



Making Archetypal Analysis Practical

Christian Thureau // christian.thureau@unbelievable-machine.com

February 11, 2015

Table of Contents



k-Means Clustering

what and why not?

Archetypal Analysis

what, why, and how?

Making Archetypal Analysis practical

Applications

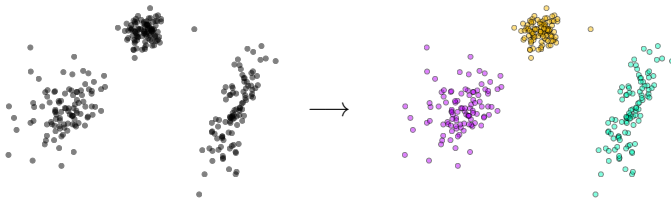
The average human has one breast and one testicle.

Des MacHale

k-Means Clustering

setting:

- assume a data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^m$
- determine k clusters $C_i \subset X$ such that
 - elements assigned to C_i are *similar*
 - $C_i \cap C_j = \emptyset$ for $i \neq j$



approach:

- define *similarity* and *clusters* in terms of **centroids** $\mu_i \in \mathbb{R}^m$

$$C_i = \left\{ \mathbf{x}_j \in X \mid \|\mathbf{x}_j - \mu_i\|^2 \leq \|\mathbf{x}_j - \mu_l\|^2 \forall l \neq i \right\}$$

- determine suitable centroids
- consider objective function

$$E(k) = \sum_{i=1}^k \sum_{\mathbf{x}_j \in C_i} \|\mathbf{x}_j - \mu_i\|^2$$

k-Means Clustering



algorithm:

at $t = 0$, initialize $\mu_1^{(t)}, \dots, \mu_k^{(t)}$

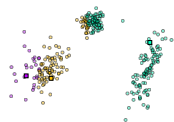
repeat until convergence

for each $\mu_i^{(t)}$ compute $C_i^{(t)}, n_i = |C_i|$

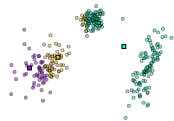
for each $C_i^{(t)}$ update the centroid as

$$\mu_i^{(t+1)} = \frac{1}{n_i} \sum_{x_j \in C_i^{(t)}} x_j$$

$t \leftarrow t + 1$

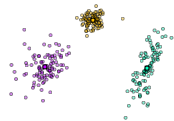


$t = 0$



$t = 1$

\vdots



$t = 5$

k -Means Clustering



note:

- k -means clustering is a matrix factorization problem

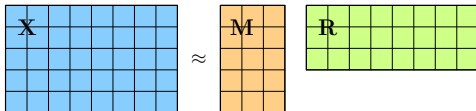
k-Means Clustering



note:

- k -means clustering is a matrix factorization problem
- upon convergence of the k -means algorithm, we have

$$\begin{aligned} \mathbf{X} &\approx \mathbf{MR} \\ \Leftrightarrow \mathbf{x}_j &\approx \mathbf{M}\mathbf{r}_j \end{aligned}$$



where

$\mathbf{X} \in \mathbb{R}^{m \times n} \Leftrightarrow$ data matrix

$\mathbf{M} \in \mathbb{R}^{m \times k} \Leftrightarrow$ centroid matrix

$\mathbf{R} \in \mathbb{R}^{k \times n} \Leftrightarrow$ indicator matrix

k-Means Clustering



note:

- k -means clustering is solving

$$\min_{M, R} \left\| \mathbf{X} - \mathbf{MR} \right\|^2$$

under several *implicit constraints*



note:

- k -means clustering is solving

$$\min_{M,R} \left\| \mathbf{X} - \mathbf{MR} \right\|^2$$

under several *implicit constraints*

- every centroid is a **convex combination** of data points

$$\mu_i = \mathbf{X} \mathbf{s}_i \quad \Leftrightarrow \quad \mathbf{M} = \mathbf{X} \mathbf{S}$$

- \mathbf{S} is **column stochastic** (\mathbf{s}_i contains 0s and n_i values of $\frac{1}{n_i}$)
- \mathbf{R} is **column stochastic** and **binary** (\mathbf{r}_j contains a single 1)

constraints on S :

$$s_{ji} \geq 0 \Leftrightarrow \mathbf{S} \succeq \mathbf{0}$$

$$\sum_{j=1}^n s_{ji} = 1 \Leftrightarrow \mathbf{1}^T \mathbf{s}_i = 1$$

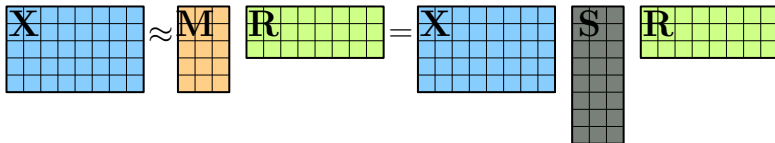
$$s_{ji} \in \left\{0, \frac{1}{n_i}\right\} \Leftrightarrow H(\mathbf{s}_i) = -\sum_{j=1}^n s_{ji} \log s_{ji} \gg 0$$

constraints on R :

$$r_{ij} \geq 0 \Leftrightarrow \mathbf{R} \succeq \mathbf{0}$$

$$\sum_{i=1}^k r_{ij} = 1 \Leftrightarrow \mathbf{1}^T \mathbf{r}_j = 1$$

$$r_{ij} \in \{0, 1\} \Leftrightarrow H(\mathbf{r}_j) = -\sum_{i=1}^k r_{ij} \log r_{ij} = 0$$



k-Means Clustering vs. archetypal analysis

data (sample)



k-means



AA



The Official Creebboby Comics Archetype Times Table

	Robot	Zombie	Astronaut	Monster	Lincoln	Vampire	T.Rex	Ninja	Alien	Platypus
Robot										
Zombie										
Astronaut										
Monster										
Lincoln										
Vampire										
T.Rex										
Ninja										
Alien										
Platypus										

Jacob Borshard 2009

Archetypes

- **Plato:**

ideals and/or pure forms that embody fundamental characteristics of a thing rather than its specific peculiarities

- **C.G. Jung:**

innate, universal forms (the *hero*, the *great mother*, the *wise old man*, ...) that channel experiences and emotions, resulting in recognizable and typical behaviors with certain probable outcomes

- **A. Cutler and L. Breiman (in Technometrics 36(4), 1994):**

archetypal analysis ⇔ new way of data analysis for multivariate data

Archetypal Analysis what?

setting:

- assume a data matrix

$$\mathbf{X} = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_n] \in \mathbb{R}^{m \times n}$$

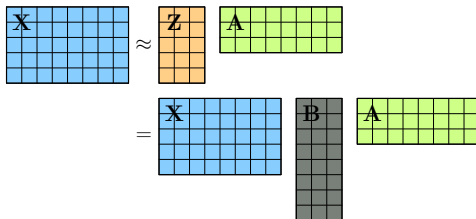
- choose an integer $k \leq \min \{m, n\}$
- determine two column stochastic factor matrices

$$\mathbf{B} \in \mathbb{R}^{n \times k}$$

$$\mathbf{A} \in \mathbb{R}^{k \times n}$$

such that

$$\begin{aligned}\mathbf{X} &\approx \mathbf{Z}\mathbf{A} \\ &= \mathbf{X}\mathbf{B}\mathbf{A}\end{aligned}$$

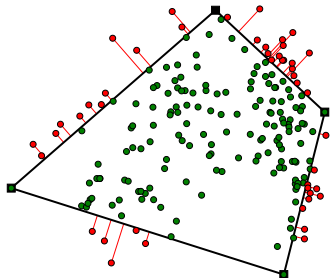


Archetypal Analysis what?

approach:

- solve

$$\begin{aligned} \min_{B,A} \quad & \|X - XBA\|^2 \\ \text{s.t.} \quad & B \succeq 0, \mathbf{1}^T b_i = 1 \\ & A \succeq 0, \mathbf{1}^T a_j = 1 \end{aligned}$$



observe:

- substituting $Z = XB \in \mathbb{R}^{m \times k}$ yields

$$z_i = \sum_{j=1}^n x_j b_{ji} \quad \text{and} \quad x_j \approx \sum_{i=1}^k z_i a_{ij}$$

Archetypal Analysis why?



therefore:

- the **archetypes** \mathbf{z}_i are sparse convex combinations of the data \mathbf{x}_j
- the data points \mathbf{x}_j are sparse convex combinations of the archetypes \mathbf{z}_i

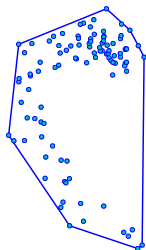


The four basic personality types

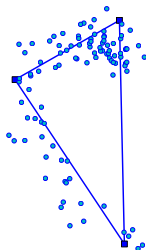
Archetypal Analysis why?

properties:

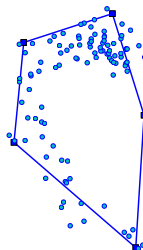
- for $k > 1$, archetypes reside on the data convex hull
- for increasing k , the archetypal hull approximates the data hull



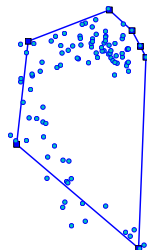
convex hull



$k = 3$



$k = 5$

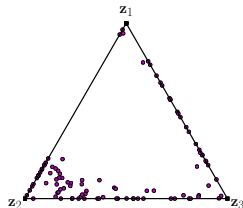
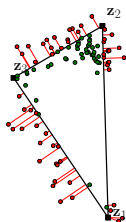


$k = 7$

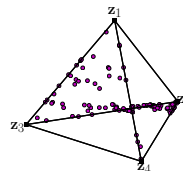
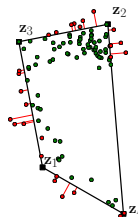
Archetypal Analysis why?

properties:

- the coefficient vectors \mathbf{a}_j in $\mathbf{x}_j = \mathbf{Z} \mathbf{a}_j$ are *stochastic*
- \Rightarrow the coefficients a_{ij} can be thought of as $p(\mathbf{x}_j \mid \mathbf{z}_i)$
- \Rightarrow (soft)clustering, classification, ranking, ...



$k = 3$

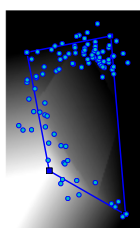


$k = 4$

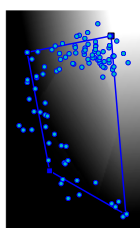
Archetypal Analysis why?

properties:

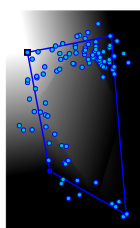
- archetypal analysis bridges geometry and probability



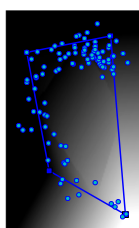
$$p(\mathbf{x} \mid \mathbf{z}_1)$$



$$p(\mathbf{x} \mid \mathbf{z}_2)$$



$$p(\mathbf{x} \mid \mathbf{z}_3)$$



$$p(\mathbf{x} \mid \mathbf{z}_4)$$

Archetypal Analysis how?



alternating least squares:

randomly initialize B

repeat until convergence

compute $Z = XB$

solve under constraints

$$A = \operatorname{argmin}_A \|X - ZA\|^2$$

compute $\hat{Z} = XA^\dagger$

solve under constraints

$$B = \operatorname{argmin}_B \|\hat{Z} - XB\|^2$$

projected gradients:

$$\|X - XBA\|^2 = \operatorname{tr} [X^T X - 2X^T XBA + A^T B^T X^T XBA]$$

$$\frac{\partial E}{\partial A} = 2 [B^T X^T XBA - B^T X^T X]$$

$$\frac{\partial E}{\partial B} = 2 [X^T XBA A^T - X^T X A^T]$$

randomly initialize A and B

repeat until convergence

$$A \leftarrow A - \eta_A \frac{\partial E}{\partial A}, \text{ project } a_j \text{ onto } \Delta^{k-1}$$

$$B \leftarrow B - \eta_B \frac{\partial E}{\partial B}, \text{ project } b_i \text{ onto } \Delta^{n-1}$$

Archetypal Analysis how?



note:

- archetypal analysis makes no implicit density assumption
- archetypal analysis is NP-hard
- both algorithms are heuristics
- both (implicitly) depend on $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{n \times n}$
 - ⇒ they scale with at least $O(n^2)$

⇒ both do not apply to BIG DATA

Archetypal Analysis how?

three key observations:

- 1 $\mathbf{Z} = \mathbf{XB}$ comprises information as to *extreme* data points
- 2 extreme data points are surprisingly easy to determine
- 3 if \mathbf{Z} were known, then

$$\|\mathbf{X} - \mathbf{ZA}\|^2 = \|\mathbf{x}_1 - \mathbf{Za}_1\|^2 + \|\mathbf{x}_2 - \mathbf{Za}_2\|^2 + \dots$$

Archetypal Analysis how?

three key observations:

- 1 $\mathbf{Z} = \mathbf{XB}$ comprises information as to *extreme* data points
- 2 extreme data points are surprisingly easy to determine
- 3 if \mathbf{Z} were known, then

$$\|\mathbf{X} - \mathbf{ZA}\|^2 = \|\mathbf{x}_1 - \mathbf{Za}_1\|^2 + \|\mathbf{x}_2 - \mathbf{Za}_2\|^2 + \dots$$

two key ideas:

- 1 perform AA on an appropriate subset
- 2 decouple the computation of \mathbf{Z} and \mathbf{A}



Making Archetypal Analysis practical how?

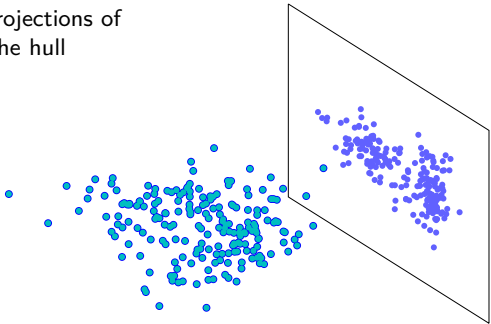
next: brief summary of

- C. Bauckhage and C. Thureau: Making Archetypal Analysis Practical, DAGM, 2009.
- C. Thureau, K. Kersting, and C. Bauckhage: Yes we can: simplex volume maximization for descriptive web-scale matrix factorization. ACM CIKM, 2010.
- C. Thureau, K. Kersting, M. Wahabzada, and C. Bauckhage: Convex non-negative matrix factorization for massive datasets. Knowledge and Information Systems, 29(2):457-478, 2011.
- K. Kersting, C. Bauckhage, C. Thureau, and M. Wahabzada: Matrix Factorization as Search. ECML/PKDD, 2012.

Making Archetypal Analysis practical how?

sampling the convex hull:

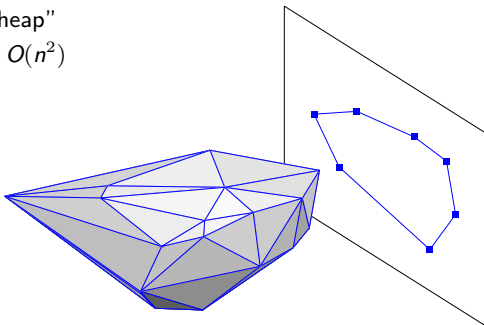
- archetypes are mixtures of points on the data convex hull
- ⇒ restrict algorithm to $X^H \subseteq X$
- in \mathbb{R}^m , convex hull computation is “expensive” ($O(n^{\lfloor m/2 \rfloor + 1})$)
- ⇒ consider (many) 2D projections of the data and sample the hull



Making Archetypal Analysis practical how?

sampling the convex hull:

- if P is a polytope and $\pi : \mathbf{x} \rightarrow \mathbf{M}\mathbf{x} + \mathbf{t}$ is an affine map, then $\pi(P)$ is a polytope
- every vertex of $\pi(P)$ corresponds to a vertex of P
- sampling the hull is “cheap”
- effort is then $O(\eta^2) \ll O(n^2)$
- η is $\Omega(\sqrt{\log n})$



Making Archetypal Analysis practical how?

using distance geometry:

- given the distances d_{ij} between k vertices of an $k - 1$ simplex S , then

$$V(S)^2 = \frac{-1^k}{2^{k-1}((k-1)!)^2} \det(\mathbf{A})$$

where $\det(\mathbf{A})$ is the CM determinant

Cayley-Menger determinant

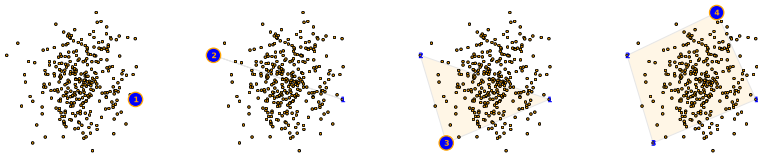
$$\begin{vmatrix} 0 & 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & d_{11}^2 & d_{12}^2 & \dots & d_{1k}^2 \\ 1 & d_{11}^2 & 0 & d_{22}^2 & \dots & d_{2k}^2 \\ 1 & d_{12}^2 & d_{22}^2 & 0 & \dots & d_{3k}^2 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & d_{1k}^2 & d_{2k}^2 & d_{3k}^2 & \dots & 0 \end{vmatrix}$$

Making Archetypal Analysis practical how?

simplex volume maximization (SiVM):

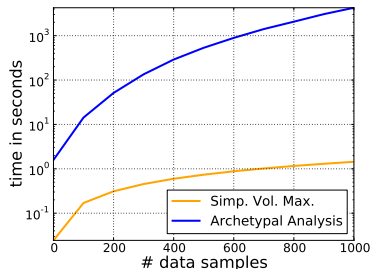
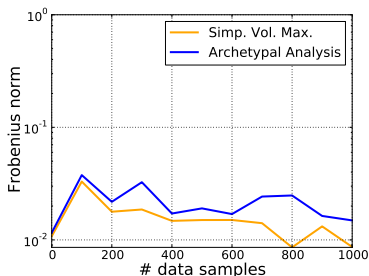
- determine \mathbf{Z} by fitting a simplex of maximal volume into the data
- ⇒ fast, greedy, iterative algorithm
- given S with $k - 1$ vertices, determine next vertex $\mathbf{z}_k \in \mathbf{X}$ s.t.

$$\mathbf{z}_k = \operatorname{argmax}_j V(S \cup \mathbf{x}_j)^2$$



Making Archetypal Analysis practical how?

simplex volume maximization (SiVM):



Applications

Applications mining World of Warcraft

idea / motivation:

- automatically characterize MMORPG teams / guilds
- understand group behavior in MMORPGs
- simplify game design process



Applications mining World of Warcraft

data:

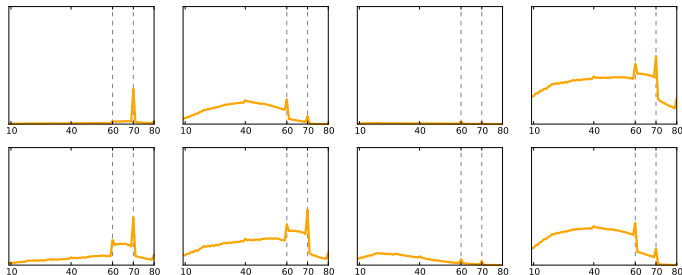
- 150 million in-game observations of
 - 18 million characters
 - in 1.4 million guilds
 - observations reveal a player's experience level (10-80)
 - for each guild, determine a 70 dim. experience histogram
- ⇒ data matrix of 98 million entries

details in:

- C. Thureau and C. Bauckhage: Analyzing the evolution of social groups in World of Warcraft. IEEE CIG, 2010.

Applications mining World of Warcraft

basis vectors obtained from k -means:

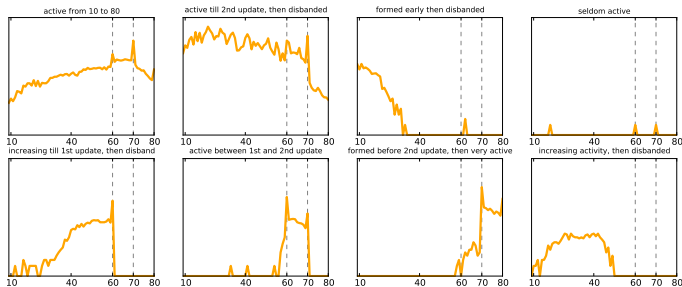


- rather similar entities
- difficult to interpret

Applications mining World of Warcraft



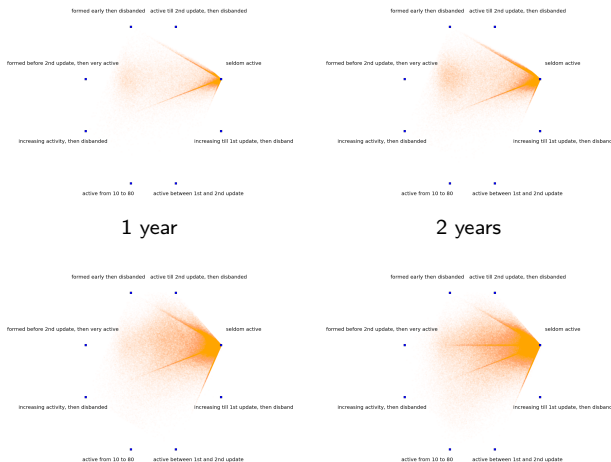
basis vectors resulting from SiVM:



- more meaningful (since more *extreme*) entities
- correspond to existing guilds

Applications mining World of Warcraft

temporal evolution:



Applications

authorship analysis of twitter tweets

idea / motivation:

- do celebrities really tweet themselves?

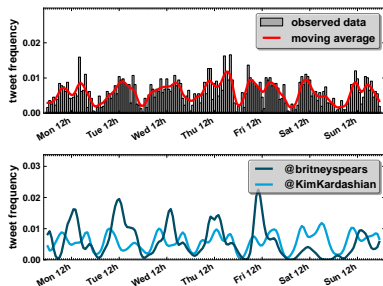
data / features:

- 300 celebrity twitterers
- 1.5 million tweets
- 30 dim. feature vectors

⇒ data matrix of 45 million entries

details in:

- C. Thureau, K. Kersting, and C. Bauckhage: Yes we can: simplex volume maximization for descriptive web-scale matrix factorization. ACM CIKM, 2010.



Applications

authorship analysis of twitter tweets

@aplusk

z ₁	I often ponder as to how similar the name Adam and the word "atom" are to one another	0.871
z ₂	@thesettlers Its R	0.051
z ₃	new blog... http://tinyurl.com/cggp43	0.024
z ₄	http://www.ustream.tv/ashton	0.020
z ₅	@toojiggy what?	0.011
z ₆	@DillonPrat amazing!	0.009
z ₇	HOIY GOOD GOATS MILK WE'VE GOT A QB!!!!!!!!!! THANK YOU BABY JESUS BUDAH ALLAH KABBALAH DARWIN WHOEVER!!!!!!	0.007
z ₈	Victory is ours!!!!!!!!!!	0.007

@britneyspears

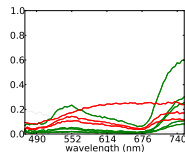
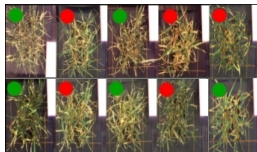
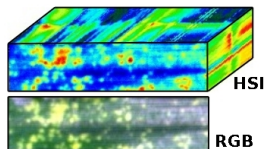
z ₁	Brit's up for a People's Choice Award for Fave Scene Stealing Guest Star from How I Met Your Mother. Please vote: http://tinyurl.com/5ns9n3	0.317
z ₂	Over a million hits in 24 hours on the new http://www.britneyspears.com . You guys really are the best fans in the world!	0.301
z ₃	New leg, new song, new outfits, new do.... -Brit	0.301
z ₄	Have you joined the Britney social network yet? Connect with other Britney fans at http://www.circusvip.com	0.032
z ₅	NYC....Here I come!!! -Brit	0.022
z ₆	HAPPY BIRTHDAY BRETT!!!!!! XO XO -Britney	0.011
z ₇	Happy Holidays from Britney!! http://tinyurl.com/7uyrr4	0.011
z ₈	Thanks @britneylive @filmmaker311 @singlemomindebt @youbetheanchor @perezhilton! Were you there? Share @ http://britneyspears.com/tour	0.005

Applications

drought stress recognition in plants

idea / motivation:

- drought / abiotic stress depreciate crop yields by up to 70%
- understanding the mechanisms of stress resistance is of vital importance; phenomic approaches are called for
- hyperspectral imaging provides an auspicious modality



Applications

drought stress recognition in plants

data / approach:

- 205 daily measurements of healthy and stressed plants
- recorded with Surface Optics Corp. SOC-700
- each image is a $640 \times 640 \times 69$ data cube

⇒ data matrix of 5.8 billion entries

- extract extreme spectra
- fit Dirichlet distributions
- perform Bayesian regression

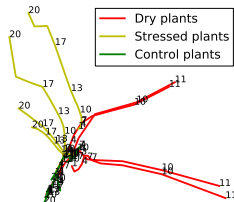
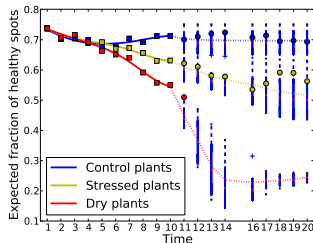
details in:

- C. Römer, M. Wahabzada, A. Ballvora, F. Pinto, M. Rossini, C. Panigada, J. Behmann, J. Leon, C. Thureau, C. Bauckhage, K. Kersting, U. Rascher, and L. Plümer: Early drought stress detection in cereals: simplex volume maximisation for hyperspectral image analysis. *Functional Plant Biology* 39(11), 2012.
- C. Bauckhage, K. Kersting: *Data Mining and Pattern Recognition in Agriculture*. KI 27(4), 2013.

Applications

drought stress recognition in plants

results:



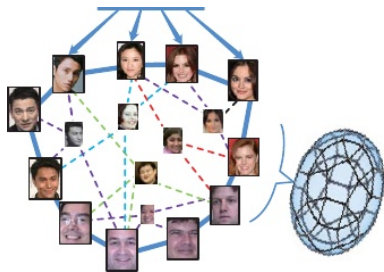
Applications

classification via archetype hulls



results:

- convex archetype regions as class model
- classify against the nearest archetype hull
- straight forward geometric classifier
- well suited for high dimensional data



details in:

- C. Thureau: Nearest archetype hull methods for large-scale data classification. ICPR 2010.
- Yuanjun Xiong, Wei Liu, Deli Zhao, and Xiaoou Tang: Face Recognition via Archetype Hull Ranking. ICCV 2013

Take home message



- k-means is a matrix factorization problem
- it is not always a good choice and lacks interpretability
- archetypal analysis is data analysis with *extremes*
- it yields intuitive and interpretable results
- it implicitly depends on distances
- it thus applies to many scenarios
- it allows for efficient computation

- Prof. Christian Bauckhage
University Bonn and Fraunhofer IAIS
- Prof. Kristian Kersting
TU Dortmund
- Dr. Mirwaes Wahabzada
Fraunhofer IAIS



Thank you for your attention!

`christian.thurau@unbelievable-machine.com`



note:

- k -means clustering implicitly fits a Gaussian mixture model
- k -means clustering is NP hard (Aloise et al., ML, 75(2), 2009)
- introducing indicator variables

$$r_{ij} = \begin{cases} 1, & \text{if } \mathbf{x}_j \in C_i \\ 0, & \text{otherwise} \end{cases}$$

its objective can be cast as

$$E(k) = \sum_{i=1}^k \sum_{j=1}^n r_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

hard k -means clustering

$$\min_{\mathbf{S}, \mathbf{R}} \left\| \mathbf{X} - \mathbf{XSR} \right\|^2 - \sum_{i=1}^k H(\mathbf{s}_i)$$

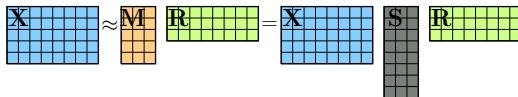
$$\mathbf{S} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{s}_i = 1$$

$$\text{s.t.} \quad \mathbf{R} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{r}_j = 1$$

$$H(\mathbf{r}_j) = 0$$



soft k -means clustering

$$\min_{\mathbf{S}, \mathbf{R}} \left\| \mathbf{X} - \mathbf{XSR} \right\|^2 - \sum_{i=1}^k H(\mathbf{s}_i)$$

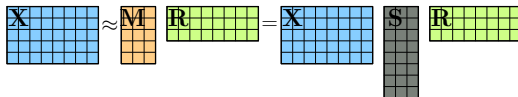
$$\mathbf{S} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{s}_i = 1$$

$$\text{s.t.} \quad \mathbf{R} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{r}_j = 1$$

$$H(\mathbf{r}_j) = 0$$



archetypal analysis

$$\min_{S, R} \left\| \mathbf{X} - \mathbf{XSR} \right\|^2 - \sum_{i=1}^k H(s_i)$$

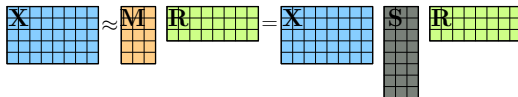
$$\mathbf{S} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{s}_i = 1$$

$$\text{s.t.} \quad \mathbf{R} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{r}_j = 1$$

$$H(r_j) = 0$$



non-negative matrix factorization

$$\min_{S, R} \left\| \mathbf{X} - \mathbf{XSR} \right\|^2 - \sum_{i=1}^k H(s_i)$$

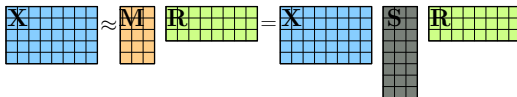
$$\mathbf{S} \succeq \mathbf{0}$$

$$\mathbf{1}^T \mathbf{s}_i = 1$$

$$\text{s.t.} \quad \mathbf{R} \succeq \mathbf{0}$$

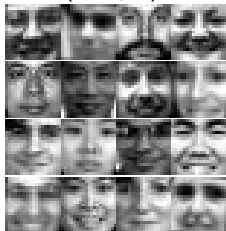
$$\mathbf{1}^T \mathbf{r}_j = 1$$

$$H(r_j) = 0$$



Appendix

data (sample)



SVD



NMF



k-means



AA

