

Unsupervised Learning

Yuyan Zhou

May 20, 2020

Types of learning

- Supervised learning

labelled data. $x \xrightarrow{f} y$

- Classification
- Regression

- Semi-supervised learning

- Active learning

- Unsupervised learning

- Clustering
- Dimension reduction

- Reinforcement learning

- ...

Cluster analysis

- Top-down
- Bottom-up
- Key questions:
 - How to measure proximity
 - How to choose the number of clusters
 - Initialization

k too large n
too small /

Proximity matrices

- N objects, N by N matrix D to describe their pairwise proximity
- Nonnegative? $D_{ij} \geq 0$?
- Symmetric? $D_{ij} = D_{ji}$
- Can be suitably convert to a dissimilarity matrix?
- Triangle inequality? $D_{ij} + D_{jk} > D_{ik}$?

Common dissimilarities

- N measurements x_{ij} for $i=1,2,\dots,N$ and $j=1,2,\dots,p$ (variables/attributes)
- Define $D(x_i, x_{i'}) = \sum_{j=1}^p d_j(x_{ij}, x_{i'j})$
- Most common $d_j(x_{ij}, x_{i'j}) = (x_{ij} - x_{i'j})^2$ minimize $\sum (x_{ij} - x_{i'j})^2 \propto 2(1 - \rho(x_i, x_{i'}))$ ↑ maximize.
- Equivalent to based on correlation (similarity)

$$\rho(x_i, x_{i'}) = \frac{\sum_j (x_{ij} - \bar{x}_i)(x_{i'j} - \bar{x}_{i'})}{\sqrt{\sum_j (x_{ij} - \bar{x}_i)^2 \sum_j (x_{i'j} - \bar{x}_{i'})^2}} \quad \bar{x}_i = \frac{1}{p} \sum_j x_{ij}$$

Weighted dissimilarity

- Convex combination

$$D(x_i, x_{i'}) = \sum_{j=1}^p \underbrace{w_j}_{\text{weight}} \cdot d_j(x_{ij}, x_{i'j}); \quad \sum_{j=1}^p w_j = 1.$$

$$\bar{D} = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N D(x_i, x_{i'}) = \sum_{j=1}^p \underbrace{w_j}_{\text{weighted influence of attribute } j} \cdot \bar{d}_j,$$

$$\underbrace{\bar{d}_j}_{\text{average dissimilarity for attribute } j} = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N d_j(x_{ij}, x_{i'j})$$

Weighted squared Euclidean distance

$$D_I(x_i, x_{i'}) = \sum_{j=1}^p \underbrace{w_j}_{\text{d}} \cdot \underbrace{(x_{ij} - x_{i'j})^2}$$

$$\underbrace{\bar{d}_j}_{\text{influence of } j \text{ over the data set.}} = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N (x_{ij} - x_{i'j})^2 = \underbrace{2 \cdot \text{var}_j}$$

1° $w_j = 1$

2° $w_j = \frac{1}{\bar{d}_j}$ ✓

Clustering algorithms

- Combinatorial algorithms

- Work directly on the observed data
- No direct inference to an underlying probability model
- Encoder, many-to-one mapping
- Local optimal, small-sized

- Mixture modeling

- MLE, Bayesian approach

- Mode seeking

- Bump hunting
- PRIM

ESL. 9.3

$$C(N) = \vec{v} \quad v \in \{1, \dots, K\}.$$
$$N \longrightarrow K \quad S_2(N, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^N.$$
$$S_2(19, 4) \approx 10^{10}.$$

Common Heuristic in Practice:

The Lloyd's method

k-means.

[Least squares quantization in PCM, Lloyd, IEEE Transactions on Information Theory, 1982]

Input: A set of n datapoints $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n$ in \mathbb{R}^d

$$d(\mathbf{x}_i, \mathbf{x}_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

Initialize centers $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ and
clusters C_1, C_2, \dots, C_k in any way.

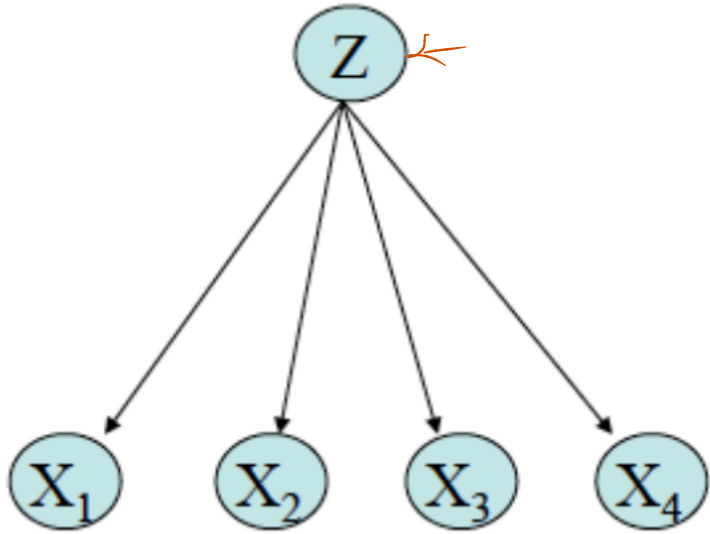
Repeat until there is no further change in the cost.

- For each j : $C_j \leftarrow \{x \in S \text{ whose closest center is } \mathbf{c}_j\}$
- For each j : $\mathbf{c}_j \leftarrow \text{mean of } C_j$

Holding $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$ fixed,
pick optimal C_1, C_2, \dots, C_k

Holding C_1, C_2, \dots, C_k fixed,
pick optimal $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$

Gaussian mixtures as soft K-means



K clusters $P(K=k | \bar{x})$

$$P(x | K=k) = \prod_i N(x_i | \mu_k, \sigma^2)$$

$$P(x) = \sum_{k=1}^K P(K=k | \bar{x}) \cdot \prod_i N(x_i | \mu_k, \sigma^2)$$

E-step: $P(z | x, \theta)$

M-step: maximize $E_{z|x, \theta} \log P(x, z | \theta')$

$\sigma^2 \rightarrow 0$ soft K-means \rightarrow K-means.

Hierarchical clustering

- Bottom-up
- Top-down
- Different linkage
 - Single linkage
 - Complete linkage
 - Average linkage
 - Ward's linkage

Spectral clustering

- Non-convex clusters
- Undirected similarity graph $G = \langle V, E \rangle$
- Graph-partition problem
 - Edges between different groups have low weight
 - Within a group have high weight
- Random walk on the graph with transition probability matrix P
 - Construct groups that the random walk seldom transitions from one group to another

PCA

$$X \in \mathbb{R}^{D \times n} \quad x_1, x_2, \dots, x_n \in \mathbb{R}^D$$

$$\underbrace{V \in \mathbb{R}^{D \times d}}_{\Delta} \quad \underbrace{d \ll D.}_{\sim}$$

- Principal components

- Projections of data ✓
- Mutually uncorrelated (orthogonal) $V_i \cdot V_j = 0 \quad i \neq j. \quad V_i \cdot V_j = 1 \quad i=j \Rightarrow \|V_i\| = 1$
- Ordered in variance

x centered. $x_i - \bar{x}$

$$\frac{1}{n} \sum_{i=1}^n (V^T x_i)^2 = \underline{V^T X X^T V.}$$

maximize $V^T X X^T V$

s.t. $\underline{V^T V = I.}$

$$\partial (V^T X X^T V - \lambda V^T V) / \partial V = 0 \Rightarrow \underline{(X X^T) V = V \Lambda}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d & & \\ & & & 0 & \dots & 0 \end{bmatrix}$$

Kernel principal components

$$\boxed{\underline{\phi}}: \mathbb{R}^D \rightarrow \mathbb{R}^P \quad \underline{F}: \mathbb{R}^P \quad D < P \quad \left| \quad \underline{k(x,y)} = \underline{\phi(x)^T \phi(y)} \quad \text{inner product} \quad \textcircled{1}: \underline{\phi(X) \phi(X)^T \phi(X) \alpha} = \lambda \phi(X) \alpha \quad \checkmark$$

$$X = [x_1, \dots, x_n] \quad \checkmark \quad \textcircled{2}: \underline{\phi(X)^T \phi(X)}$$

$$\phi(X) = [\phi(x_1), \phi(x_2), \dots, \phi(x_n)] \in \mathbb{R}^{P \times N} \quad \checkmark$$

kernel

$$\textcircled{2} \quad \phi(X)^T \phi(X) \phi(X)^T \phi(X) \alpha = \lambda \phi(X)^T \phi(X) \alpha$$

$$\underline{K = \phi(X)^T \phi(X)} \in \mathbb{R}^{N \times N} \quad K_{ij} = \phi(x_i)^T \phi(x_j)$$

$$\textcircled{2} \quad K \cdot K \alpha = \lambda K \alpha \Rightarrow \underline{K \alpha = \lambda \alpha} \quad \checkmark$$

$$\text{cov matrix: } C_F = \underline{\frac{1}{N} \phi(X) \phi(X)^T} = \underline{\frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^T}$$

$$\underline{C_F} \cdot \underline{P} = \lambda \underline{P} \quad \textcircled{1}$$

$$\lambda \neq 0 \Rightarrow \underline{P} = \underline{\frac{1}{\lambda} \sum_{i=1}^N \phi(x_i) \phi(x_i)^T P} = \underline{\sum_{i=1}^N \alpha_i \phi(x_i)} = \underline{\phi(X) \cdot \alpha} \quad P = \phi(X) \alpha$$

$$\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix}$$

ICA

- PCA finds directions of maximum variation ✓
- ICA would find directions most “aligned” with data
 - Statistically independent linear transformation
 - Optimize based on mutual information