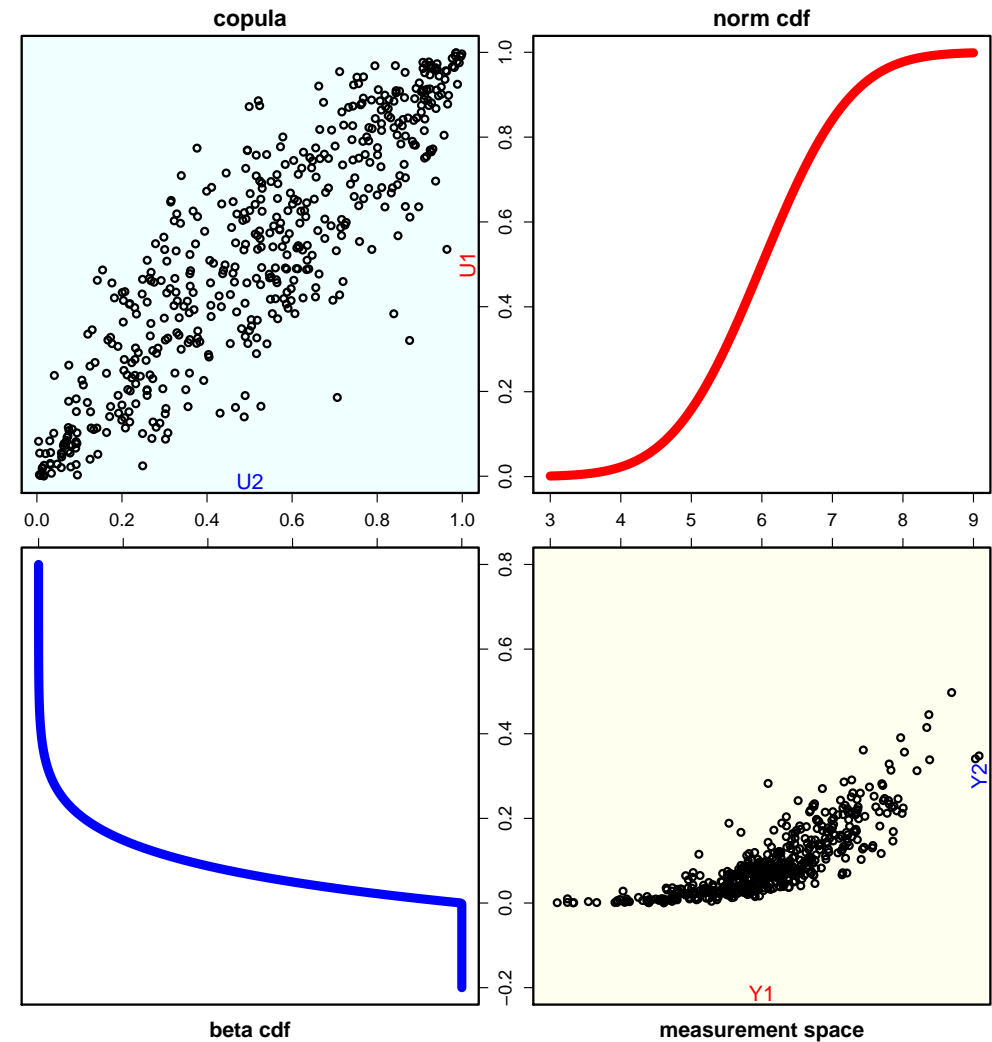
- $p$ -dim pdf on  $[0, 1]^p$
- uniform marginals
- defines dependency structure

## Property

- construct arbitrary multivariate distribution

$$y_j = F_j^{-1}(u_j)$$

(Sklar 1959)



## Gaussian Copula

# Semi-parametric Gaussian Copula

## Hierarchical model:

$$\begin{array}{c} R \\ | \\ \mathbf{x} \sim \mathcal{N}(\mathbf{0}, R) \\ | \\ y_i = f_i(x_i) \end{array}$$

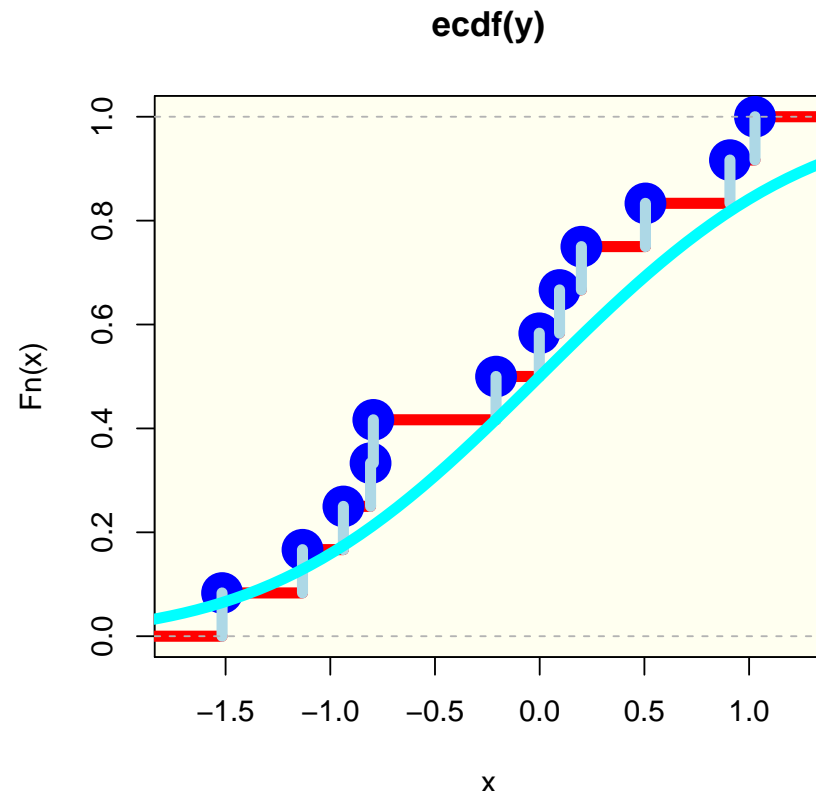
- ATs invariant against monotone marginal transformations
- ATs depend only on ranks  
 $\leadsto$  insensitive to outliers.

- Idea: reconstruct latent space

$$x_j = (\Phi^{-1} \circ F_j)(y_j)$$

- Use empirical cdf's

$$\hat{F}_j = \frac{\text{ranks}(Y[\cdot, j])}{n+1}$$



# Mixed Continuous / Discrete Data

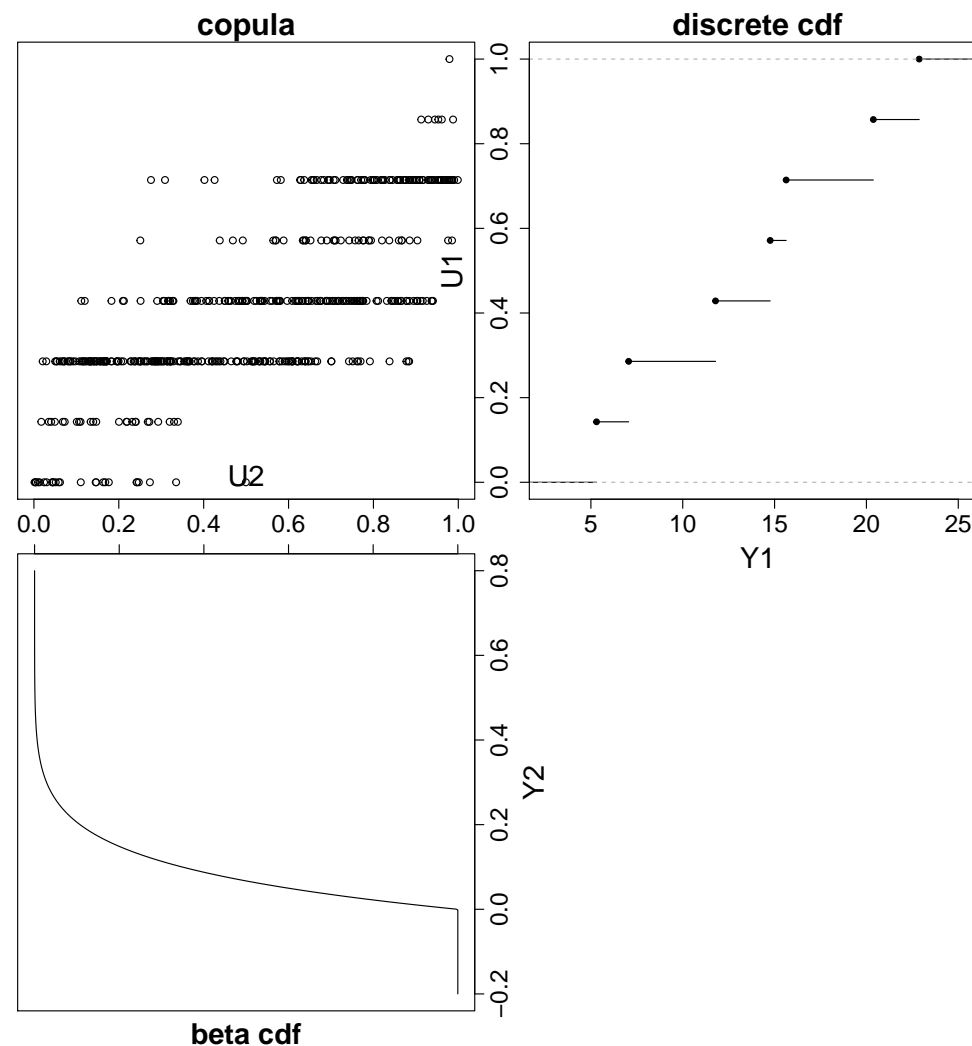
**Problem with discrete data:** ties, empirical copula no longer uniform

**Extended rank likelihood** (Hoff, 2007):

stochastic association-preserving mapping

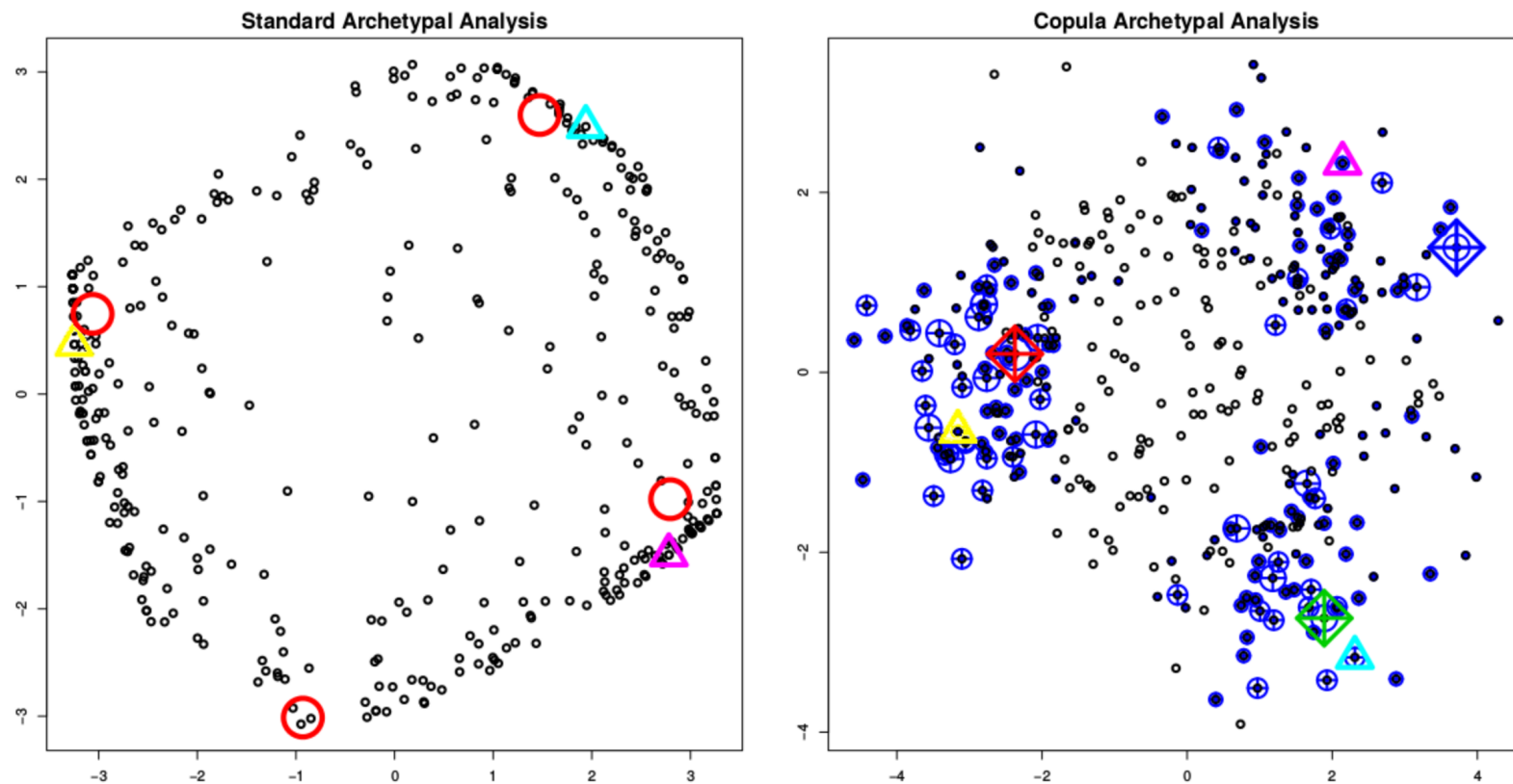
**Algorithm:** Gibbs sampler

1. sample latent variables  $x, R|x$
2. compute archetypes  $z$



# Artificial Data

- Monotone transformation with beta marginals
- Variables quantised to 5 levels



# Why use the Gaussian copula?

## Generative model

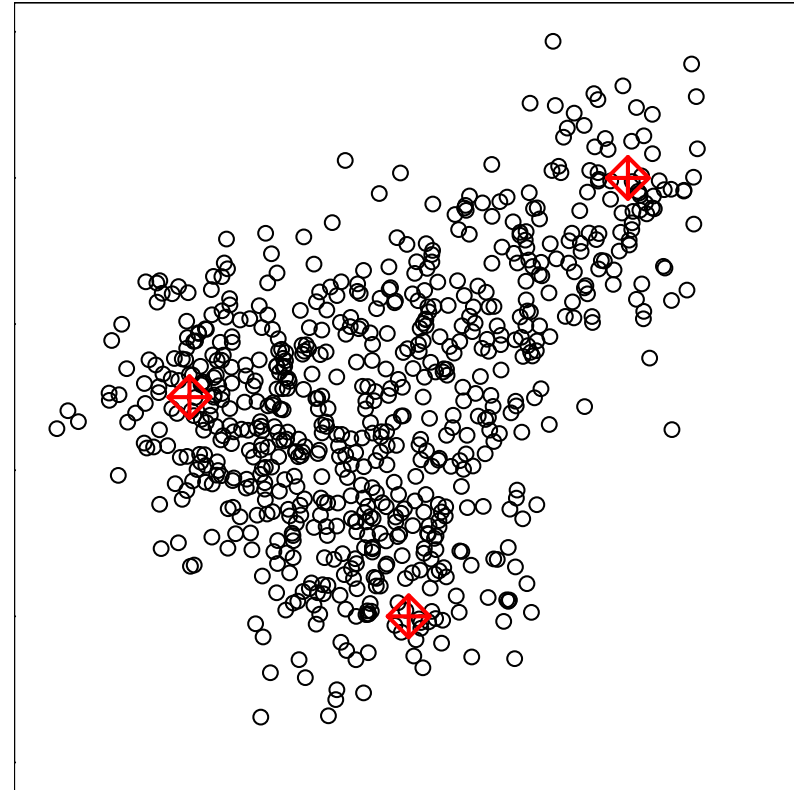
$$\mathbf{a}_i \sim \text{Dir}_K(\boldsymbol{\alpha})$$

$$\mathbf{x}_i | \mathbf{Z}, \mathbf{a}_i \sim \mathcal{N}(\mathbf{Z}^t \mathbf{a}_i, \eta I_p)$$

$$X | Z, A \sim \mathcal{MN} \left( \overbrace{AZ}^M, I, \eta I \right)$$

$$X^t X \sim \mathcal{W}_{\text{nc}}(n, nI, M^t M) \\ \approx \mathcal{W}_c \left( n, \frac{1}{n} M^t M + \eta I \right)$$

(Steyn and Roux, 1972)



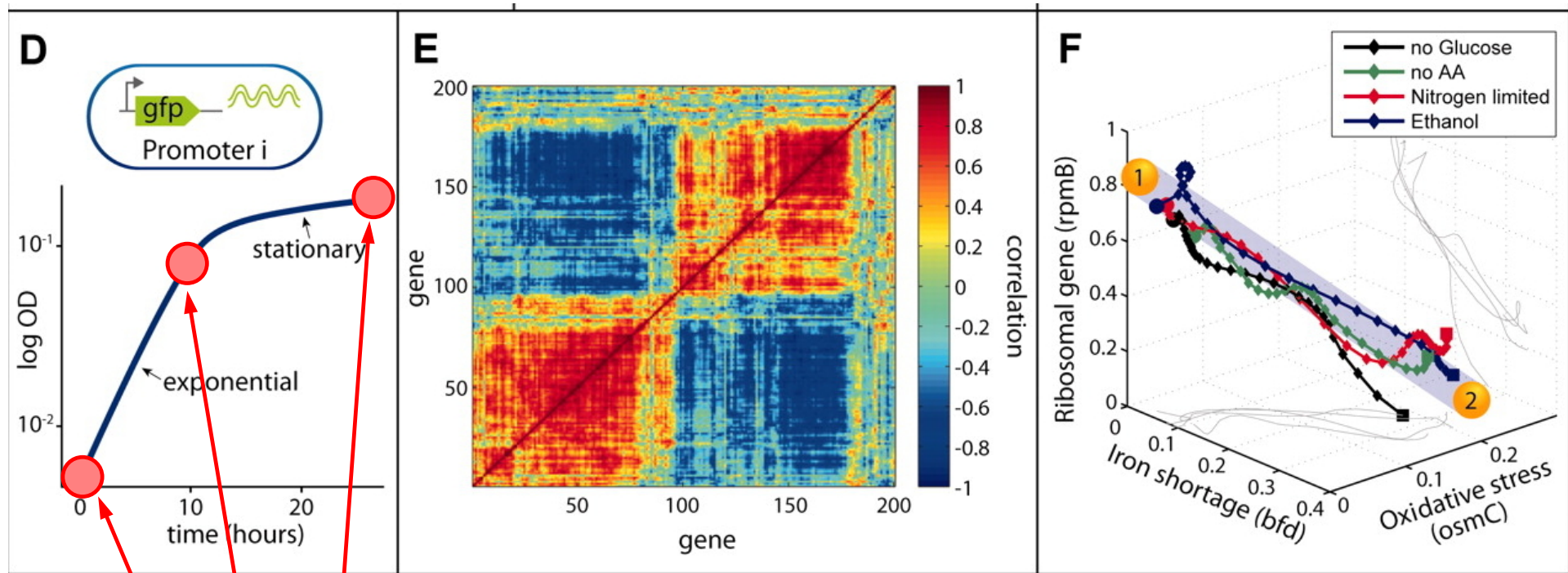
Under the assumed generative model,  
a Gaussian covariance structure is plausible.



# Archetype Analysis: Applications

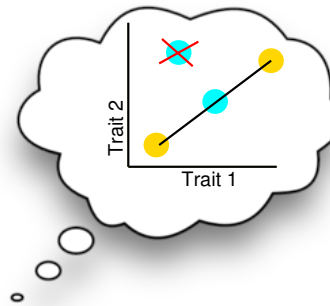
**What is it good for?**

# Analysis of E.coli Data



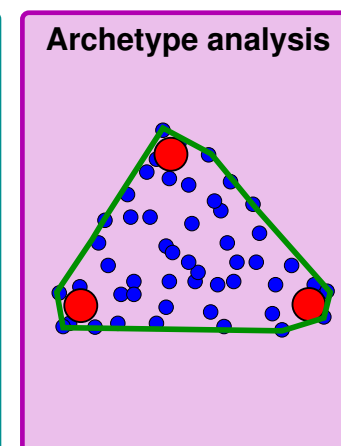
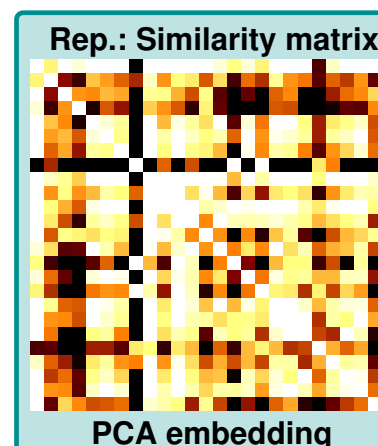
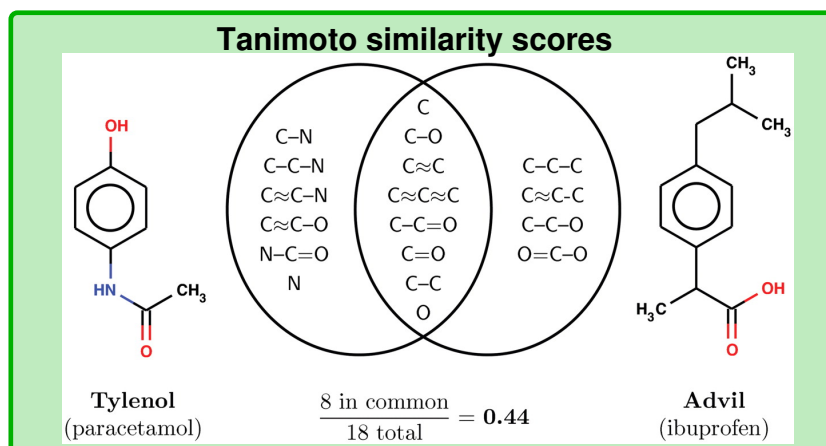
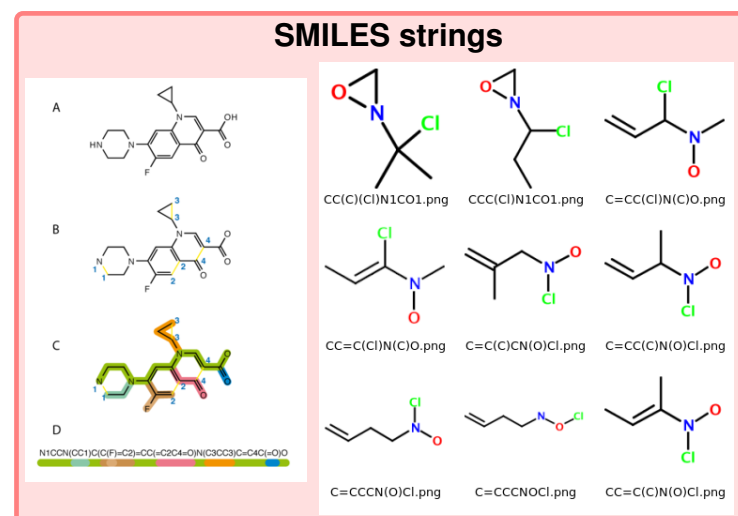
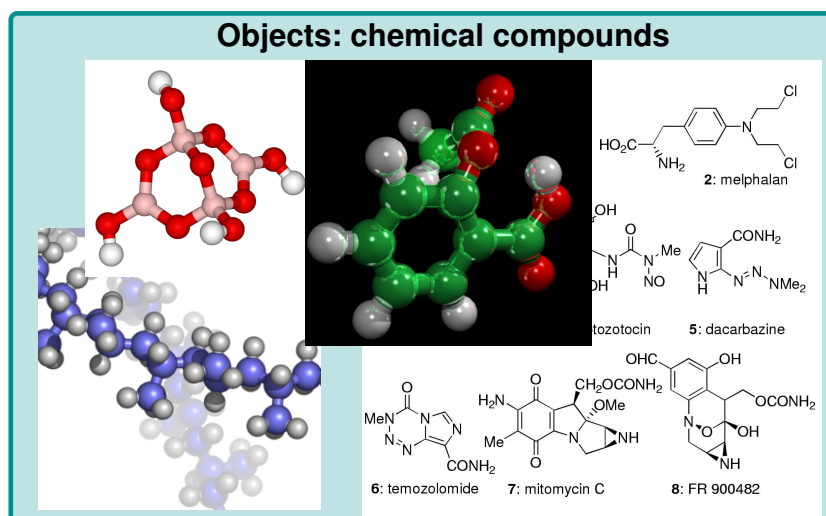
Shoval et al., Science, 2012

**Identified Archetypes**

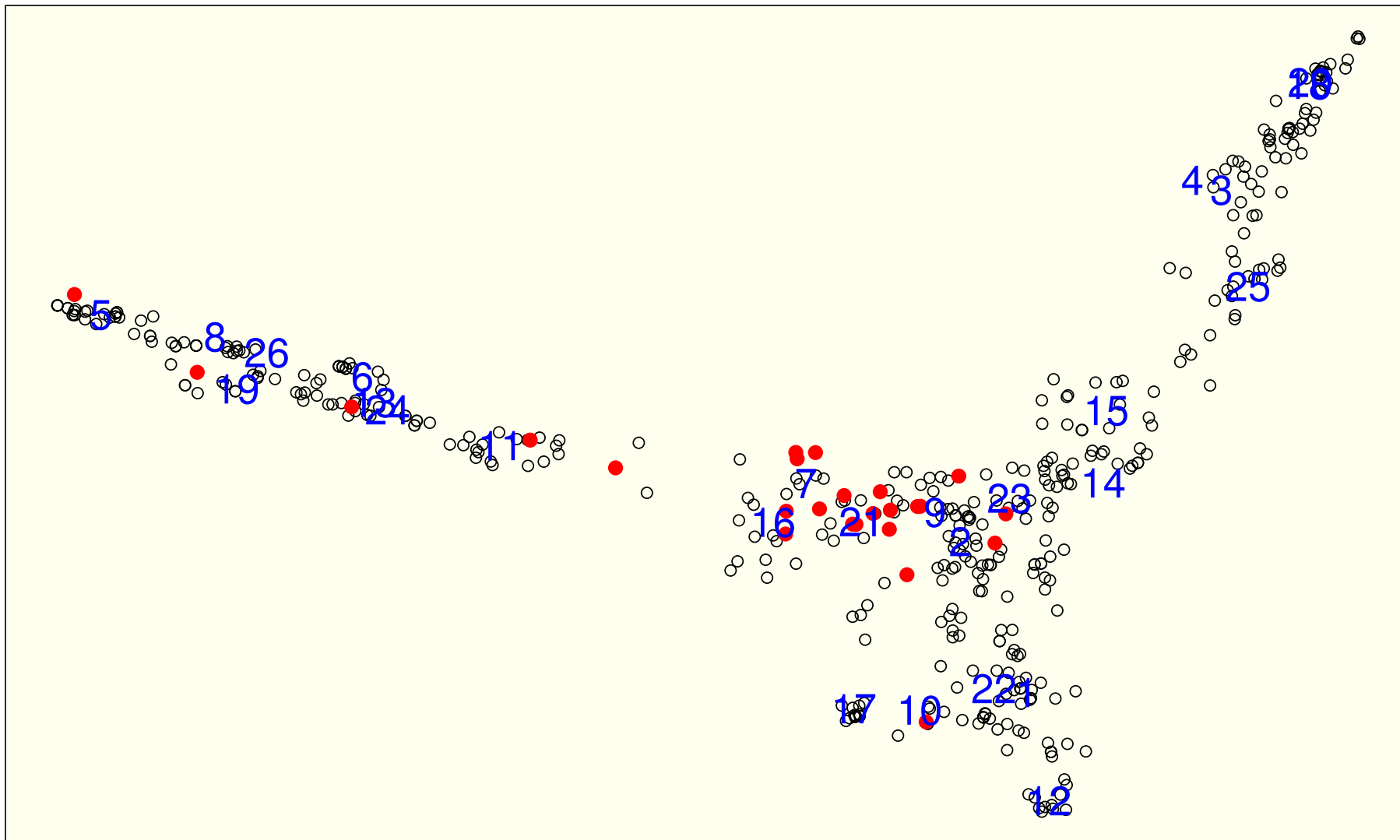


# Archetypical Chemical Compounds

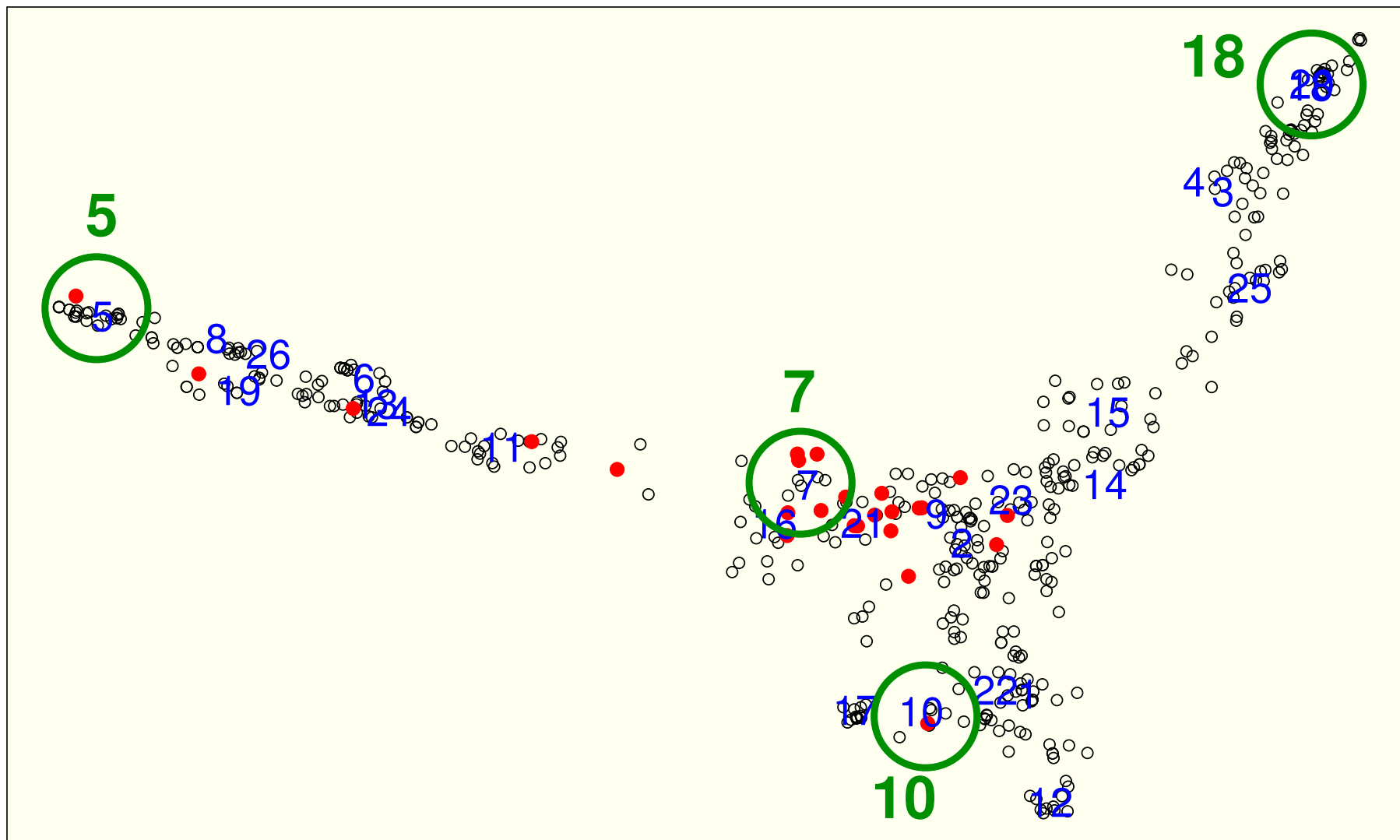
Idea: find ATs in list of compounds identified in an AIDS antiviral screen performed by the Developmental Therapeutics Program of the NCI/NIH, enriched with all available anti-HIV drugs



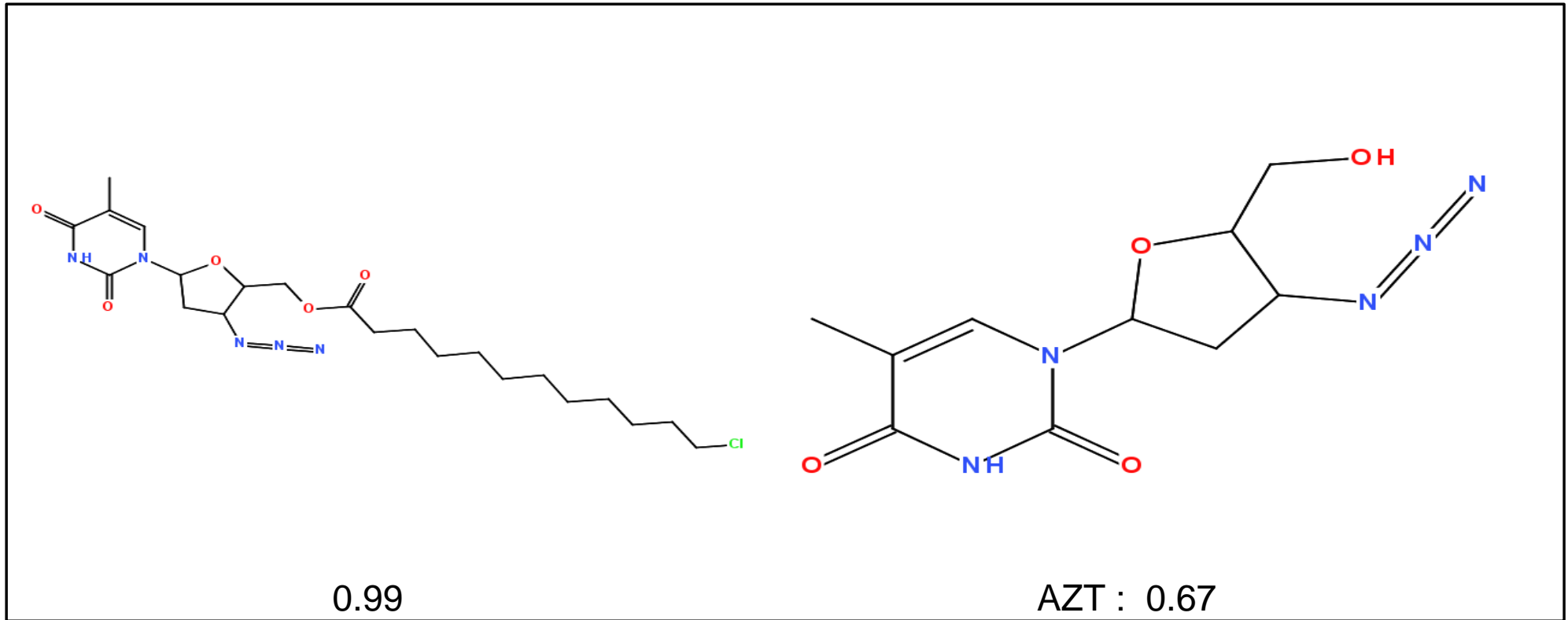
# Archetypical Chemical Compounds



# Archetypical Chemical Compounds

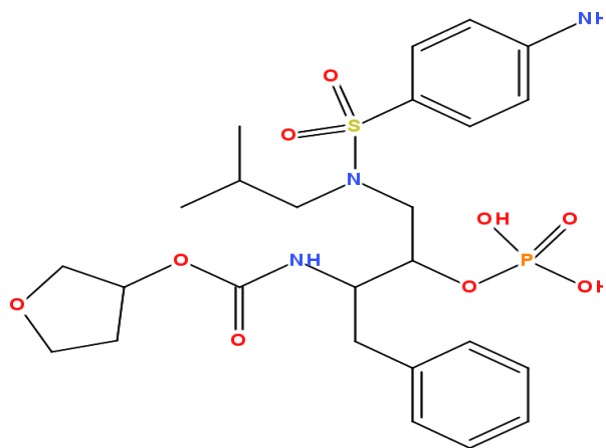


# Compounds explained by AT 5

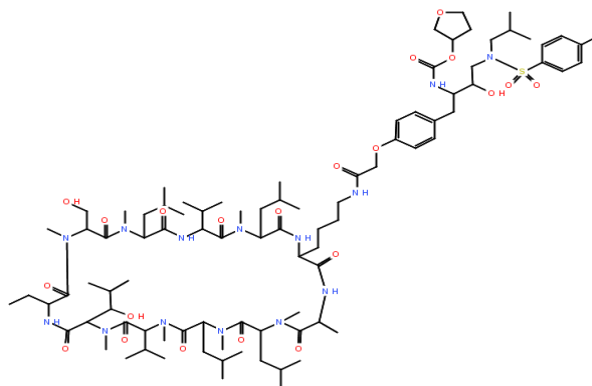


AZT It is of the nucleoside reverse-transcriptase inhibitor (NRTI) class. It inhibits the enzyme (reverse transcriptase) that HIV uses to synthesize DNA.

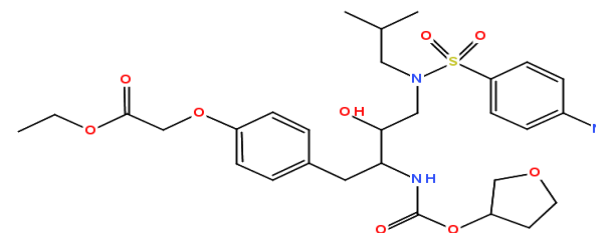
# Compounds explained by AT 7



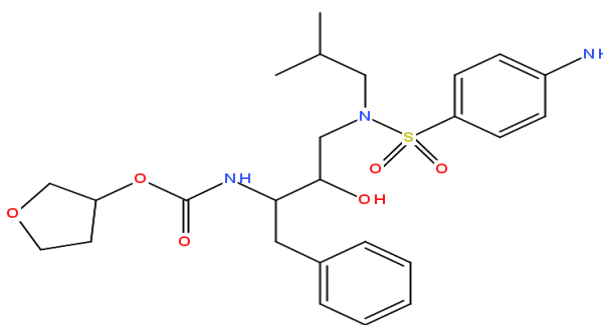
Fosamprenavir : 0.94



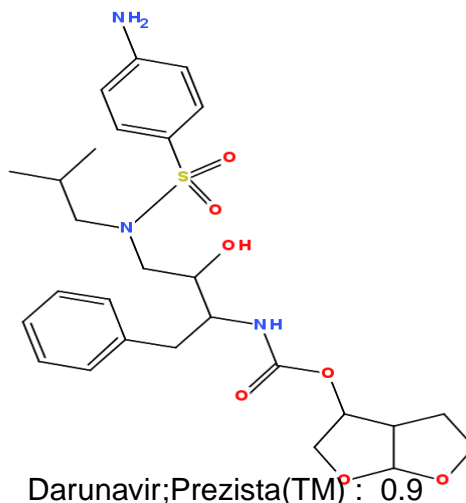
0.99



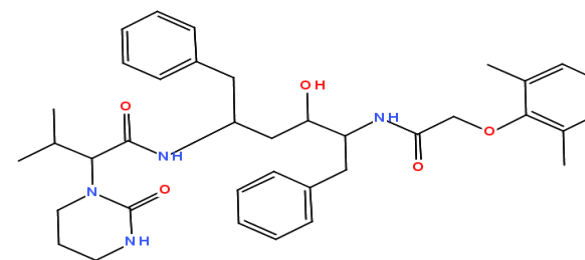
0.99



amprenavir : 0.99



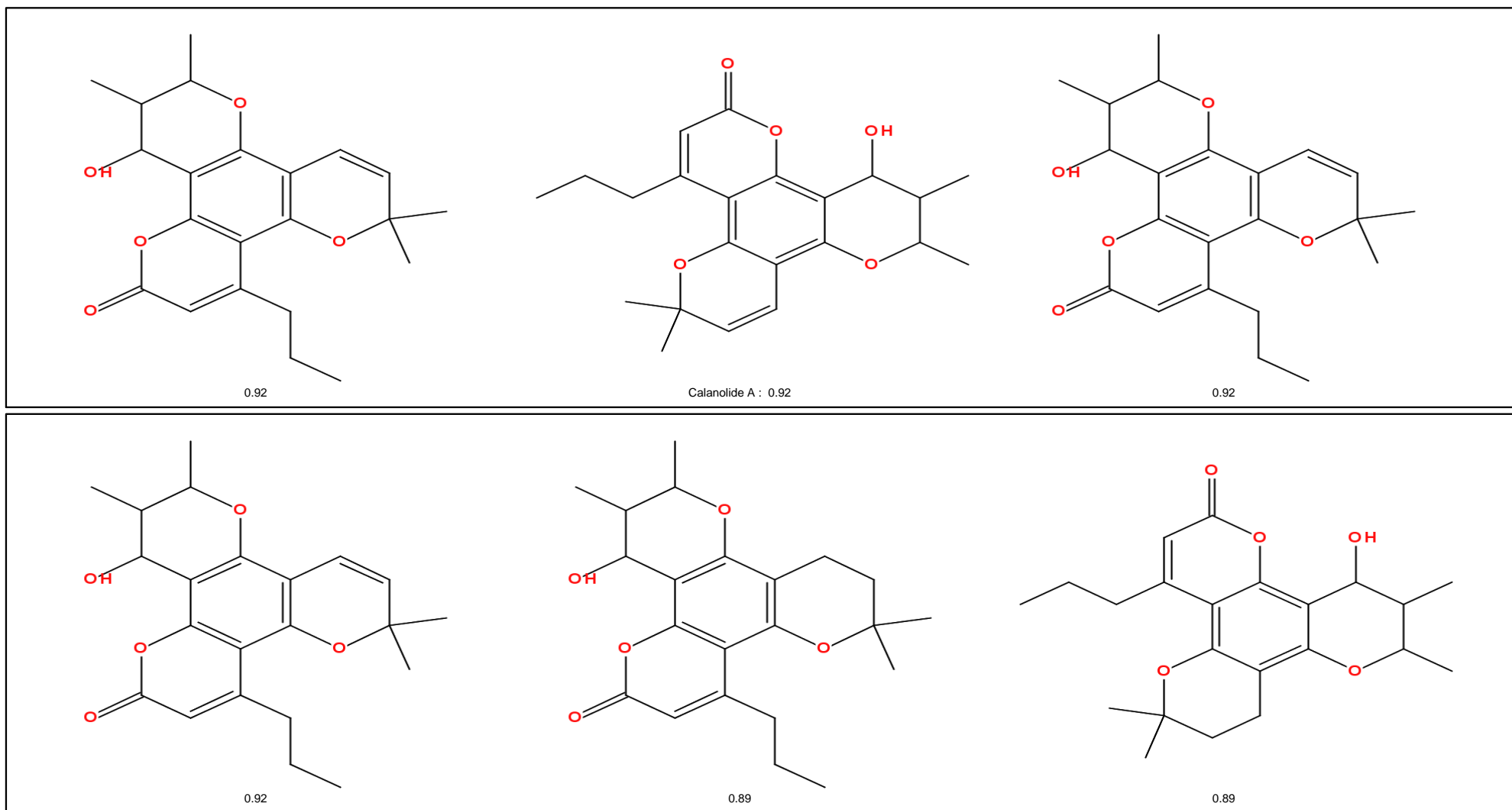
Darunavir;Prezista(TM) : 0.9



Aluviran;Lopinavir : 0.7

HIV protease inhibitors (they block a peptide cleaving enzyme).

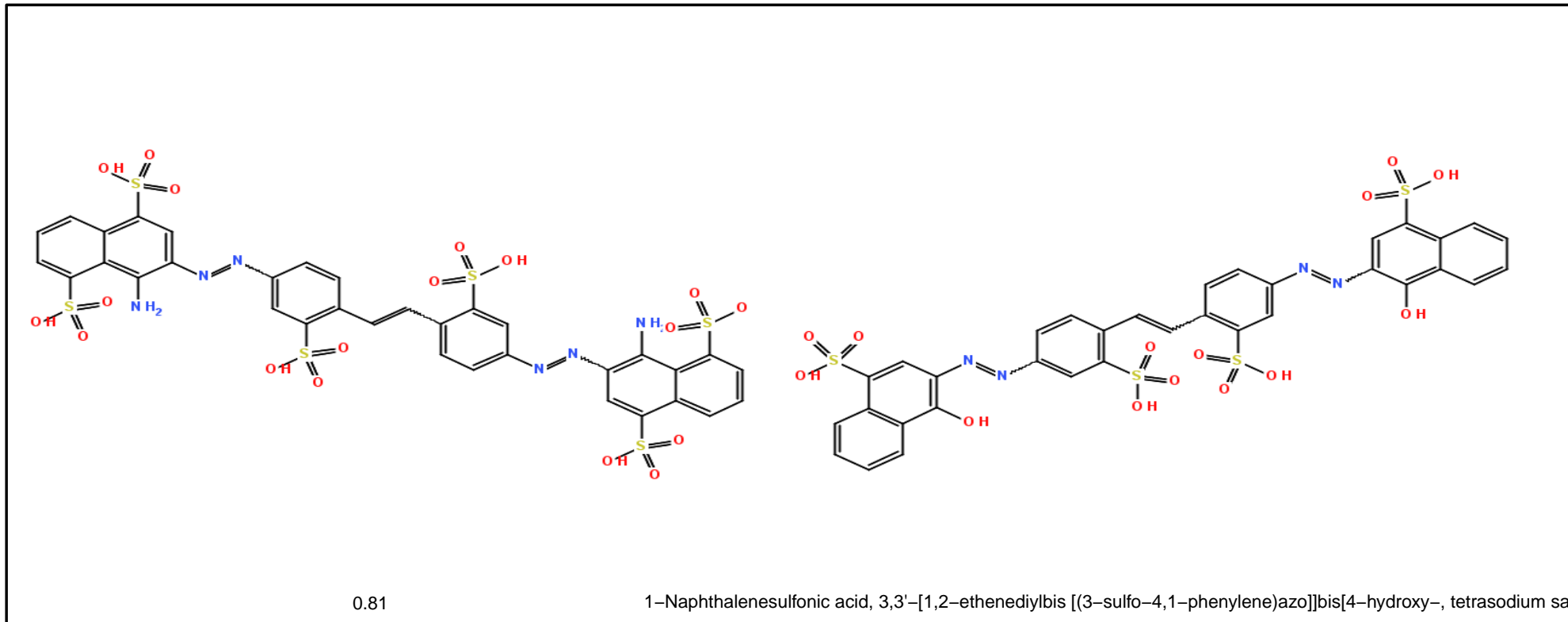
# Compounds explained by AT 10



Calanolide A is an experimental non-nucleoside reverse transcriptase inhibitor (NNRTI).

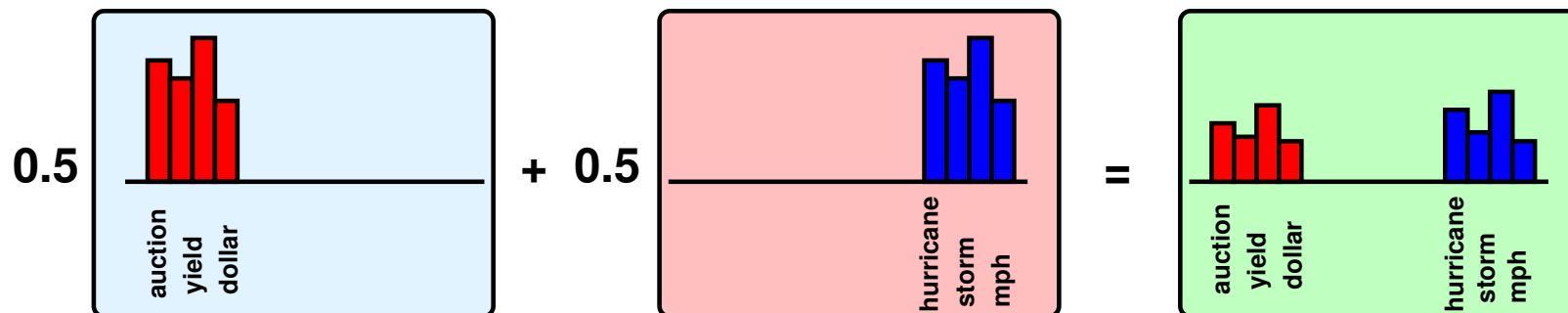


# Compounds explained by AT 18

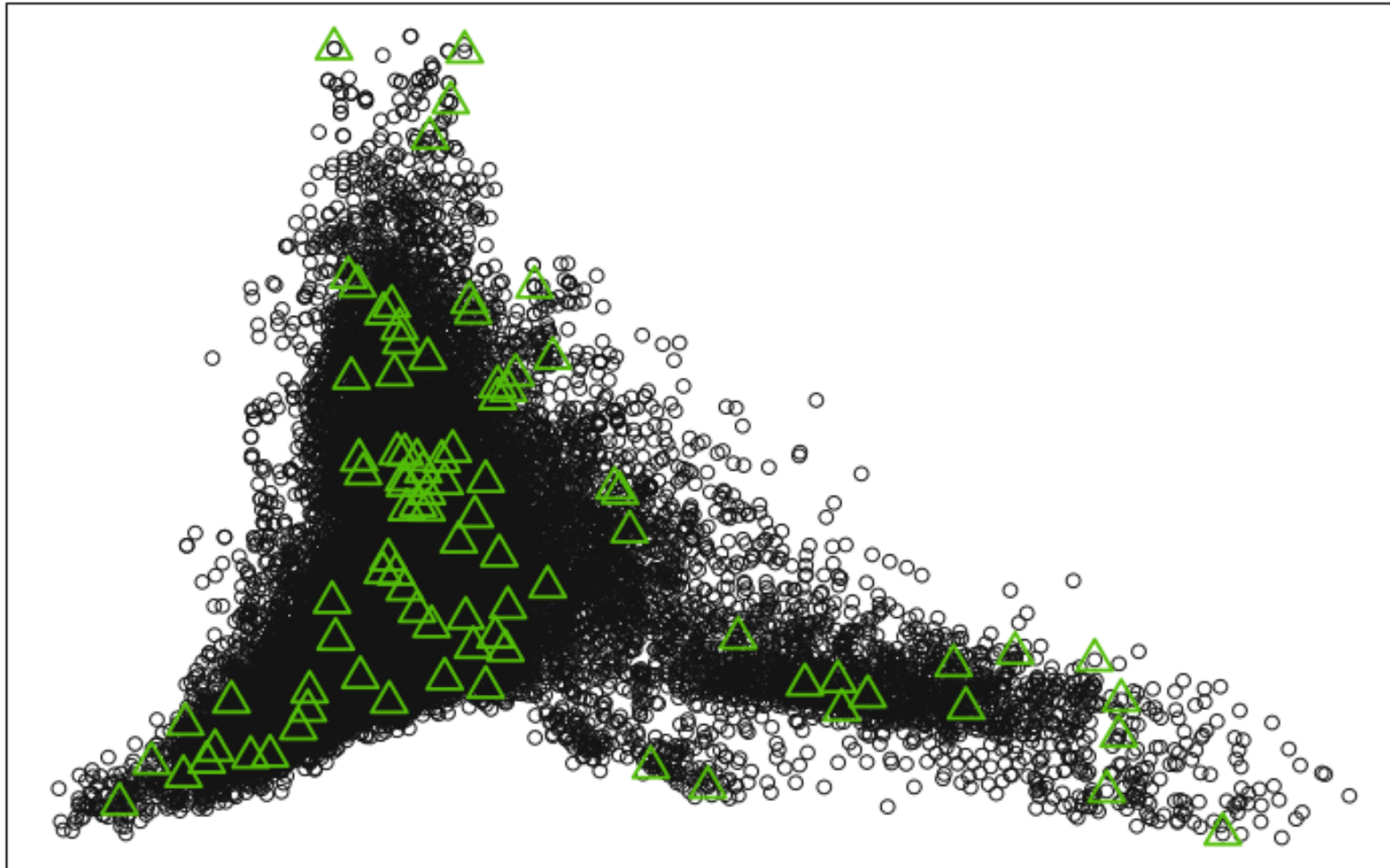


# Text categorization

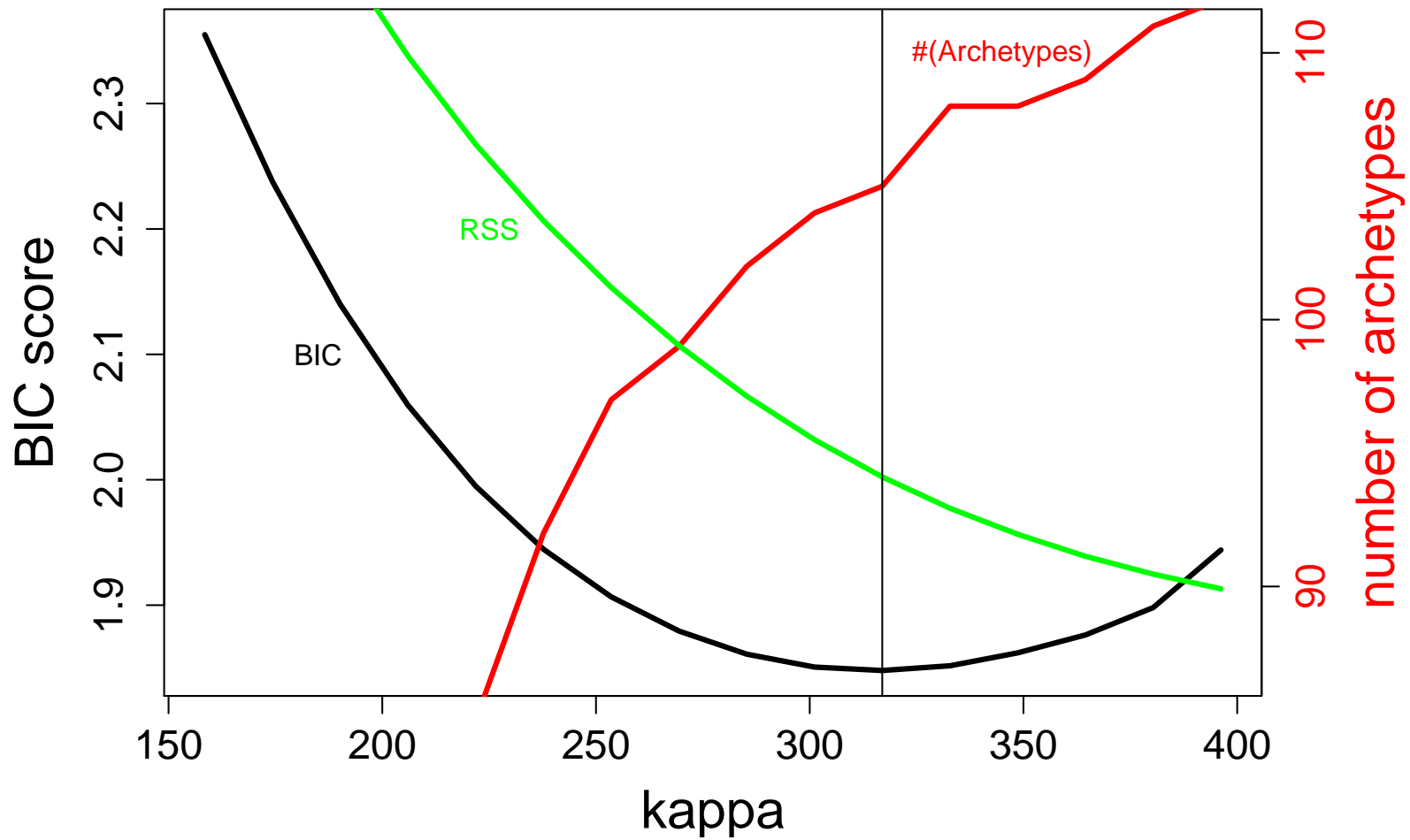
- Reuters Corpus Volume 1 (RCV1): an archive of news documents
- 4 categories: Economics, Government, Corporate, Markets
- 23149 documents, vocabulary of 57180 words.
- Documents represented by word frequencies:  
**Term frequency (TF) times Inverse Document Frequency (IDF).**
- Automatically detect “pure” or archetypical documents (which might represent “pure” topics)



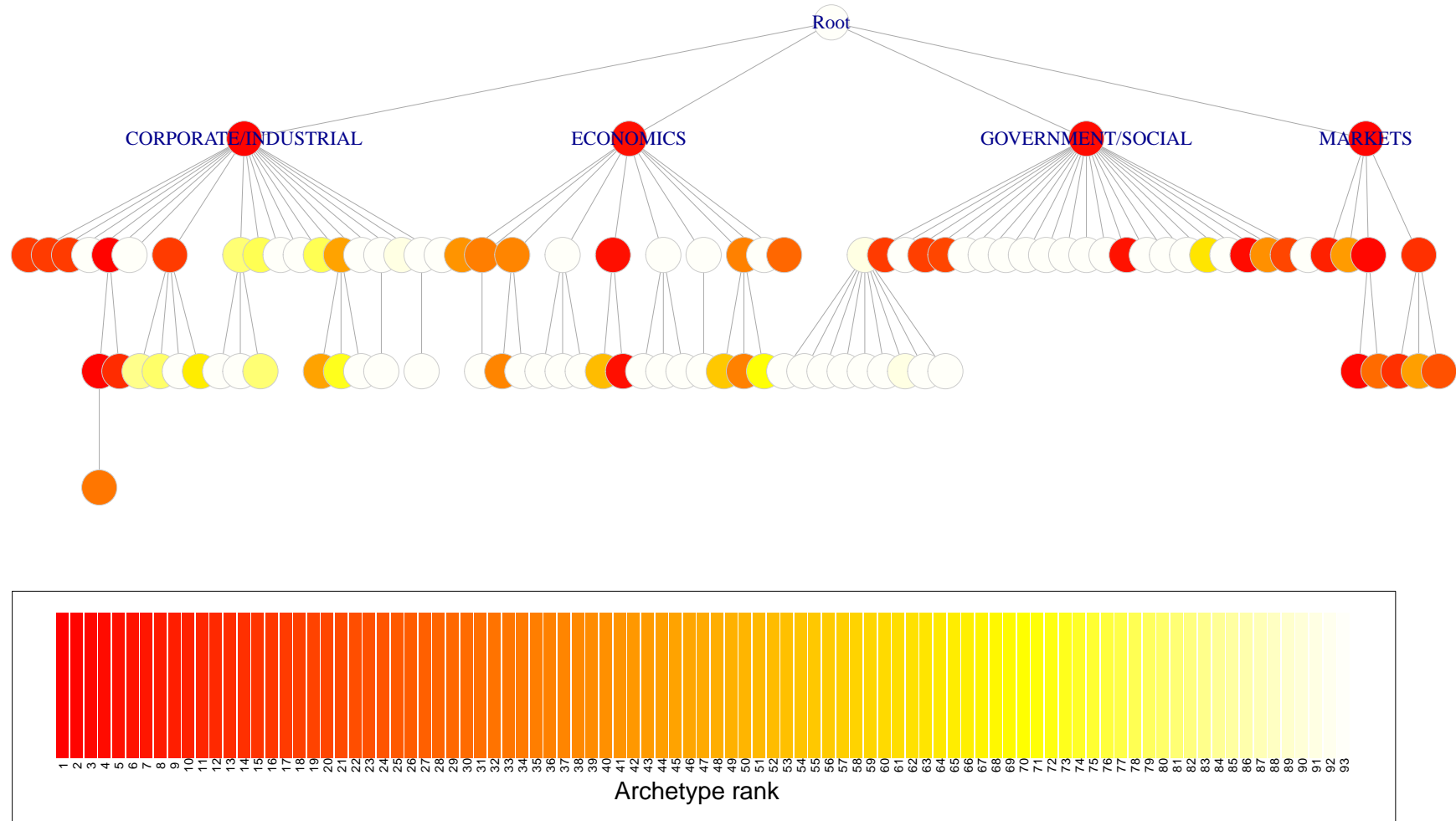
# Text categorization



# Text categorization



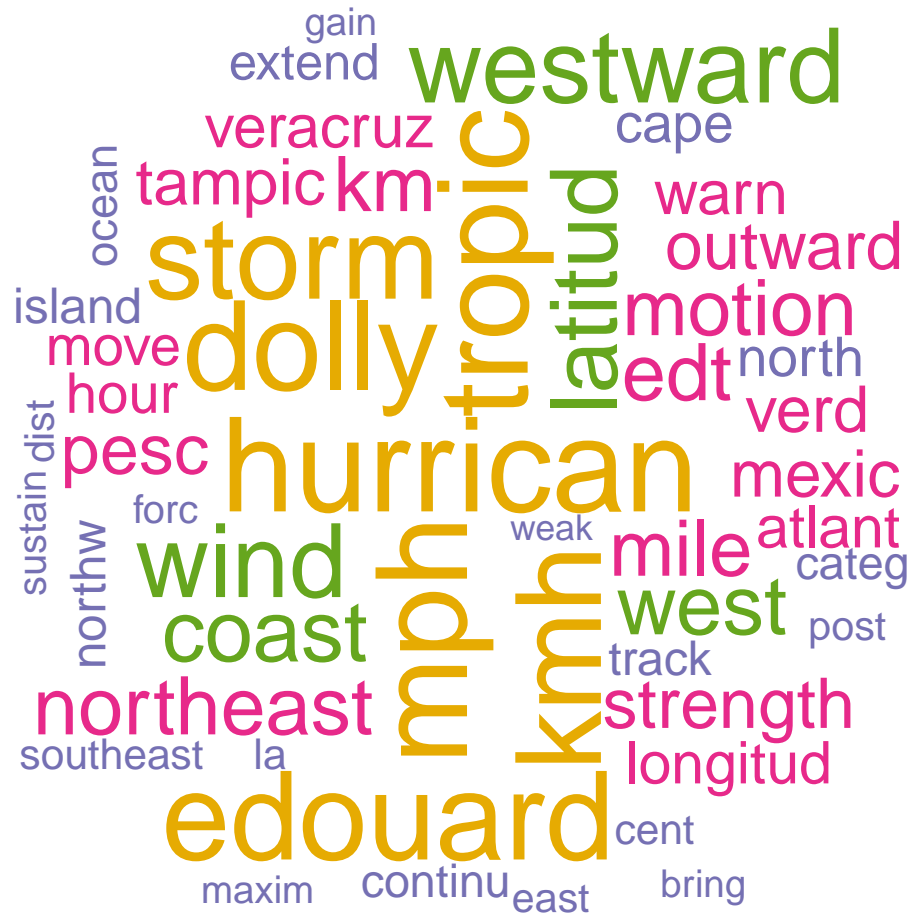
# Text categorization



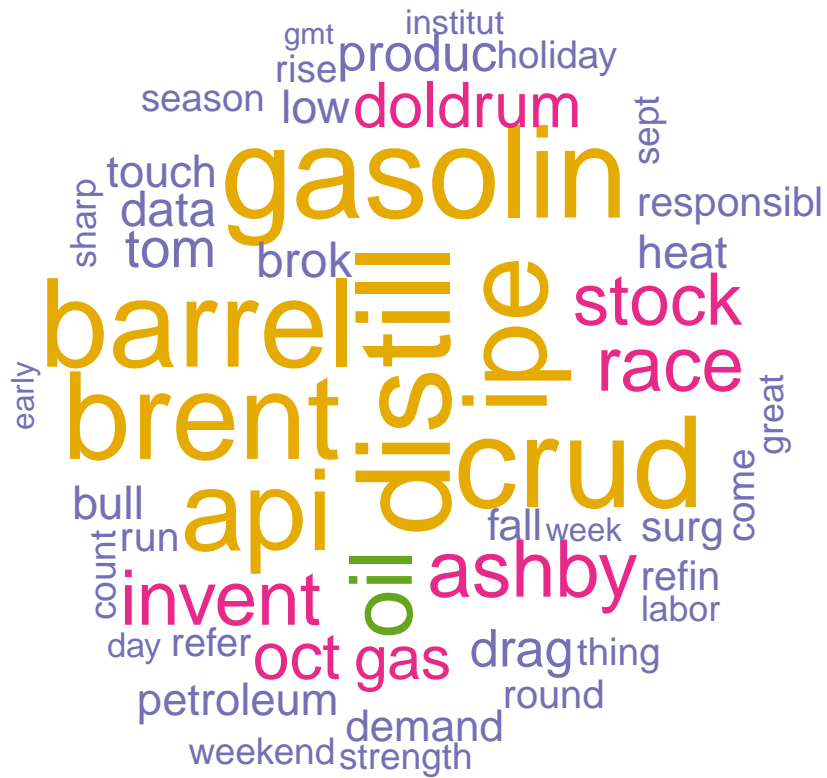
# Text categorization



# Text categorization



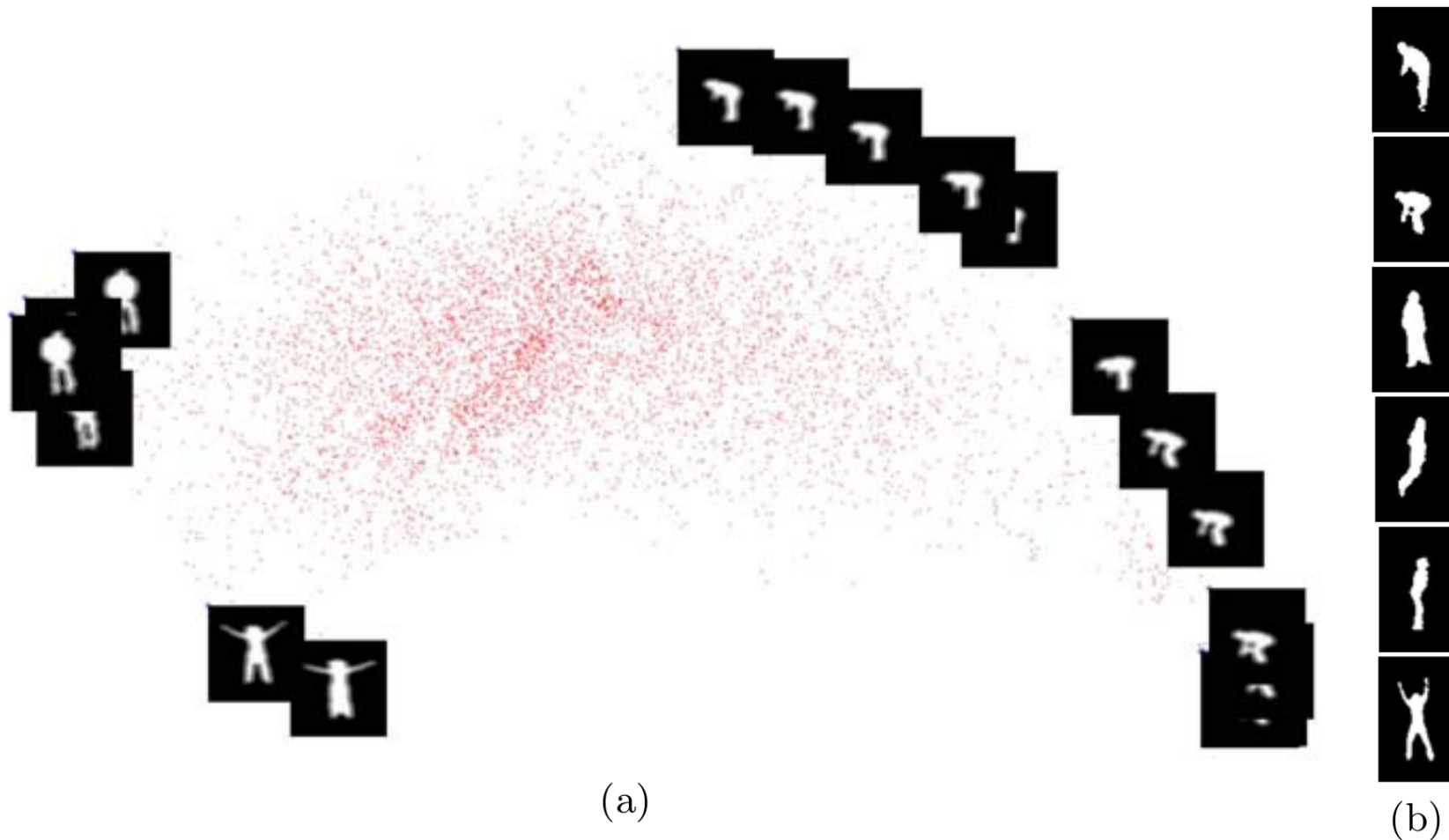
# Text categorization





# Pose analysis

Making Archetypal Analysis Practical 279



**Fig. 5.** (a) 2D projection of the Weizman set containing 5.000 body poses; points on the convex hull are shown as pictures. (b) 6 archetypal poses extracted from the data.

# Image Encoding

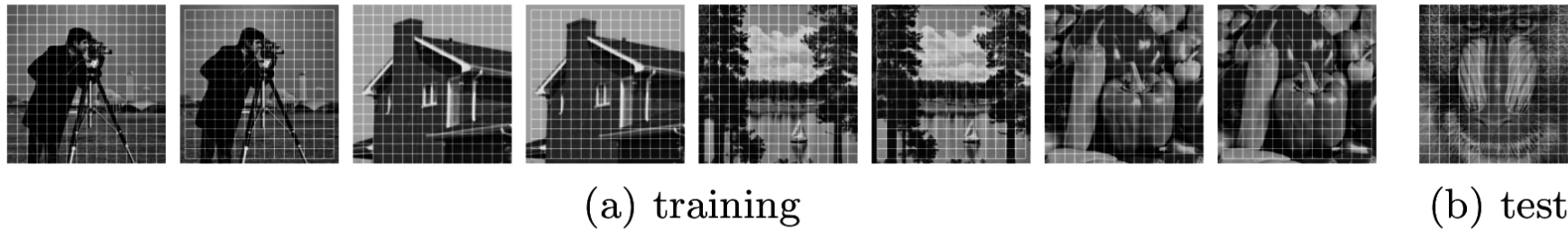
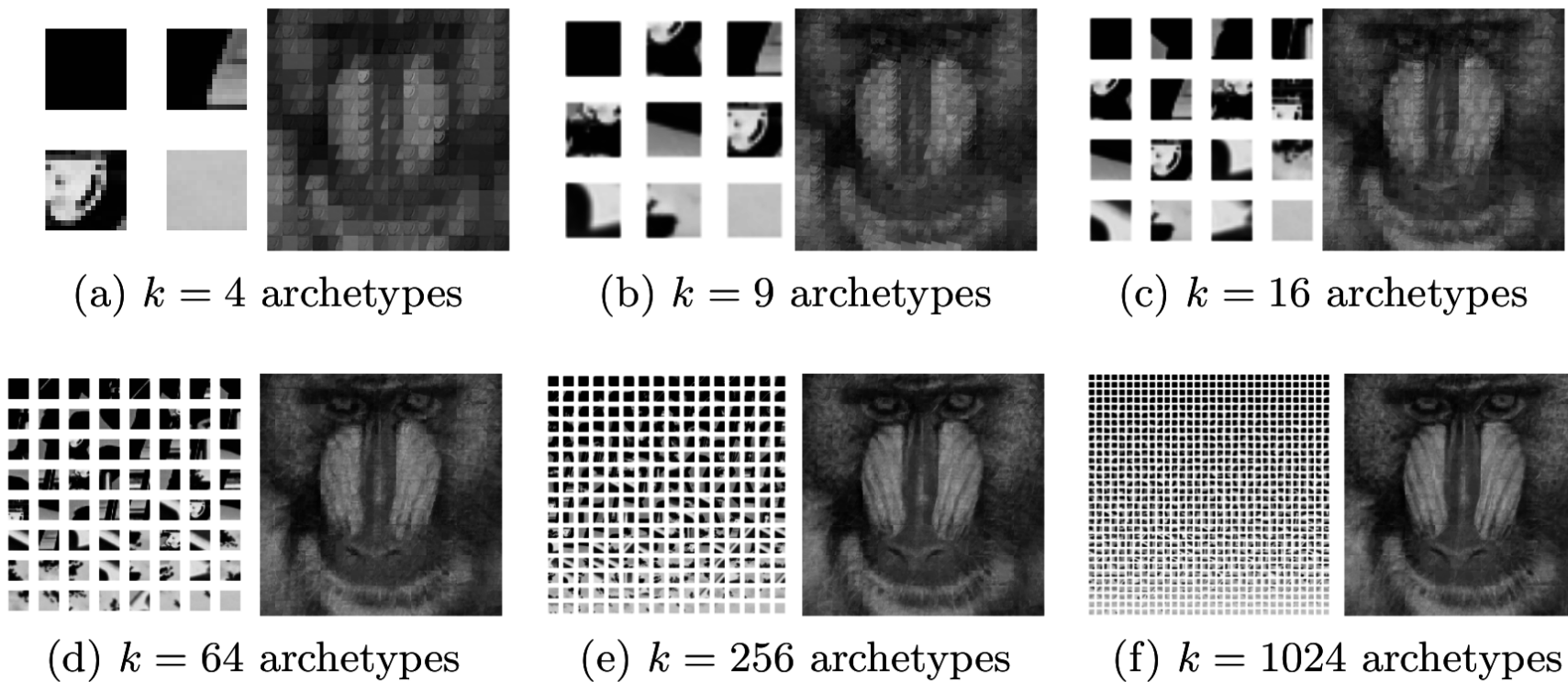


Fig. 3: Image patches ( $16 \times 16$  pixels) for archetypal autoencoding experiments.



# Conclusion

- Archetype analysis is a powerful tool identifying **representative objects** in large data collections.
- Many **technical challenges**:
  - data types, representation
    - ∼ invariances due to semi-parametric copula construction
  - computational/memory complexity
    - ∼ stagewise forward algorithms, convex hull approximations, etc.
- Many **application areas**: objects can be bacteria, genes, documents, chemical compounds, images, etc. etc.
- **Open questions**: AT analysis essentially is an auto-encoding technique
  - ∼ neural implementations,
  - ∼ use as building-blocks in deep belief networks, etc.

# Acknowledgments

Department of Mathematics and Computer Science, University of Basel:  
**Dinu Kaufmann, Sebastian Keller, Damian Murezzan, Sonali Parbhoo, Sandhya Prabhakaran, Mélanie Rey, Aleksander Wieczorek, Mario Wieser**

University Hospital Zurich: **Francesca Di Giallonardo, Yannick Duport, Christine Leemann, Stefan Schmutz, Nottania K. Campbell, Beda Joos, Osvaldo Zagordi, Huldrych F. Günthard, Karin J. Metzner**

Department of Biosystems Science and Engineering, ETH Zurich:  
**Armin Töpfer, Christian Beisel, Niko Beerenwinkel**

Functional Genomics Center Zurich: **Maria R. Lecca, Andrea Patrignani**

Inst. Medical Virology, U Zurich: **Peter Rusert, Alexandra Trkola**