Introduction to Machine Learning. CSCI-UA 9473, Lecture 7.

Augustin Cosse

#### Ecole Normale Supérieure, DMA & NYU Fondation Sciences Mathématiques de Paris.



2018

## Supervised Learning (I)

- Linear regression
  - Bias variance trade-off (Linear and non linear data)
  - Regularization (Ridge, Lasso, Subset Selection)
- Linear classification
  - Separating hyperplane, LDA, logistic regression
  - Perceptron
  - Discriminative vs Generative classifiers
- Non parametric regression/classification
  - Kernel methods
  - Support vector machines
- Neural Networks, convolutional neural networks

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

## Today: Unsupervised Learning

- So far : predictions based on training samples for which joint values {(𝑥<sub>i</sub>, y<sub>i</sub>)}<sup>N</sup><sub>i=1</sub> are known.
- Problem: costly. Most datasets are not labeled.
- Today: Unsupervised learning = learning without a teacher
- ► In unsupervised Learning we are given samples (x<sub>1</sub>, x<sub>2</sub>,..., x<sub>N</sub>) from a distribution P(X) and the goal is to infer the properties of the distribution without the help of the teacher.
- One approach is to characterize the values X for which P(X) is large. This includes a variety of methods which attempt to locate low dimensional manifolds within the X space

## Samsung AI Forum Offers a Roadmap for the Future of AI

on September 18, 2018

#### **Unsupervised Learning Takes Center Stage**



LeCun used training self-driving cars as a key example of unsupervised learning's potential. "A lot of people who are working on autonomous driving are hoping to use reinforcement learning to get cars to learn to drive by themselves by trial and error," said LeCun. "The problem with this is that, because of [reinforcement learning's inherent inefficiencies], you'd have to get a car to drive off a cliff

## Science Home News Journals Topics Careers

# What artificial brains can teach us about how our real brains learn

By Matthew Hutson | Sep. 29, 2017, 3:10 PM

#### Q: Why focus on unsupervised learning, which is much less common in Al?

A: With supervised learning, you are assuming that you have a teacher providing the correct label at each learning event. Think about how we humans learn. This very rarely happens.

#### Alexa is Now Even Smarter-



naturally; and used active learning and unsupervised learning to improve foundational wake word detection, speech recognition, and natural language understanding," said Rohit Prasad, Vice President and Head Scientist, Amazon Alexa. "We've only scratched the surface of A.I.-powered inventions and we'll continue to invent ways to make Alexa more useful for our customers."

#### several thousand times before it figures out 7 factors that will push implementation of AI in healthcare

"As new unsupervised learning algorithms are discovered, the data efficiency of deep learning will be greatly augmented in the years ahead, and its potential applications in healthcare and other fields will increase rapidly," according to Hinton.

August 31, 2018 Melissa Rohman Artificial Intelligence

## Cortica Will Apply 'Unsupervised Learning' AI Tech to Help Self-Driving Cars Get Smarter

This tech company wants autonomous cars to figure things out on their own.

BY STEPHEN EDELSTEIN JULY 13, 2017

ARTIFICIAL INTELLIGENCE

TECH

AUTONOMOUS CARS SELF-DRIVING CARS





# The Missing Link of Artificial Intelligence

We don't know how to make software that learns without explicit instruction—but we need to if dreams of humanlike Al are to come true.

Chnology by Tom Simonite February 18, 2016

THE ORIVE MIT Technology Review

The Download Magazine Events More+

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

## Today: Unsupervised Learning

Finding lower dimensional representation of the data gives an intuition on the association among the variables and whether the variables can be considered as functions of a smaller set of "latent" variables, i.e.  $X = X(\theta)$  for some  $\theta$ .

Examples include

- Associations rules,
- Factor Analysis, including principal and independent component Analysis,
- Manifold learning methods including Multidimensional scaling, Locally Linear Embedding, Self organizing maps, ..

Principal curves,...

#### Association Rules

- ▶ Let us use X = (X<sub>1</sub>,...,X<sub>d</sub>) to encode the set of items purchased in a store with X<sub>i</sub> = {0,1} depending on whether item *i* is purchased (1) or not (0).
- The idea underlying association rules is to try to find values of the variables X<sub>1</sub> to X<sub>d</sub> that frequently appear simultaneously.

This approach is found in product organization, cross marketing (i.e. sales promotion) or even finance.

#### Association Rule Analysis

In Association rule analysis, we want to find subsets s<sub>j</sub> of the prototypes that maximize the probability

$$P\left(\cap_{j=1}^{p}X_{j}\in s_{j}
ight)$$

The intersection  $\bigcap_{i=1}^{p} (X_i \in s_i)$  is called conjunctive rule.

- As an example, the association rule could hep us detect two clusters
  - Cluster 1 is defined from the observation that a (sufficiently large) group of customers always purchase items 1 and 3 simultaneously ((1,0,1,0))
  - Cluster 2 is defined from the observation that a (sufficiently large) group of customers always purchase items 2 and 4 simultaneously (X<sub>i</sub> = (0, 1, 0, 1)).

- General approaches at finding subsets of variable values (s<sub>1</sub>,..., s<sub>p</sub>) with relatively high occurence are not feasible for large databases
- Instead, we turn to a more tractable algorithm known as Market Basket Analysis
- ► Market Basket Analysis simplifies the problem by only considering two types of subsets. Either s<sub>j</sub> consists of a single value of X<sub>j</sub>, or it consists of the entire set of values that X<sub>j</sub> can assume (s<sub>j</sub> = S<sub>j</sub>)
- The problem is then reduced to finding subsets of indices J and associated values v<sub>0i</sub> such that the probability

$$P\left(\bigcap_{j\in\mathcal{J}}(X_j=v_{0j})\right)$$



**FIGURE 14.1.** Simplifications for association rules. Here there are two inputs  $X_1$  and  $X_2$ , taking four and six distinct values, respectively. The red squares indicate areas of high density. To simplify the computations, we assume that the derived subset corresponds to either a single value of an input or all values. With this assumption we could find either the middle or right pattern, but not the left one.

#### From HTF, The Elements of Statistical Learning.

- ► Market Basket Analysis also relies on the use of dummy variables Z<sub>k</sub> ∈ {0,1} that each represent one possible value of the variables X<sub>ℓ</sub>.
- ► The first value v<sub>1,1</sub> that X<sub>1</sub> can take is encoded through the variable Z<sub>1</sub> (Z<sub>1</sub> = 1 if this value appears and 0 otherwise), and so on for every value and every variable
- We thus have K = ∑<sub>j=1</sub><sup>p</sup> |S<sub>j</sub>| dummy variables and we can try to determine subsets of indices K for which the probability

$$P\left(igcap_{k\in\mathcal{K}}(Z_k=1)
ight)=P\left(\prod_{k\in\mathcal{K}}Z_k=1
ight)$$

is maximized.

► The set *K* is called an item set. Together those two assumptions define the general Market Basket Analysis formulation.

- Note that two dummy variables Z<sub>k</sub> associated to the same variable X<sub>j</sub> cannot simultaneously take the value 1
- One way to estimate the probability of occurence of given item sets is to count the number of instances in each subset *K*, i.e

$$\hat{P}\left(\prod_{k\in\mathcal{K}}(Z_k=1)
ight)=rac{1}{N}\sum_{i=1}^N\prod_{k\in\mathcal{K}}z_{ik}$$

• Where the variable  $z_{ik}$  is 1 if the corresponding variable  $X_k^i$  in the *i*<sup>th</sup> prototype  $X^i = (X_1^i, \ldots, X_d^i)$  takes the value encoded by this  $z_{ik}$ .

- ► Given a model for the probability P, the clusters are defined by setting a threshold {K<sub>ℓ</sub> | T(K<sub>ℓ</sub>) > t}
- How do we solve the Market Basket Analysis problem in practice? One possibility is the A priori algorithm
  - Start with single item sets and discard the sets for which the support is less than the threshold
  - As a second step, form all subsets of size 2 that can be formed by pairing subsets of size 1. Then discard those size 2 subsets that are less than the threshold.
  - In other words, forming subsets K of size M, only requires looking at the points whose size (M − 1) ancestors all share the same value (i.e are part of a same subset)
  - In terms of computational complexity, the A priori algorithm requires one pass over the data

- Once the A priori algorithm has returned high support item sets *K*, Association Rule Analysis partition each set *K* into disjoint subsets *A* and *B* from which it defines association rules of the form *A* ⇒ *B* where *A* is called antecedent and *B* is called the consequent.
- Association rules are ways of studying how much the purchase of one item influence the purchase of another item.

- To characterize those association rules, we introduce 3 quantities,
  - The support T(A ⇒ B) is the fraction of items in the union of antecedent and the consequent (that is the number of item in the item set K)
  - ► The confidence or predictability C(A ⇒ B) is its support divided by its antecedent.

$$C(A \Rightarrow B) = \frac{T(A \Rightarrow B)}{T(A)}$$

The confidence  $C(A \Rightarrow B)$  can be considered as an estimate of the probability P(B|A)

► The Lift of the rule A ⇒ B is defined as the ratio of the confidence over the expected confidence

$$L(A \Rightarrow B) = \frac{C(A \Rightarrow B)}{T(B)}$$

- The expected confidence T(B) can be computed as T(B) = Pr(∏<sub>k∈A</sub> Z<sub>k</sub> = 1) and is an estimate of the unconditional probability P(B).
- ► The Lift can be considered as an estimate of the dependence/association measure Pr(A ∩ B)/P(A)P(B)
- Once the item sets K with a sufficiently large support have been returned by the A priori algorithm, a confidence threshold c is set and all the association rules with confidence above the threshold are formed, i.e.

$$\{A \Rightarrow B \mid C(A \Rightarrow B) > c\}$$

The output of the Analysis is thus a set of association rules satisfying the constraints

$$T(A \Rightarrow B) > t$$
 and  $C(A \Rightarrow B) > c$ 

- Once the analysis is completed the results are stored in a database which can be accessed to receive specific information on particular items
- Most often we will be interested in retrieving all the transactions (i.e prototypes) in which one particular item appeared as consequent
- The analysis will then indicate the antecedents which might reveal valuable in predicting future sales for the consequent.

#### **Clustering Algorithms**

- There are three main types of clustering algorithms (HTF)
  - Combinatorial algorithms (Work directly on the data without assuming an underlying probability distribution)
  - Mixture Models (Assume that the prototypes are i.i.d. samples from some probability distribution. The probability distribution is assumed to be a mixture of simpler densities, each one of the corresponding to the distribution of a cluster. The distribution is fit to the data through MLE)
  - The last type, called Mode seekers or Bump hunting algorithms try to estimate the modes of the distribution (and hence the clusters) directly from the data (non-parametric)

#### K-means and K-medoid

- The most popular clustering algorithms are combinatorial algorithms which assign every observation to a given cluster without regard to any predefined probabilistic model.
- ▶ The number of clusters *K* is usually predefined
- One approach is to introduce a loss that will drive the assignment. If we let C<sub>k</sub> to denote the k<sup>th</sup> cluster, we get

$$\ell(\mathcal{C}) = rac{1}{2} \sum_{k=1}^{K} \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} d(x_i, x_j)$$

When the dissimilarity is chosen to be the Euclidean distance,

$$d(x_i, x_i') = \sum_{j=1}^{p} (x_{ij} - x_{i'j})^2 = \|x_i - x_{i'}\|^2$$

The loss then reads as

$$\ell(\mathcal{C}) = \frac{1}{2} \sum_{k=1}^{K} \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} \|x_i - x_j\|^2$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

► In particular, developing, we get

$$\begin{split} &\frac{1}{2} \sum_{k} \sum_{i \in \mathcal{C}_{k}} \sum_{j \in \mathcal{C}_{k}} \|x_{i} - x_{j}\|^{2} \\ &= \frac{1}{2} \sum_{k} \sum_{i \in \mathcal{C}_{k}} \sum_{j \in \mathcal{C}_{k}} \langle x_{i}, x_{i} \rangle + \langle x_{j}, x_{j} \rangle - 2 \langle x_{i}, x_{j} \rangle \\ &= \frac{1}{2} \sum_{k} \sum_{i \in \mathcal{C}_{k}} 2 \langle x_{i}, x_{i} \rangle N_{k} - \frac{1}{2} \sum_{k} \sum_{i \in \mathcal{C}_{k}} 2 \langle x_{i}, \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle \\ &= \left( \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \langle x_{i}, x_{i} \rangle - \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \langle x_{i}, \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle \frac{1}{N_{k}} \right) \end{split}$$

 $\frac{1}{2}\sum_{k}\sum_{i\in\mathcal{C}_{k}}\sum_{i\in\mathcal{C}_{k}}\|x_{i}-x_{j}\|^{2}$  $=\left(\sum_{k=1}^{K} N_k \sum_{i \in \mathcal{C}_k} \langle x_i, x_i \rangle - \sum_{k=1}^{K} N_k \sum_{i \in \mathcal{C}_k} \langle x_i, \sum_{j \in \mathcal{C}_k} x_j \rangle \frac{1}{N_k}\right)$  $=\sum_{k=1}^{K}N_{k}\sum_{i\in\mathcal{C}_{k}}\left(\langle x_{i},x_{i}\rangle-\langle x_{i},\sum_{i\in\mathcal{C}_{k}}x_{j}\rangle\frac{1}{N_{k}}\right)$  $=\sum_{k=1}^{K}N_{k}\sum_{i\in\mathcal{C}_{k}}\left(\langle x_{i},x_{i}\rangle-\langle x_{i},\sum_{j\in\mathcal{C}_{k}}x_{j}\rangle\frac{1}{N_{k}}+\frac{1}{N_{k}^{2}}\langle\sum_{j\in\mathcal{C}_{k}}x_{j},\sum_{j\in\mathcal{C}_{k}}x_{j}\rangle\right)$  $-\sum_{k=1}^{K}\sum_{i\in\mathcal{C}_{k}}\left(\frac{1}{N_{k}}\sum_{i\in\mathcal{C}_{k}}\langle x_{j},x_{i}\rangle\right)$ 

$$\begin{split} &\frac{1}{2} \sum_{k} \sum_{i \in \mathcal{C}_{k}} \sum_{j \in \mathcal{C}_{k}} \|x_{i} - x_{j}\|^{2} \\ &= \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \left( \langle x_{i}, x_{i} \rangle - \langle x_{i}, \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle \frac{1}{N_{k}} + \frac{1}{N_{k}^{2}} \langle \sum_{j \in \mathcal{C}_{k}} x_{j}, \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle \right) \\ &- \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \left( \frac{1}{N_{k}} \sum_{j \in \mathcal{C}_{k}} \langle x_{j}, x_{i} \rangle \right) \\ &= \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \left( \langle x_{i}, x_{i} \rangle - 2 \langle x_{i}, \frac{1}{N_{k}} \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle + \langle \frac{1}{N_{k}} \sum_{j \in \mathcal{C}_{k}} x_{j}, \frac{1}{N_{k}} \sum_{j \in \mathcal{C}_{k}} x_{j} \rangle \right) \\ &= \sum_{k=1}^{K} N_{k} \sum_{i \in \mathcal{C}_{k}} \|x_{i} - \frac{1}{N_{k}} \sum_{j \in \mathcal{C}_{k}} x_{j} \|^{2} \end{split}$$

In other words, when using the Euclidean distance, one can write the clustering objective/loss as

$$\ell(\mathcal{C}) = \sum_{k=1}^{K} N_k \sum_{i \in \mathcal{C}_k} \|x_i - \overline{x}^k\|^2$$

Where  $\overline{x}^k$  is the center of mass of the  $k^{th}$  cluster.

the optimal clustering in that framework is thus the clustering that minimizes the average dissimilarity

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Given a set of centers x
k, the optimal clustering is thus defined by solving the minimization

$$C^* = \min_{C} \sum_{k=1}^{K} N_k \sum_{i \in C_k} \|x_i - \overline{x}_k\|^2$$

On the other hand, for a subset of observations x<sub>S</sub>, the center of mass can also read as a minimization

$$x_{S} = \underset{m}{\operatorname{argmin}} \sum_{i \in S} \|x_{i} - m\|^{2}$$

We can thus solve the clustering problem by iterating between the two steps. This gives the Kmeans algorithm.

► 1. Given a cluster assignement C, compute the center of mass of each cluster

$$ar{x}_{\mathcal{S}} = \mathop{\mathrm{argmin}}_{m} \sum_{i \in \mathcal{S}} \|x_i - m\|^2$$

▶ 2. Update the assignement by setting

$$x_i \in \mathcal{C}_k$$
 if  $k = \underset{k}{\operatorname{argmin}} \|x_i - m_k\|^2$ 

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Repeat Steps 1 and 2 until the assignment does not change



source: Bishop, Pattern Recognition and Machine Learning.

### K-means and K-medoid

- The use of the Euclidean distance provides an easy way to compute the center of each cluster.
- It is however possible to get an extension to general more general similarity by using explicit optimization.
- This gives the K-medoid algorithm which iterates over the two steps
  - ► For a given assignement, find the points that are minimizing the distances to any other points in the cluster (the centroid is thus taken as one of the prototypes)

$$c_k \leftarrow \operatorname*{argmin}_{x_i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} d(x_i, x_j)$$

► For a given set of centroids, c<sub>k</sub>, k = 1,..., K, assign each points to its closest centroid,

$$x_{\ell} \in \mathcal{C}_i$$
 with  $i = \operatorname*{argmin}_{1 \leq k \leq K} d(x_i, c_k)$ 

## K-means and K-medoid

- Because it relies on the Euclidean distance, K means is often much more sensitive to outliers (i.e it places the highest influence on the largest distances)
- ► Conversely, the number of operations needed to derive the centroid in K-means was only O(N<sub>k</sub>). For K-medoid, it is now O(N<sup>2</sup><sub>k</sub>) (i.e. comparing every pair of points in the clusters).
- K-means and K-medoid both suffer from several drawbacks among which
  - Both methods assume the number of clusters to be known beforehand (not true in practice)
  - As iterative techniques, they are both sensitive to initial conditions
  - Finally, both K-means and K-medoid converge to local minimas.

## K-means and K-medoid: Initialization

- Random partitioning. The approach divides the dataset in K distinct clusters chosen at random.
- ► The Forgy or *Lloyd-Forgy* approach (FA). The method picks *K* feature vectors at random to define the centroids and assigns remaining feature vectors to their nearest centroid.
- MacQueeen's approach. Instead of considering a batch algorithm where the whole dataset if used to update the centroids, the MacQueeen approach is the extension of the Forgy approach to a recalculation of the centroids after each assignment.
- Kaufman's approach. The Kauffman initialization is a deterministic initialization which places the centroids in the areas where there is a higher density of prototypes.

#### MacQueen's approach

- Initialize the centroids, m<sub>k</sub> and set the cluster sizes to zero, n<sub>k</sub> = 0, for all k = 1,..., K.
- Repeat

Pick an observation x<sub>i</sub> and determine the cluster following

$$k = \underset{x}{\operatorname{argmin}}_{\ell} \|\boldsymbol{x}_{i} - \boldsymbol{m}_{k}\|^{2}$$

Update the centroid *m<sub>k</sub>* following

$$\boldsymbol{m}_k \leftarrow \boldsymbol{m}_k + \frac{1}{n_k} (\boldsymbol{x}_i - \boldsymbol{m}_k)$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

where  $n_k$  is the number of prototypes in cluster k.

## Kauffman's approach

- Set x
  <sub>0</sub> to be the "median" feature vector (i.e the one that is the most centrally located)
- Set S = {x
  <sub>0</sub>} and let S<sup>c</sup> denote the set of remaining feature vectors.
- ▶ For all the remaining feature vectors  $x_i \in S^c$ ,

• compute 
$$\overline{d}_i \equiv \min_{k \in S} d(\overline{x}_k, x_i)$$

- Set  $C_{i,j} \equiv \max(\overline{d}_i d(\mathbf{x}_j, \mathbf{x}_i), 0)$
- ► Choose the next centroid x̄<sub>k+1</sub> from the index j that solves max<sub>i</sub> C<sub>ij</sub> (The result will be a prototype with high density of points around him (i.e the prototype with the largest number of points closer to him than to the previous centroids))
- If there are K points in S stop. Otherwise set S ← S ∪ x̄<sub>k+1</sub> and go back to step 1.

## Hierarchical clustering (I)

- As we saw, an inconvenient aspect of K-means or K-medoid is that those methods require the user to explicitly provide a number of clusters.
- Hierarchical clustering does not require such initialization. Instead it only relies on a measure of dissimilarity between groups of observations
- The algorithm then iteratively defines the clusters by either merging (bottom up) or dividing (top down) the set of measurements
  - ► Agglomerative strategies start at the bottom (*n* clusters) and successively merge a selected pair of clusters (*n* − 1 clusters)
  - Divisive approaches start at the top (1 cluster) and successively split the previous clusters into two new clusters.
- It is then up to the user to decide which level represents the most natural clustering.

## Hierarchical clustering (II)

- Agglomerative and divisive Hierarchical Clustering approaches can be defined to possess a monotonicity property which means that the dissimilarity between the clusters that are merged increase monotonically with the level in the hierarchy.
- Recursive splitting/agglomeration happening in hierarchical clustering can be represented by a tree (a.k.a dendogram) in which the nodes represent the clusters. The root node encode the whole dataset and the leafs represent each of the prototypes.
- Furthermore, when the clustering approach satisfies the monotonicity property, the height in the tree represents the dissimilarity between merged clusters.

## Hierarchical clustering (III)



FIGURE 14.12. Dendrogram from agglomerative hierarchical clustering with average linkage to the human tumor microarray data.

≣⇒

## Hierarchical clustering (IV)



#### from H., T., F., The elements of Statistical Learning

## Agglomerative Clustering

- Agglomerative clustering starts with every prototype representing a singleton.
- At each step, it then chooses the closest two clusters and merge them into a single cluster.
- ► For any two clusters *H* and *G*, there are three common ways to define the dissimilarity between those clusters

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

## Agglomerative Clustering

In Single linkage (SL), the dissimilarity between two clusters H and G is defined from the closest pair of points

$$d_{SL}(G,H) = \min_{i \in G, j \in H} d(i,j)$$

 In Complete Linkage (CL), the inter-cluster dissimilarity is defined from the dissimilarity of the furthest pair,

$$d_{CL}(G,H) = \max_{i \in G, j \in H} d(i,j)$$

 Finally, in Group Average clustering (GA), the criterion is the average dissimilarity

$$d_{GA}(G,H) = \frac{1}{N_G N_H} \sum_{i \in G} \sum_{j \in H} d(i,j)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

## A couple of observations (I)

- Single linkage only requires the dissimilarity to be small for one pair of prototypes from each clusters
- It therefore leads to clusters that violates the compactness assumption (i.e clusters are compact if all the observations within them are relatively close to each other)
- If we define the diameter of a group of observations as

$$D_G = \max_{i,j\in G} d(i,j),$$

single linkage can produce large diameters clusters

## A couple of observations (II)

- On the opposite, complete linkage will generate compact clusters with small diameters
- However CL can violate the closeness property. Two points from one cluster can be closer to prototypes from another cluster than they are to the prototypes in their own cluster.
- ► E.g. Assume C<sub>1</sub> and C<sub>2</sub> where C<sub>2</sub> has a single prototype that is very far from the points in C<sub>1</sub> and all the others which are very close (i.e closer than the closest distance between prototypes in C<sub>1</sub>)



FIGURE 14.13. Dendrograms from agglomerative hierarchical clustering of human tumor microarray data.

from HTF, The elements of Statistical Learning

#### How about divisive clustering?

- One possibility is to use a K-means algorithm (with K=2) to split one cluster at each step. But this approach would depend on the initialization of K – means and does not satisfy the monotonicity property
- ► One alternative is the Macnaughton-Smith algorithm. Starting from a single cluster C<sub>0</sub> with |C<sub>0</sub>| = n, remove that prototype x<sub>1</sub> which has the largest average dissimilarity,

$$oldsymbol{y}^* = rgmax_{oldsymbol{y}} \overline{D_0}(oldsymbol{y}) = rgmax_{oldsymbol{y}} rac{1}{n-1} \sum_{oldsymbol{x}_i \in \mathcal{C}_0} d(oldsymbol{x}_i,oldsymbol{y})$$

Then define  $C_1 = \{y^*\}$ ,  $C_0 = C_0 \setminus \{y^*\}$ . For the second step, remove the prototype  $y \in C_0$  that maximizes the difference

$$D_{C_0} - D_{C_1} = rac{1}{|C_0| - 1} \sum_{\mathbf{x}_j \in C_0} d(\mathbf{y}, \mathbf{x}_j) - rac{1}{|C_1|} \sum_{\mathbf{x}_j \in C_1} d(\mathbf{x}_j, \mathbf{y})$$

#### How about divisive clustering?

- Other alternatives include Kaufman and Rousseeuw: choose at each step the cluster that maximizes the diameter and split this particular cluster
- Largest average dissimilarity: split the cluster that maximizes

$$\bar{d}_G = rac{1}{N_G^2} \sum_{oldsymbol{x}_i \in G} \sum_{oldsymbol{x}_j \in G} d(oldsymbol{x}_i, oldsymbol{x}_j)$$

In any of these approaches, recursive splitting is then applied until each cluster is singleton or members of the cluster have zero dissimilarity.

## Spectral clustering I

- ► Given N prototypes x<sub>i</sub>, i = 1,..., N, spectral clustering starts with the N × N matrix encoding the pairwise similarities between points, s(x<sub>i</sub>, x<sub>j</sub>)
- Spectral clustering might be particularly interesting on non convex data where traditional clustering methods such as K-means might underperform
- The data is then represented by a graph G = (V, E) where E denotes the set of edges and V denotes the set of vertices.

• Each edge  $w_{ij}$  in the graph is given the weight  $s(x_i, x_j)$ .

## Spectral clustering II

- The whole idea of spectral clustering is to partition the graph into clusters, such that edges between clusters have low weights and edges inside each cluster have higher weights
- Given a set of prototypes x<sub>i</sub>, let d(x<sub>i</sub>, x<sub>j</sub>) to denote the Euclidean distance between those prototypes.
- We can define the similarity from the Radial kernel s(x, x<sub>j</sub>) = exp(-d(x<sub>i</sub>, x<sub>j</sub>)γ) where γ > 0 is a scale parameter.
- Given the similarity matrix, one can either keep all interactions into account, or only retain those interactions corresponding to the K nearest neighbors

## Spectral clustering III

- ► The matrix, A<sub>ij</sub> of edge weights, A<sub>ij</sub> = w<sub>i,j</sub> is called the adjacency matrix of the graph. We call degree of a vertex *i*, the sum d<sub>i</sub> = ∑<sub>j</sub> A<sub>i,j</sub>.
- ▶ If we define the diagonal matrix **D** encoding the degrees as

$$oldsymbol{D} = \left[egin{array}{ccc} d_1 & & & \ & \ddots & & \ & & d_N \end{array}
ight]$$

Then the (unnormalized) Laplacian of the graph is defined as the matrix  $\boldsymbol{L} = \boldsymbol{D} - \boldsymbol{A}$ 

• One can equivalently work with the normalized graph Laplacian  $\tilde{L} = I - D^{-1}A$ 

## Spectral clustering IV

- Spectral clustering then works by first computing the eigendecomposition of the graph Laplacian
- It then selects the *m* smallest eigenvalues and their corresponding eigenvectors
- This gives a matrix V of size N × m on which one can apply traditional clustering algorithms such as K-means to find a clustering of the original set of prototypes
- ► Finding the smallest eigenpair implies that the representation given by V will preserve the similarity encoded in the w<sub>ij</sub>.

$$\mathbf{v}^{T} \mathbf{L} \mathbf{v} = \sum_{i=1}^{N} d_{i} v_{i}^{2} - \sum_{i=1}^{N} \sum_{j=1}^{N} v_{i} v_{j} w_{ij}$$
$$= \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} w_{ij} (v_{i} - v_{j})^{2}$$