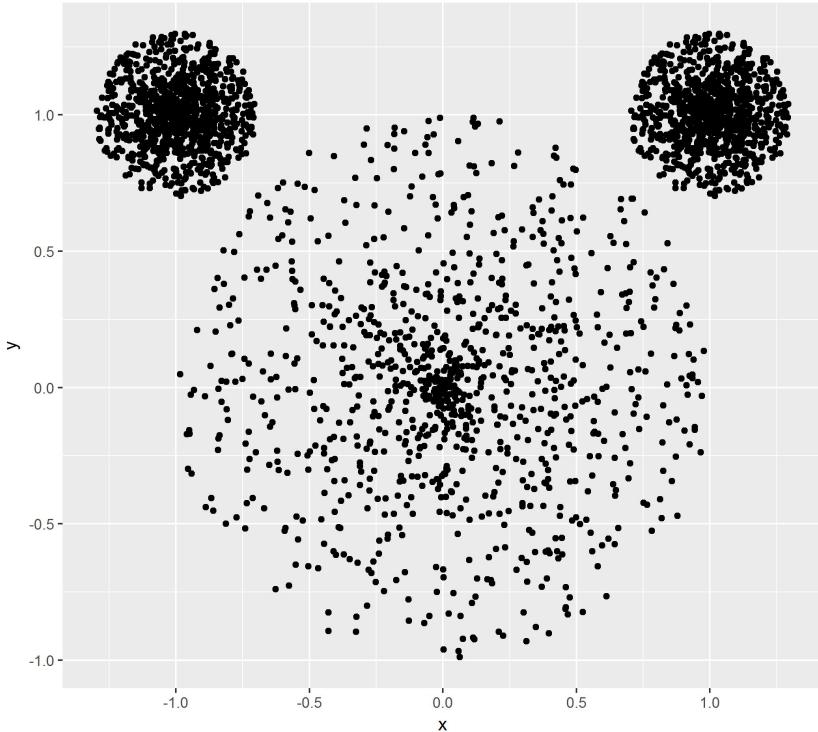
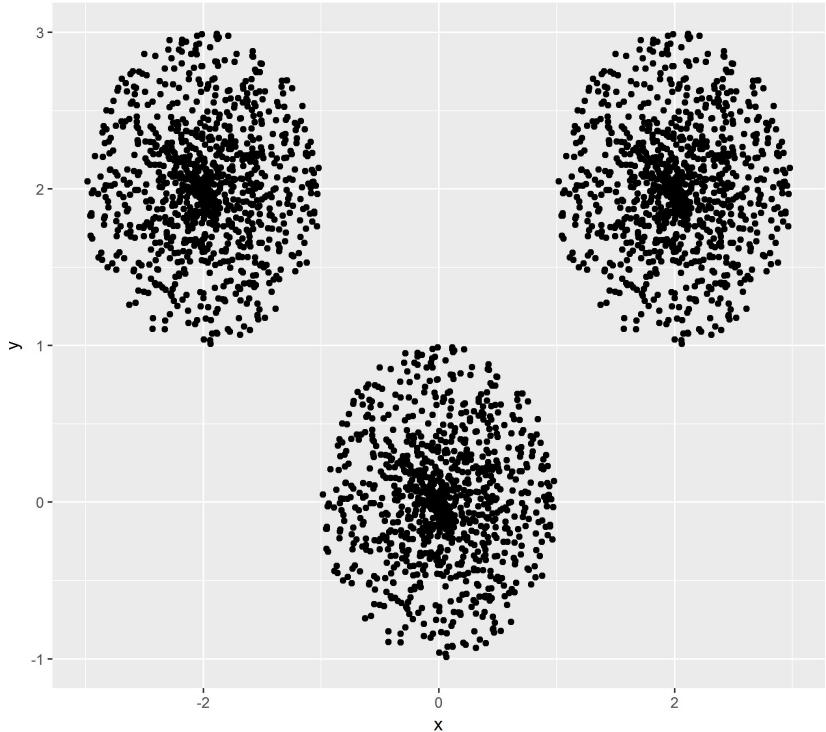


Clustering

(Datos Extensos = High dim. Data?)

Federico Zertuche
DIDE, UTA

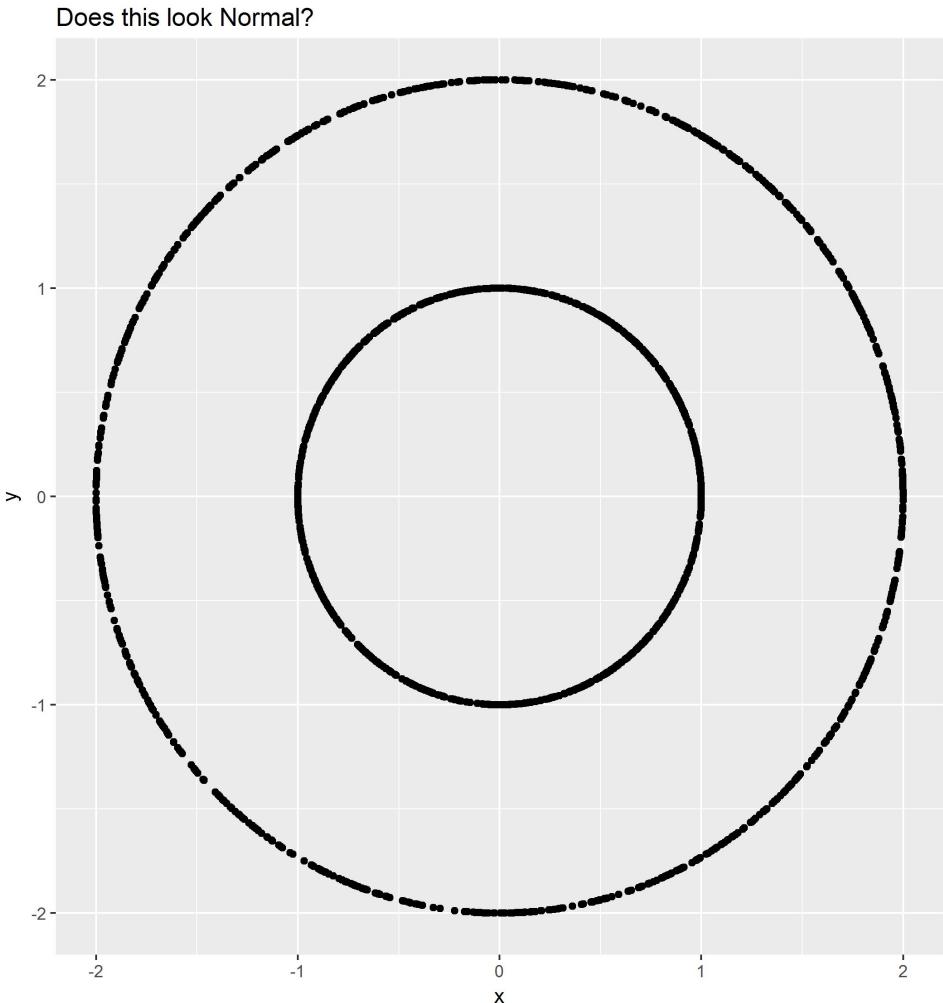
What is a cluster?



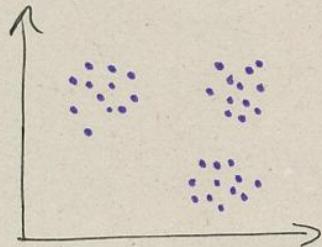
Normal?

(A sample from two Normals in high dimensions)

- Can get unexpected shapes.
- We need a definition (if possible).
- Deal with these shapes (if possible).



K-means



(K is given)

Stats Pt of View

①

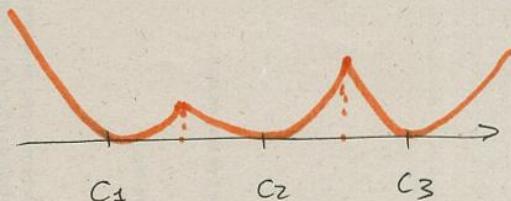
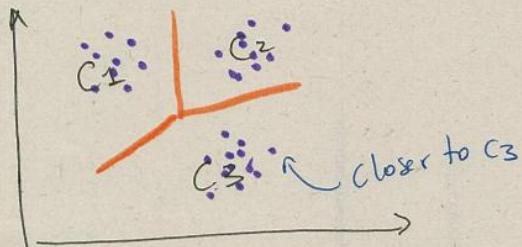
$$x_1, \dots, x_n \sim P$$

Given P :

$$\text{loss} = E \left[\min_j \|x - c_j\|^2 \right]$$

$$R(c) :=$$

- Find $P = \{c_1, c_2, c_3\}$
- Form Voronoi tessellation



minimize $R(c)$

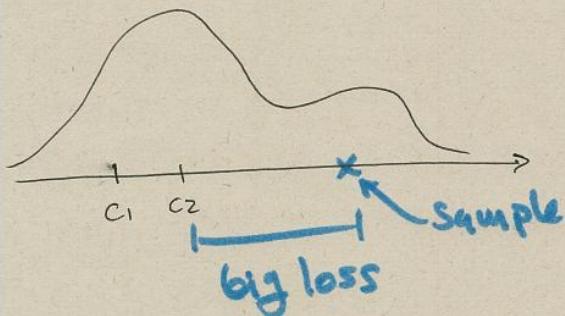
c

Non-convex optimization problem
(T.. Part NP-Hard)

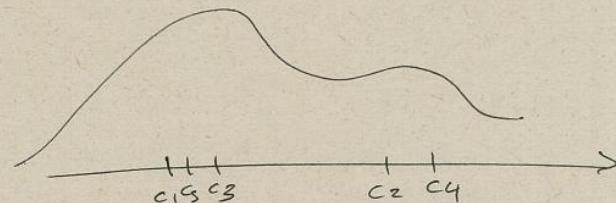
k-means

Guessing the answer

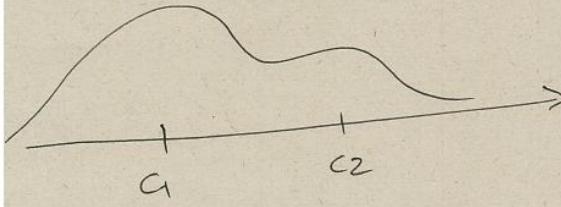
(2)



5 c_j 's:



$R(c)$ decreases with # of c_j 's.



[3]

Lloyd's Algo

Good Approx:

(Lloyd's algo)

① Choose c_1, \dots, c_k .

② Form the clusters.

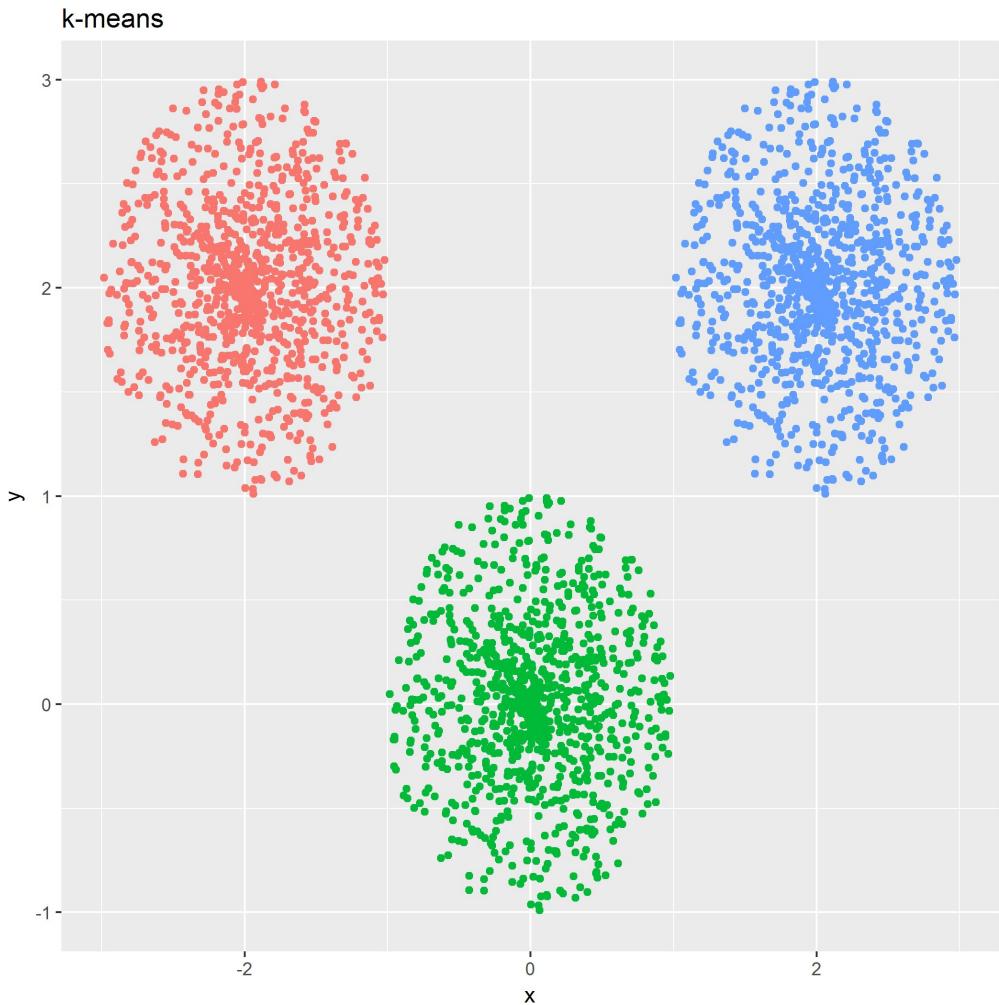
$$\text{compute } d_j(x_k) = \|x_k - c_j\|^2$$

③ Reestimate centers:

$$c_j \leftarrow \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

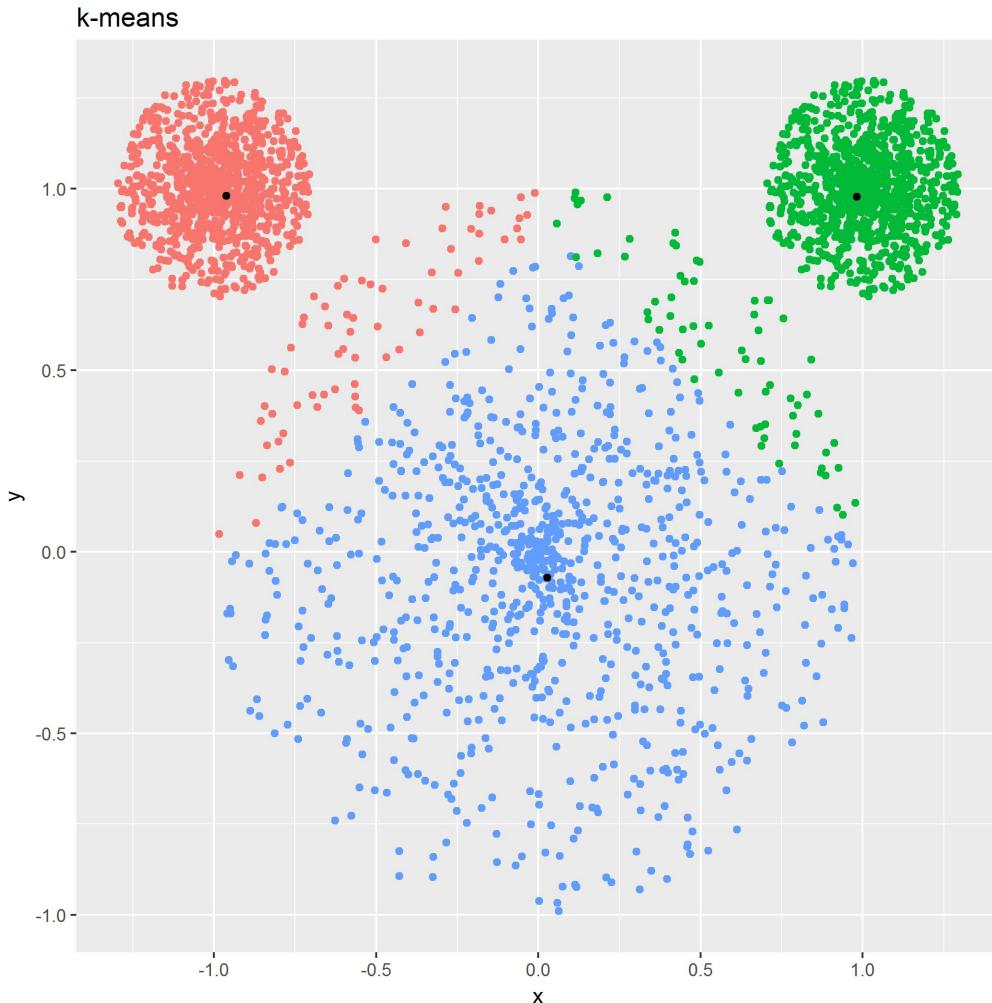
Repeat until
convergence.
(local min)

Mickey Mouse Example

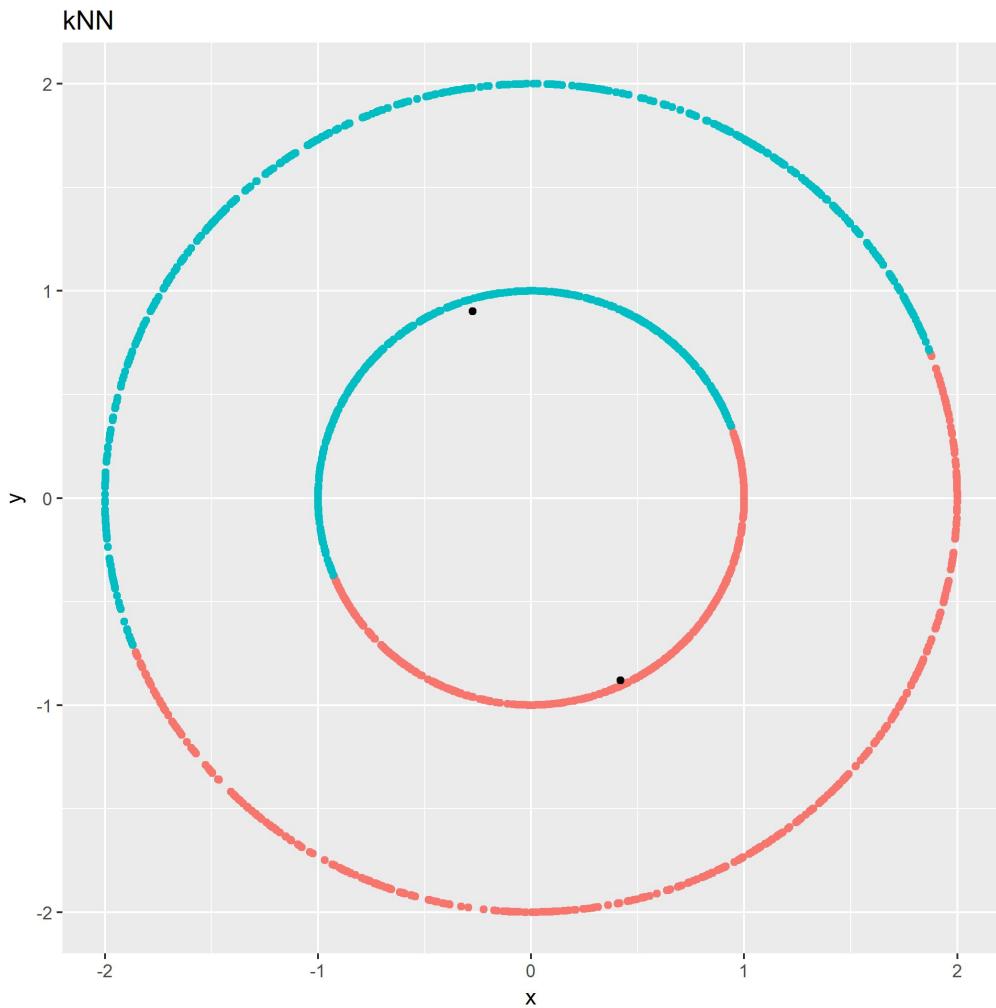


Mickey Mouse Example

- k-means is supposed to do this.
- k is NOT the # of clusters.
- It makes sense to divide the density as in the fig (otherwise big loss).

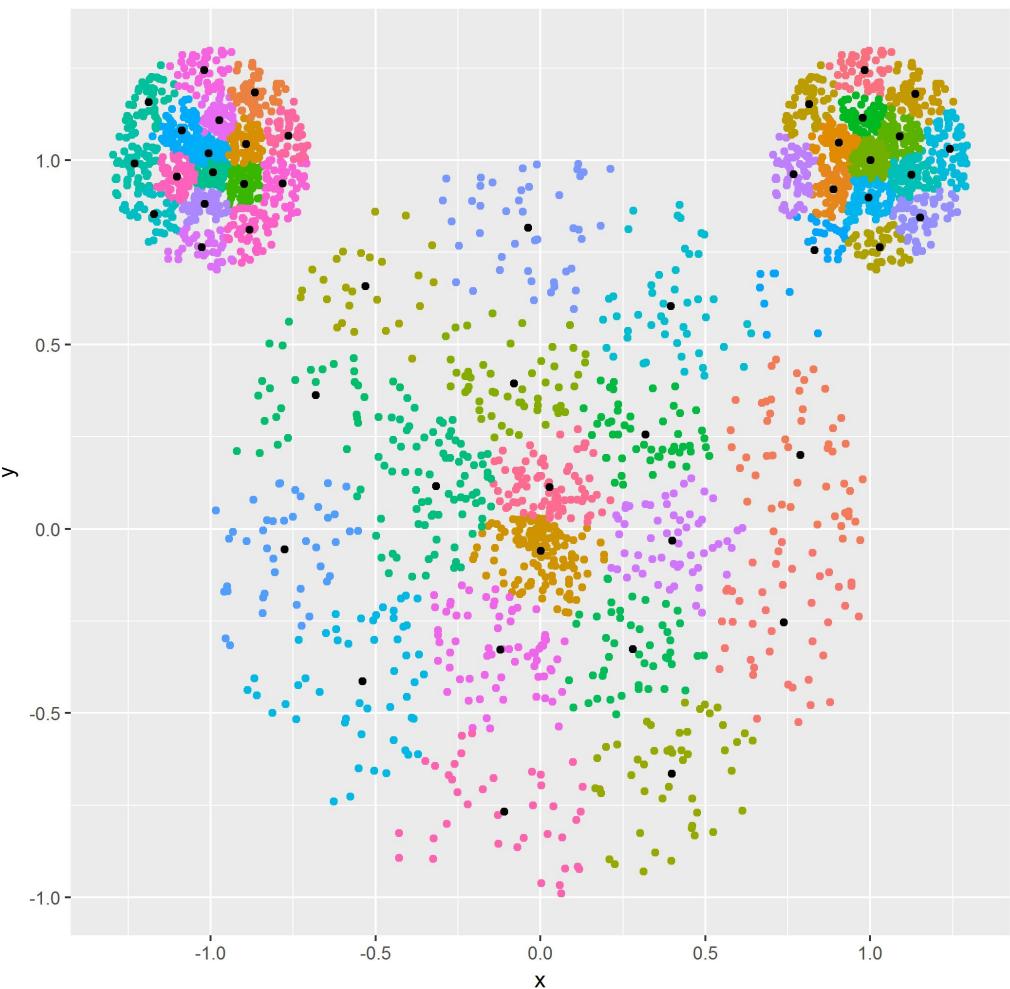


Two Rings Example



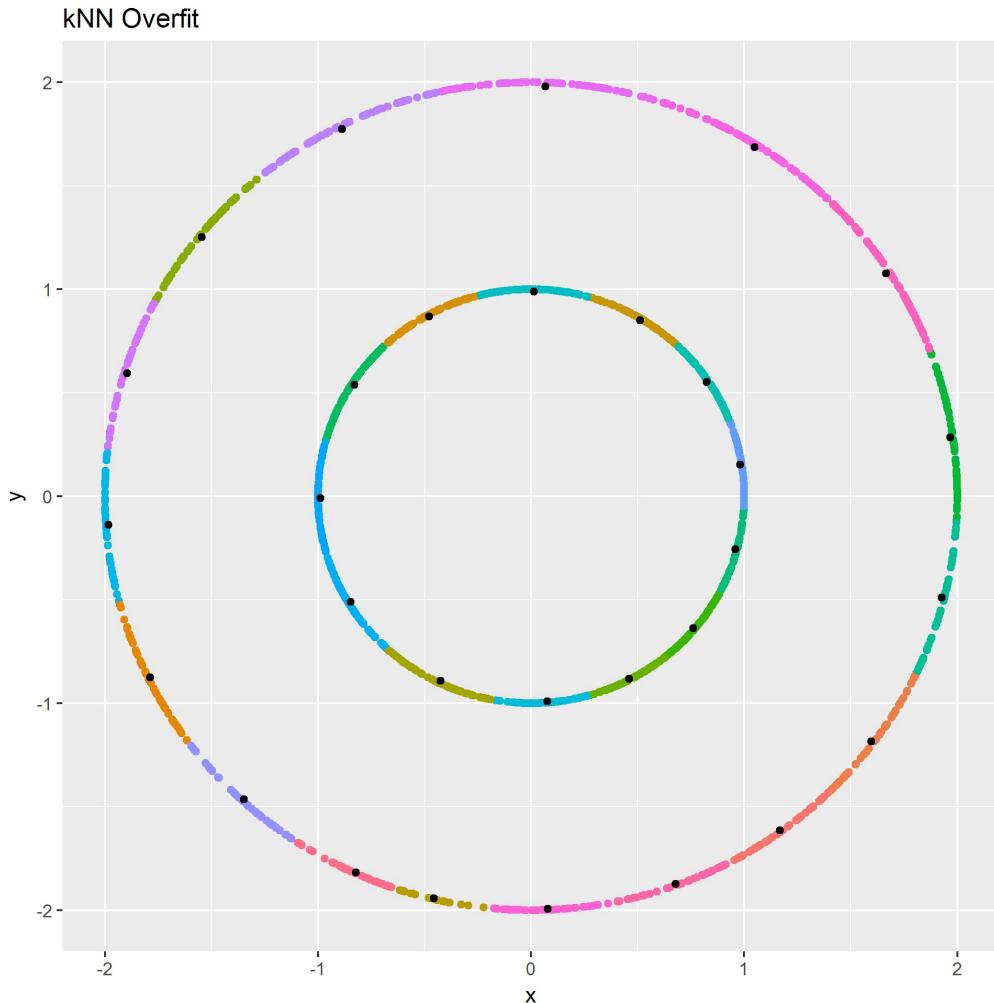
Overfit and Merge (L. Wasserman)

kNN Overfit



Overfit and Merge (L. Wasserman)

- Open problem:
 - Define this proc.
 - $k = n^{1/2}$?
 - Merging not clear.



Pros and Cons

- Easy to implement.
 - Easy to modify.
 - Seems to work ok in practice.
-
- Heuristic (local min).
 - No guarantees.
 - Not a clear notion of cluster.

Mixture Models (~Generalized-Bayesian k-means)

- Generative point of view.
- Data comes from a mixture of Gaussians.
 - Many parameters $O(kd^2)$
 - E.M. or Gibbs sampling.

Many unknown things:

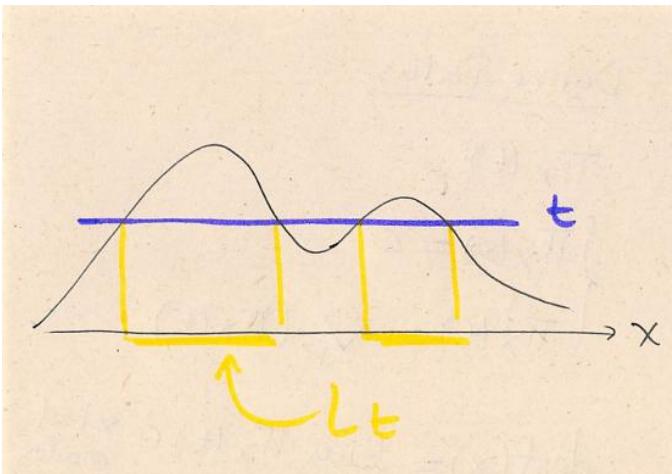
- Identifiability, infinite posteriors, non-intuitive groups, etc.
- Use with -some- caution.

Density-Based Clustering

1. Estimate the density.
2. Use the density to cluster.

Problems in high dimensions when estim. the density.

Level-Set Clustering



Density tree

$$C = \bigcup_{t \geq 0} C_t$$

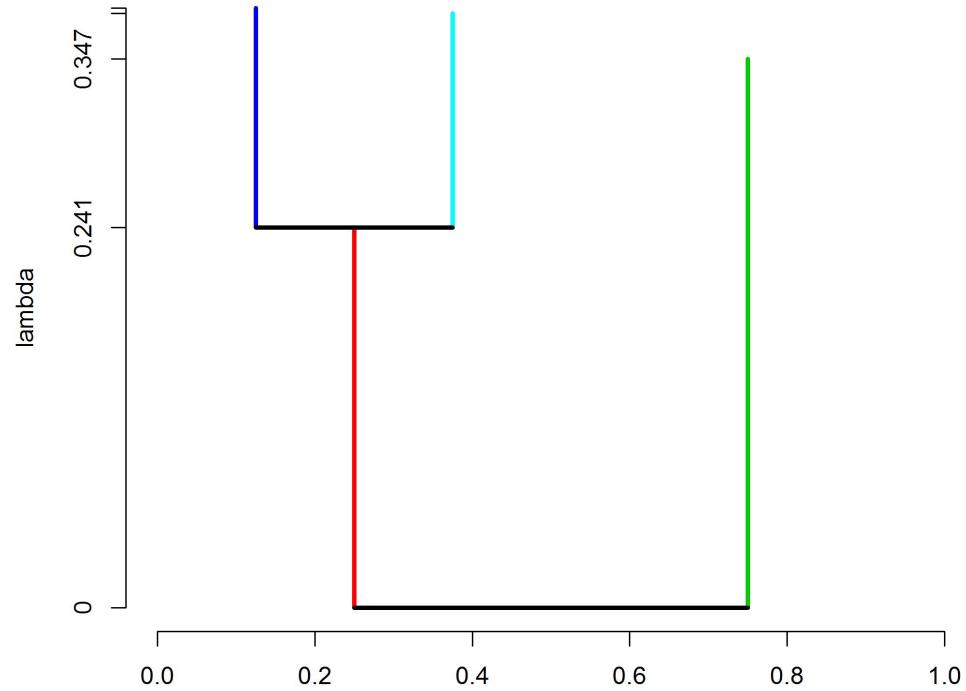
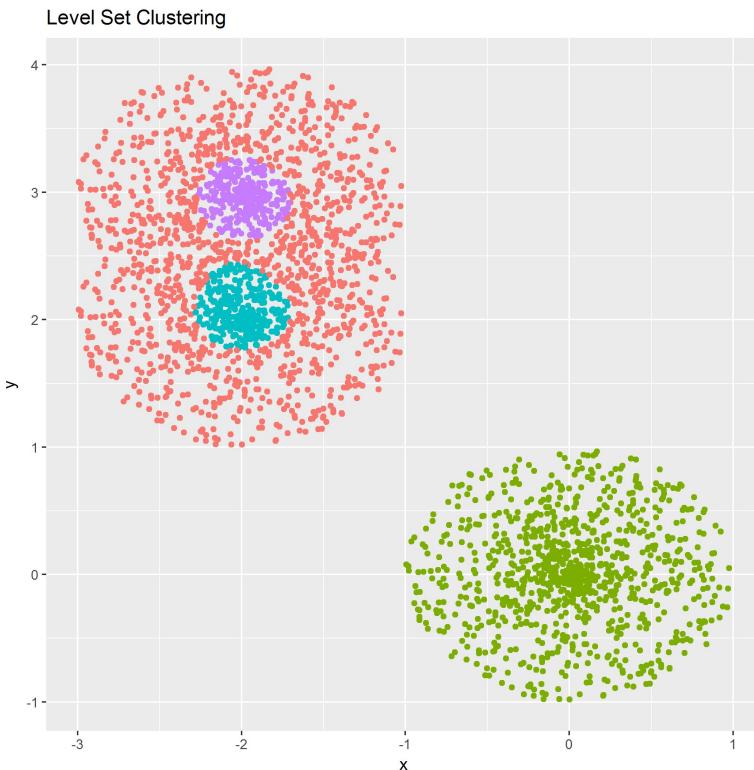
is a tree.

Goal: Build the tree.

$$L_t := \{x : p(x) > t\} = \bigcup_j C_j$$

- High density region is an intuitive notion of cluster

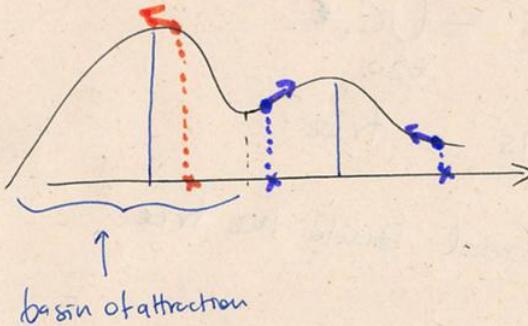
Level-Set Clustering



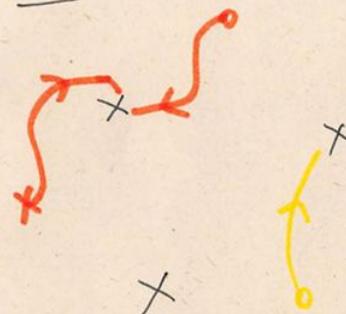
Mean-Shift Clustering

Mode = local max.

1D:



2D:



Define Paths:

$$\Pi_X(t)$$

$$\left\{ \begin{array}{l} \Pi_X(0) = 0 \\ \end{array} \right.$$

$$\Pi_X'(t) = \nabla_p (\Pi_X(t))$$

$$\text{dest}(x) = \lim_{t \uparrow \infty} \Pi_X(t) \in \underset{\text{set of modes}}{\text{set of}}$$

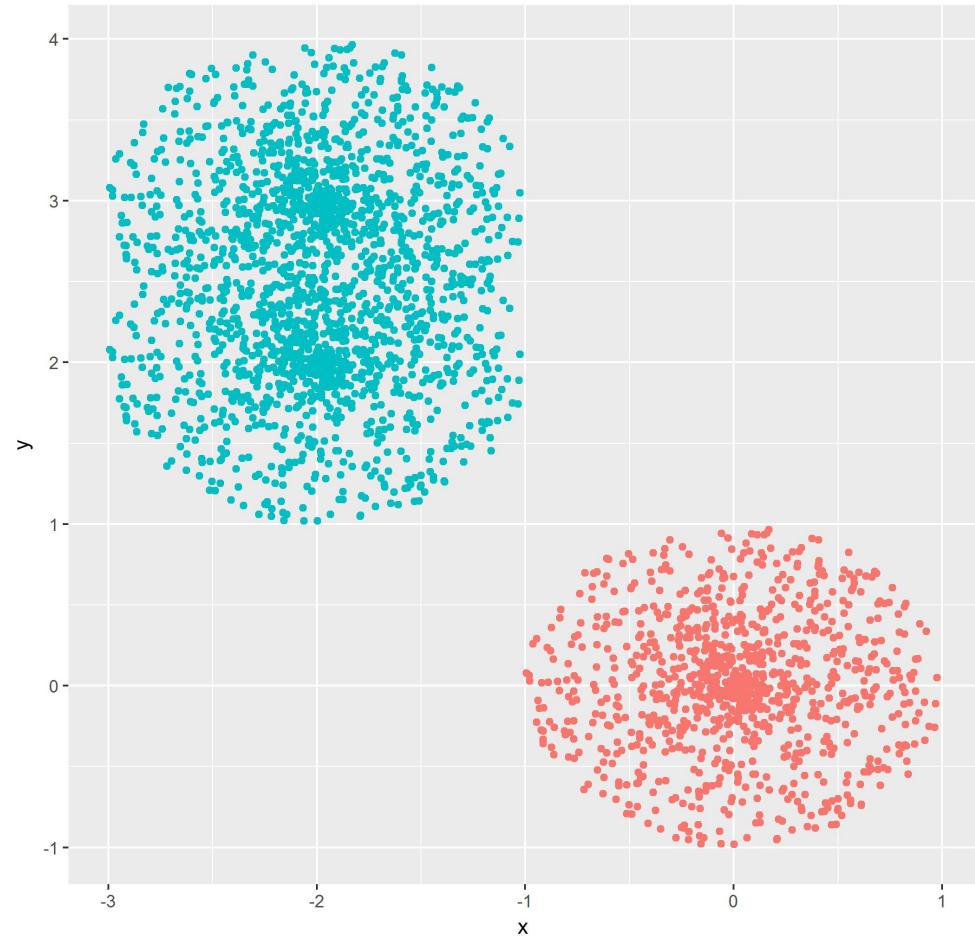
MS-Algo

① Pick any point $a = a^{(0)}$.

$$a^{(1)} \leftarrow \frac{\sum_{i=1}^n x_i K \left(\frac{x_i - a^{(0)}}{h} \right)}{\sum_{j=1}^n K \left(\frac{x_j - a^{(0)}}{h} \right)}$$

Mean-Shift Clustering

Mean Shift Clustering



Other Clustering Techniques

- Hierarchical clustering.
 - Not clear what it does.
 - Inconsistent.
- Spectral Clustering.
 - Cool math.
 - Good for graphs.

Spectral Clustering

Graph Laplacian

(~ derivative Operator for graphs)



$$w_{i,j} = \begin{cases} 1 & \text{if } \|x_i - x_j\| < \varepsilon, \\ 0 & \text{otherwise.} \end{cases}$$

$$L = D + W$$

where

$$D[i,j] = \sum_j w_{i,j}$$

$$f := (f_1, \dots, f_n)$$

$$f^T L f = \frac{1}{2} \sum_{i,j} w_{i,j} (f_i - f_j)^2$$

↑
small weights to pts.
far appart.

Differences on the graph.

$$\text{eigen}(L)$$

- $0 = \lambda_1 < \lambda_2 < \dots < \lambda_n$
- $|\lambda_j = 0| = \# \text{ of connected components}$

Spectral Clustering

Usually you need to:

- Find the graph

$$-\|x_i - x_j\|^2 / 2h^2$$

$$w_{ij} = e$$

- Def-

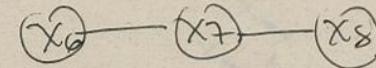
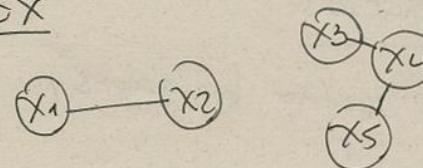
$$L = \bar{D} W$$



like a kernel
density estimation.

- Pick the first S eigenvectors.
& run k-means clustering
in the subspace

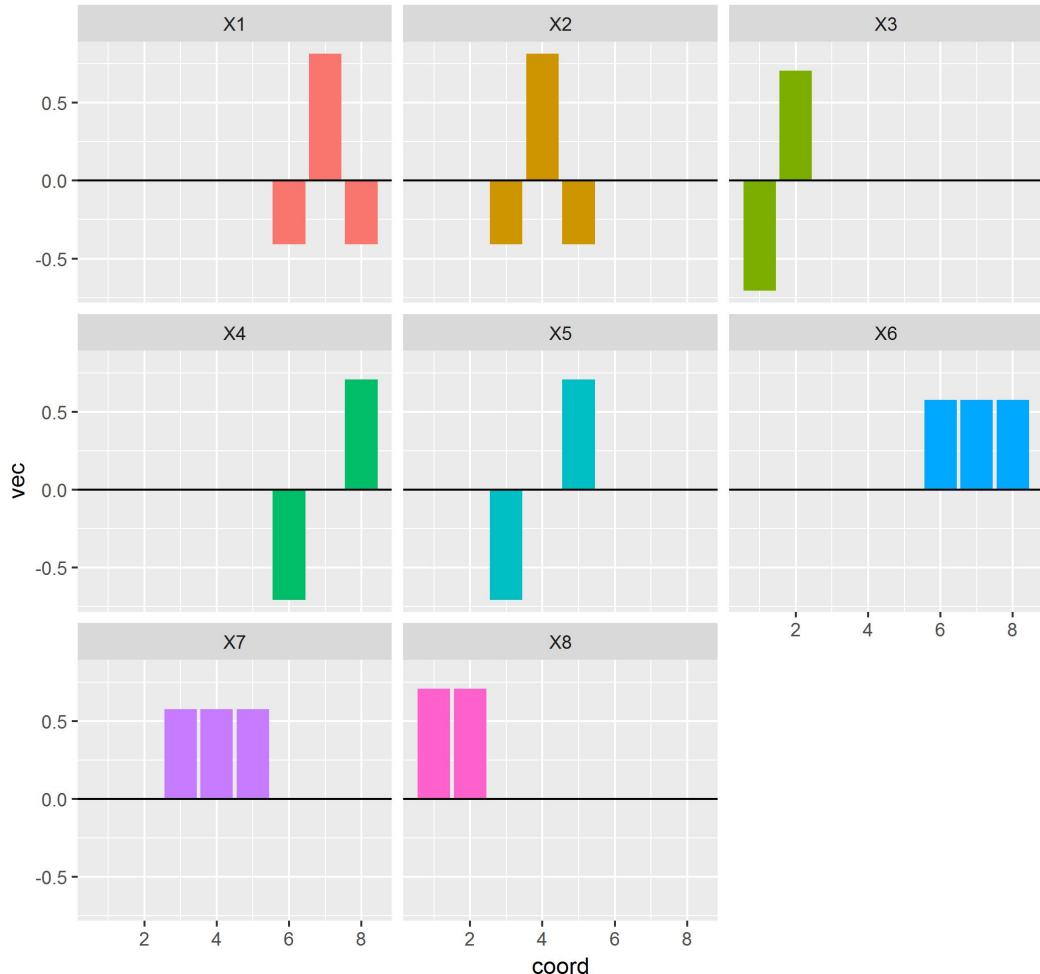
Ex



$$L = \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Spectral Clustering

Eigenvectors of the graph Laplacian
 $e_n > \dots > e_1$



High Dim. Clustering

$$X \sim N(0, I) ; X \in \mathbb{R}^d$$

Not hard to show that:

$$\|x\|^2 = d$$

i.e. points lie on rings.



One fix:
Modify the distance

$$p(x, y) := \frac{1}{n-2} \sum_{z \neq x, y} |\|x-z\| - \|y-z\||$$

[Sarkar, & Ghosh]

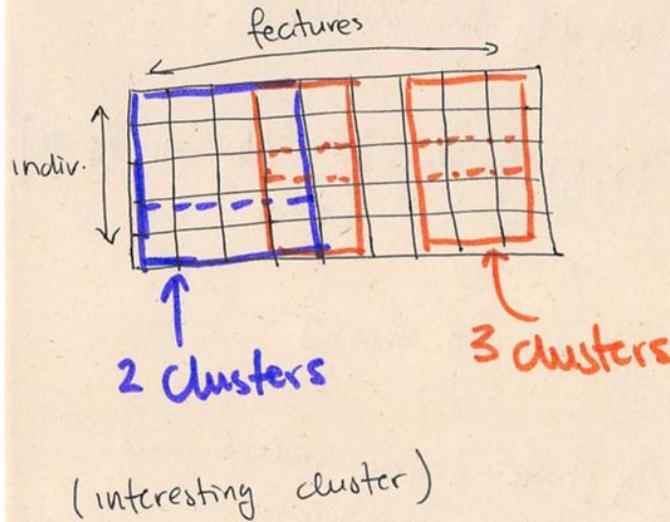
⚠ Proofs assume independence between coordinates --

Open problem:

Adjust k-means to deal with high dim. effects

Mosaics

Find Features that define clusters.



A First Idea:

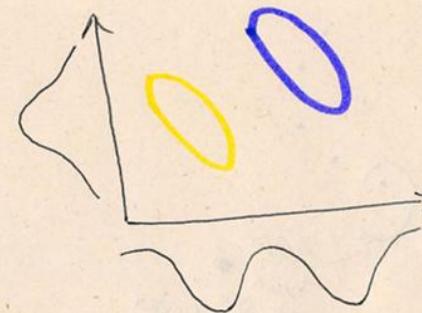
Screening: Are there marginal clusters?

$$X = (X[1], \dots, X[n]).$$

$$X_1[1], \dots, X_m[n] \sim F$$

H_0 : F is unimodal.

H_1 : F is not.



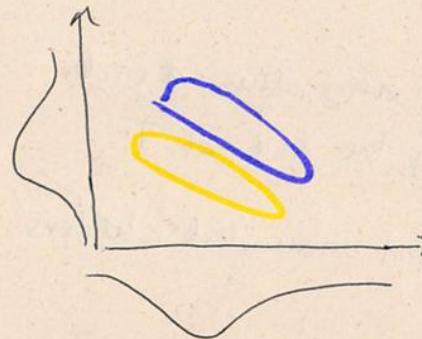
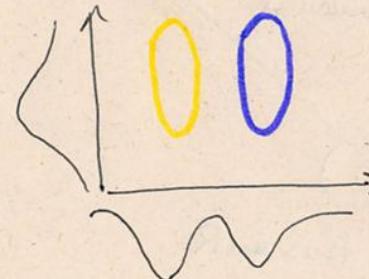
Screening

$$D_{IP}(F) = \inf_{G \text{ unimodal}} \sup_x |F(x) - G(x)|$$

$$T := D_{IP}(\hat{F}_n)$$

↑
can compute this in R

- Do this in every dim. & order the T's.
- Select least unimodal?



Alternating Sparse Clustering

$$R_n = \frac{1}{n} \sum_{j=1}^n \min_c \|x_j - c\|^2$$
$$= \dots =$$
$$= \sum_{j=1}^k \frac{1}{|C_j|} \sum_{i \in C_j} \|x_i - x_j\|^2$$

Choose $|S| = L$.

$a = 1, \dots, d$ index of features

$$\delta_a(i, j) := (x_i[a] - x_j[a])^2$$

Rescale so that $\sum_{i,j} \delta(i, j) = 1$

Do k-means with an alternative step:

- Find the centers.

- Find the features



+ minimize R_n to find the features.

A Model for Mosaics?

Each object is represented by
a vect. of latent features f .

$$f_i \mapsto x_i$$

properties are generated
by the latent features

N objects ; $F = \begin{bmatrix} f_1^T \\ \vdots \\ f_N^T \end{bmatrix}$

$$f_{ij}[j] = \begin{cases} 1 & \text{if } x_i \text{ has prop. } j, \\ 0 & \text{otherwise.} \end{cases}$$

Idea:

Def. a a priori dist $P(F)$,
compute $P(X|F)$.

Finite-dim model

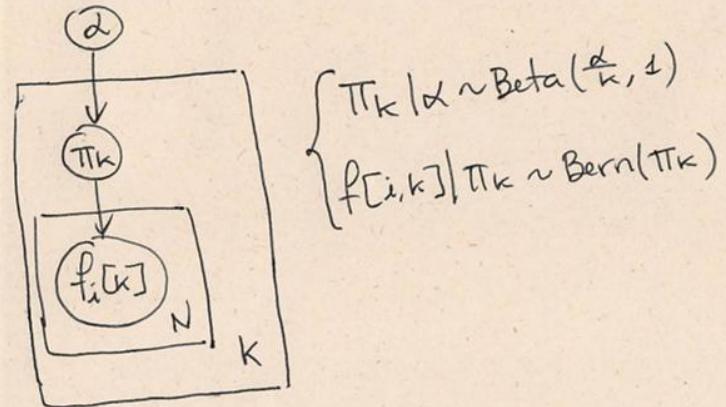
$$P(F|\pi) = \prod_{j=1}^K \pi_j^{m_j} (1-\pi_j)^{N-m_j}$$

m_j = # of obj's possessing feature j .

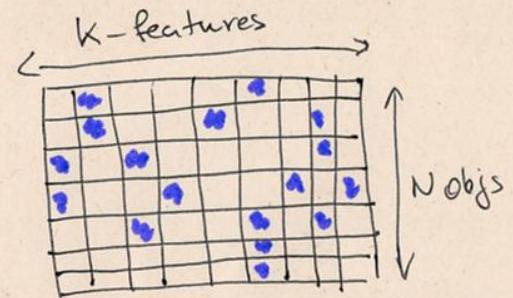
$$\pi_j = P(\text{possessing feature } j)$$

$$P([100100]) = \pi_1^5 (1-\pi_1)^1 + \dots + \pi_6^5 (1-\pi_6)^1$$

A Model for Mosaics?

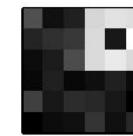
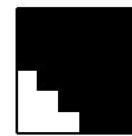
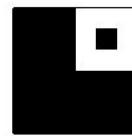
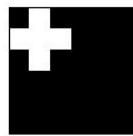


Given the model :
 $P(\text{obs} | F) \propto P(F | \text{obs}) P(F)$

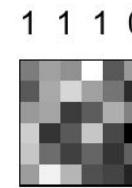
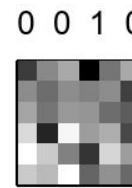
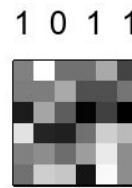
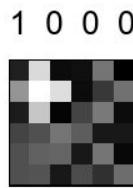


A Mosaic Model

(a)



(b)



A Mosaic Model

(a)



(Positive)



(Negative)



(Negative)



(Negative)

(b)



0 0 0 0



0 1 0 0



1 1 1 0



1 0 1 1

(c)



Some references:

1. *Statistical Machine Learning*, Ryan Tibshirani, Larry Wasserman.
2. *Functional Data Clustering: a Survey*, Julien Jacques, Cristian Preda.
3. *A Simple Approach to Sparse Clustering*, Ery Arias-Castro, Xiao Pu.
4. *On Perfect Clustering of High Dimension, Low sample Size data*, Soham Sarkar, Anil K. Ghosh.
5. *Tree-Structured Stick Breaking Process for Hierarchical Data*, Ryan P. Adams, Zoubin Ghahramani and Michael I. Jordan
6. *The Indian Buffet Process: An Introduction and Review*, Thomas Griffiths and Zoubin Ghahramani.

G		A	C	I	A	
---	--	---	---	---	---	--

