

Sapienza University of Rome

Master in Artificial Intelligence and Robotics
Master in Engineering in Computer Science

Machine Learning

A.Y. 2020/2021

Prof. L. Iocchi, F. Patrizi

14. Unsupervised Learning

L. Iocchi, F. Patrizi

with contributions from Valsamis Ntouskos

Overview

- Learning without a teacher
- K-means algorithm
- Gaussian Mixture Model
- Expectation Maximization algorithm
- General EM problem

Reference

C. Bishop. Pattern Recognition and Machine Learning. Chapter 9.

Unsupervised Learning

Input data available $D = \{\mathbf{x}_n\}$, but target values not available.

Unsupervised data clustering: finding multiple classes from data.

Modelling input data useful when combined with supervised learning.

Gaussian Mixture Model

Gaussian Mixture Model (GMM)

Mixed probability distribution P formed by k different Gaussian distributions

$$P(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

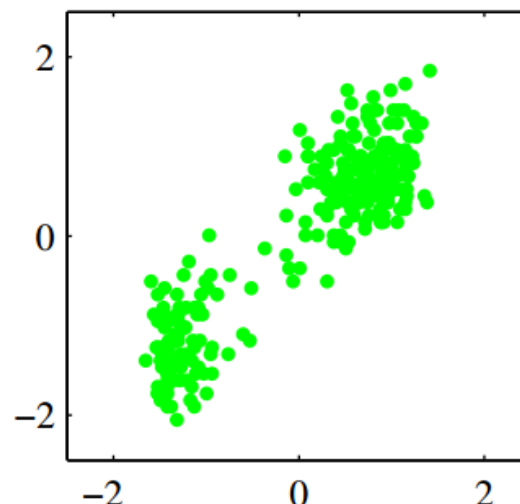
- π_k , prior probability
- $\boldsymbol{\mu}_k$, mean
- $\boldsymbol{\Sigma}_k$, covariance matrix

Unsupervised learning algorithms determine mixed probability distributions from data.

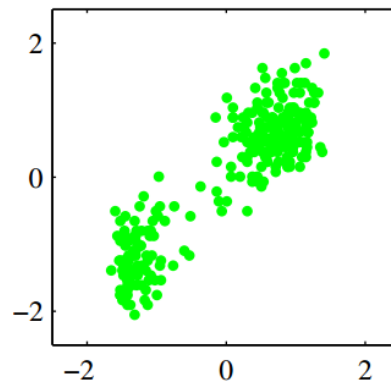
Generating data from mixture of Gaussians

Each instance \mathbf{x}_n generated by

- ① Choosing Gaussian k according to prior probabilities $[\pi_1, \dots, \pi_K]$
- ② Generating an instance at random according to that Gaussian, thus using $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$



K-means



Computing K means of data generated from K Gaussian distributions.

Input: $D = \{\mathbf{x}_n\}$, value K Output: μ_1, \dots, μ_K

K-means

Step 1. Begin with a decision on the value of $k = \text{number of clusters}$

Step 2. Put any initial partition that classifies the data into k clusters. You may assign the training samples randomly, or systematically as follows

- ① Take the first k training samples as single-element clusters
- ② Assign each of the remaining $(N-k)$ training samples to the cluster with the nearest centroid. After each assignment, recompute the centroid of the new cluster.

K-means

Step 3. Take each sample in sequence and compute its distance from the centroid of each of the clusters. If a sample is not currently in the cluster with the closest centroid, switch this sample to that cluster and update the centroid of the two clusters involved in the switch.

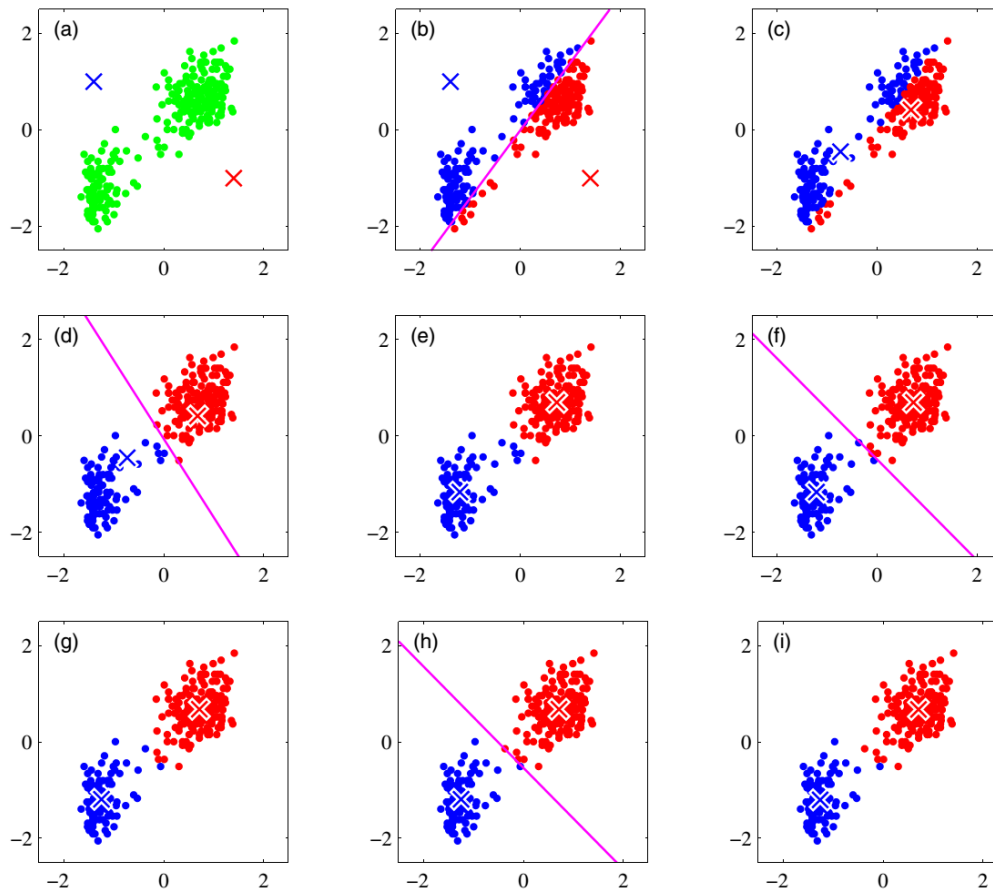
Step 4. Repeat step 3 until convergence is achieved, that is until a pass through the training sample causes no new assignments.

K-means convergence

The convergence will always occur if the following conditions are satisfied:

1. For each switch in step 2, the sum of distances from each training sample to that training sample's group centroid is decreased.
2. There are only finitely many partitions of the training examples into k clusters.

K-means example



L. Iocchi, F. Patrizi

14. Unsupervised Learning

11 / 38

Sapienza University of Rome, Italy - Machine Learning (2020/2021)

Remarks on K-means

- The number of clusters K must be determined before hand.
- Sensitive to initial condition (local optimum) when a few data available.
- Not robust to outliers. Very far data from the centroid may pull the centroid away from the real one.
- The result is a circular cluster shape because it is based on distance.

Remarks on K-means

Some solutions:

- use K-means clustering only if there are many data available
- use median instead of mean
- define better *distance* functions

Gaussian Mixture Model

$$P(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Introduce new variables $z_k \in \{0, 1\}$, with $\mathbf{z} = (z_1, \dots, z_K)^T$ using a 1-out-of- K encoding (only one component is 1, all the others are 0).

Let's define

$$P(z_k = 1) = \pi_k$$

thus

$$P(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

Gaussian Mixture Model

For a given value of \mathbf{z} :

$$P(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Thus

$$P(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

Joint distribution: $P(\mathbf{x}, \mathbf{z}) = P(\mathbf{x} | \mathbf{z})P(\mathbf{z})$ (chain rule).

Gaussian Mixture Model

When \mathbf{z} are variables with 1-out-of- K encoding and $P(z_k = 1) = \pi_k$

$$P(\mathbf{x}) = \sum_{\mathbf{z}} P(\mathbf{z})P(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

GMM distribution $P(\mathbf{x})$ can be seen as the marginalization of a distribution $P(\mathbf{x}, \mathbf{z})$ over variables \mathbf{z} .

Gaussian Mixture Model

Given observations $D = \{(\mathbf{x}_n)_{n=1}^N\}$, each data point \mathbf{x}_n is associated to the corresponding variable \mathbf{z}_n which is unknown.

Note: $z_{nk} = 1$ denotes \mathbf{x}_n sampled from Gaussian k

\mathbf{z}_n are called **latent variables**.

Analysis of latent variables allows for a better understanding of input data (e.g., dimensionality reduction).

Gaussian Mixture Model

Let's define the posterior

$$\gamma(z_k) \equiv P(z_k = 1 | \mathbf{x}) = \frac{P(z_k = 1) P(\mathbf{x} | z_k = 1)}{P(\mathbf{x})}$$

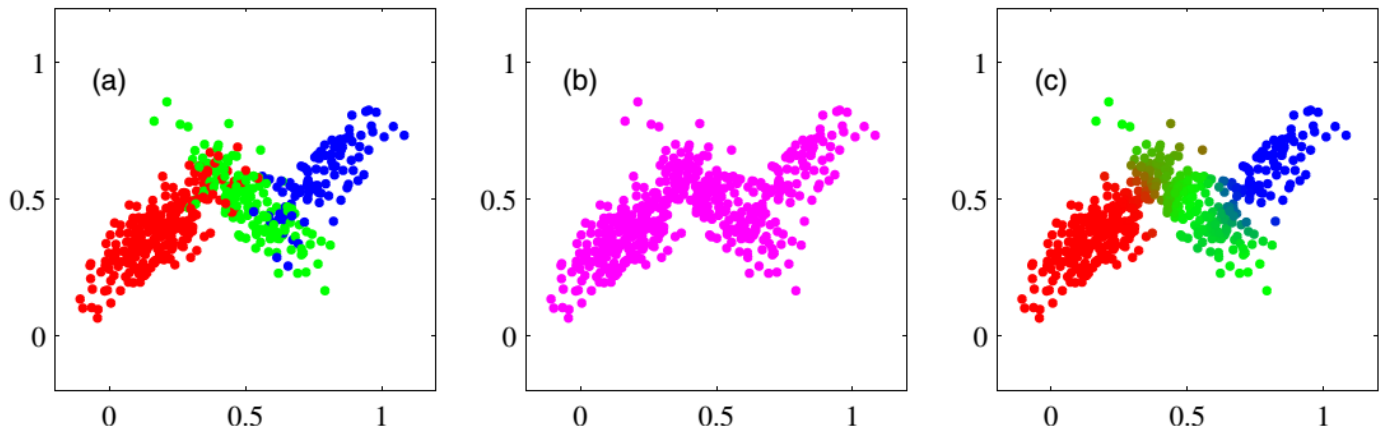
$$\gamma(z_k) = \frac{\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Note:

π_k : prior probability of z_k

$\gamma(z_k)$: posterior probability after observation of \mathbf{x} .

Gaussian Mixture Model example



- a) $P(\mathbf{x}, \mathbf{z})$ with 3 latent variables \mathbf{z} (red, green, blue)
 b) $P(\mathbf{x})$ marginalized distribution
 c) $\gamma(z_{n,k})$ posterior distribution

Expectation Maximization (EM)

Given data set $D = \{(\mathbf{x}_n)_{n=1}^N\}$ and GMM

$$P(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

determine $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$

Note: generalization of K-means algorithm

Expectation Maximization (EM)

Maximum likelihood

$$\operatorname{argmax}_{\pi, \mu, \Sigma} \ln P(\mathbf{X} | \pi, \mu, \Sigma)$$

At maximum:

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}, \quad \text{with } N_k = \sum_{n=1}^N \gamma(z_{nk})$$

Expectation Maximization (EM)

- **E step**

Given π_k, μ_k, Σ_k , compute $\gamma(z_{nk})$

- **M step**

Given $\gamma(z_{nk})$, compute π_k, μ_k, Σ_k

Expectation Maximization (EM)

- Initialize $\pi_k^{(0)}, \mu_k^{(0)}, \Sigma_k^{(0)}$
- Repeat until termination condition $t = 0, \dots, T$
 - **E step**

$$\gamma(z_{nk})^{(t+1)} = \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j^{(t)} \mathcal{N}(\mathbf{x}_n; \mu_j^{(t)}, \Sigma_j^{(t)})}$$

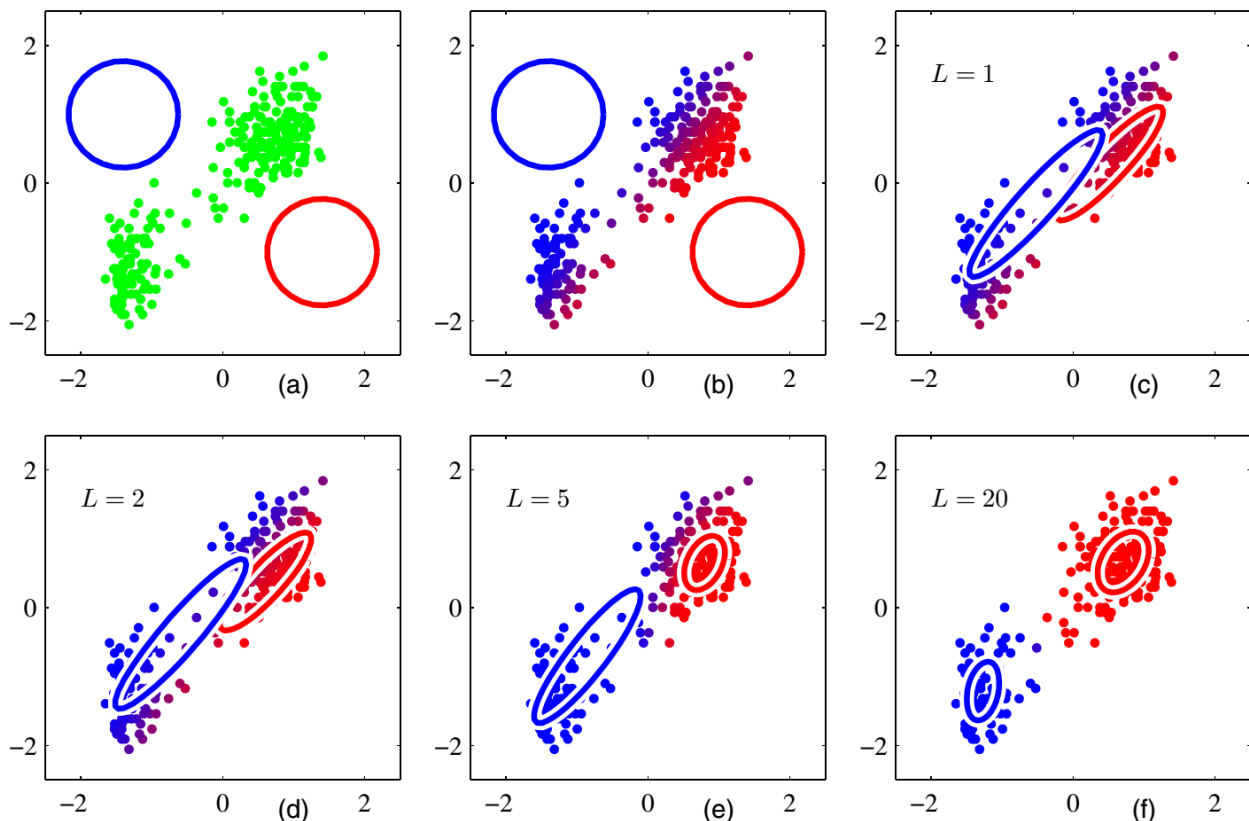
- **M step**

$$\mu_k^{(t+1)} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk})^{(t+1)} \mathbf{x}_n$$

$$\Sigma_k^{(t+1)} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk})^{(t+1)} (\mathbf{x}_n - \mu_k^{(t+1)}) (\mathbf{x}_n - \mu_k^{(t+1)})^T$$

$$\pi_k^{(t+1)} = \frac{N_k}{N}, \quad \text{with } N_k = \sum_{n=1}^N \gamma(z_{nk})^{(t+1)}$$

EM example



Remarks on EM Algorithm

- Converges to local maximum likelihood
- Provides estimates of the latent variables z_{nk}
- Extended version of K-means (probabilistic assignment to a cluster z_{nk})
- Can be generalized to other distributions (not only Gaussians)

General EM Problem

Given:

- Observed data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Unobserved latent values $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$
- Parametrized probability distribution $P(\mathbf{Y}|\boldsymbol{\theta})$, where
 - $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ is the full data $\mathbf{y}_n = \langle \mathbf{x}_n, \mathbf{z}_n \rangle$
 - $\boldsymbol{\theta}$ are the parameters

Determine:

- $\boldsymbol{\theta}^*$ that (locally) maximizes $E[\ln P(\mathbf{Y}|\boldsymbol{\theta})]$

Many uses:

- Unsupervised clustering
- Bayesian Networks
- Hidden Markov Models

General EM Method

Define likelihood function $Q(\theta'|\theta)$ defined on variables $\mathbf{Y} = \mathbf{X} \cup \mathbf{Z}$, using observed \mathbf{X} and current parameters θ to estimate \mathbf{Z}

EM Algorithm:

Estimation (E) step: Calculate $Q(\theta'|\theta)$ using current hypothesis θ and observed data \mathbf{X} to estimate probability distribution over \mathbf{Y}

$$Q(\theta'|\theta) \leftarrow E[\ln P(\mathbf{Y}|\theta')|\theta, \mathbf{X}]$$

Maximization (M) step: Replace hypothesis θ by the hypothesis θ' that maximizes this Q function

$$\theta \leftarrow \operatorname{argmax}_{\theta'} Q(\theta'|\theta)$$

Bayesian networks



Syntax:

- a set of nodes, one per variable
- a directed, acyclic graph (link \approx “directly influences”)
- a conditional distribution for each node given its parents:
 $P(X_i | \text{Parents}(X_i))$

In the simplest case, conditional distribution represented as a *conditional probability table* (CPT) giving the distribution over X_i for each combination of parent values.

Burglar BN Example

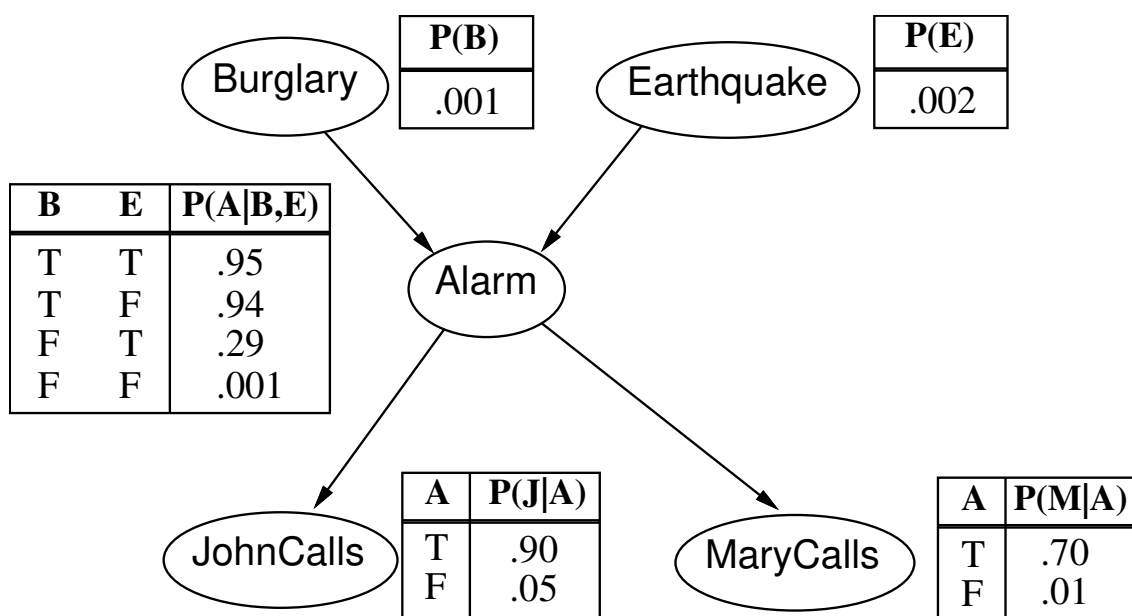
I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes the alarm is set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

Network topology reflects "causal" knowledge:

- A burglar can set the alarm
- An earthquake can set the alarm
- The alarm can cause Mary to call
- The alarm can cause John to call

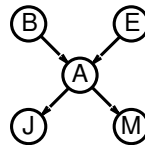
Burglar BN Example



Computing joint probabilities

All joint probabilities computed with the chain rule:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i))$$

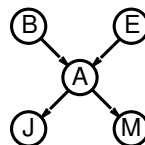


e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned}
 &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\
 &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\
 &\approx 0.00063
 \end{aligned}$$

Compactness

A CPT for Boolean variable X_i with k Boolean parents has 2^k rows for the combinations of parent values



Each row requires one number p for $X_i = \text{true}$ (the number for $X_i = \text{false}$ is just $1 - p$)

If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers

I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution

For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)

Learning BN with observable variables

Bayesian Networks model dependencies among random variables.

When structure **known** and all variables **observable**

conditional probabilities can be estimated with maximum likelihood.

Example

$$P(j|a) \approx \frac{|\{\langle a, j, \dots \rangle\}|}{|\{\langle a, \dots \rangle\}|}$$

Learning BN with hidden variables

Unsupervised learning can be seen as learning with a BN with one hidden variable (the *class* of the instances).

This can be generalized to general BN with multiple hidden variables.

Learning BN with hidden variables: example

Consider three random variables $X \in \{0, 1\}$, $A \in \{a_1, a_2\}$, $B \in \{b_1, b_2\}$, with X unobservable.

How can we learn BN parameters for $P(X)$, $P(A|X)$, and $P(B|X)$ from instances $D = \{d_1, \dots, d_n\}$, with $d_k = \langle a_k, b_k \rangle$?

Define:

$$P(X = 0) = \theta_0, P(A = a_1|X = 0) = \theta_1, P(A = a_1|X = 1) = \theta_2, \\ P(B = b_1|X = 0) = \theta_3, P(B = b_1|X = 1) = \theta_4$$

$$\theta = \langle \theta_0, \theta_1, \theta_2, \theta_3, \theta_4 \rangle$$

Apply EM method to find maximum likelihood wrt θ from D .

Learning BN with hidden variables: example

Estimation of BN parameters:

$$P(X = x_j) = \frac{1}{n} E[\hat{N}(X = x_j)] \\ P(A = a_i|X = x_j) = \frac{E[\hat{N}(A = a_i, X = x_j)]}{E[\hat{N}(X = x_j)]}$$

Note that

$$E[\hat{N}(\cdot)] = E\left[\sum_k I(\cdot|d_k)\right] = \sum_k P(\cdot|d_k)$$

Learning BN with hidden variables: example

Estimation of BN parameters:

$$\begin{aligned}
 P(X = x_j) &= \frac{1}{n} \sum_{k=1}^n P(X = x_j | d_k) \\
 P(A = a_i | X = x_j) &= \frac{\sum_{k=1}^n P(A = a_i, X = x_j | d_k)}{\sum_{k=1}^n P(X = x_j | d_k)} \\
 P(B = b_l | X = x_j) &= \dots
 \end{aligned}$$

Apply Bayes rule

$$\begin{aligned}
 P(x_j | d_k) &= P(x_j | \langle a_k, b_k \rangle) = \frac{P(a_k | x_j) P(b_k | x_j)}{\sum_i P(a_k | x_i) P(b_k | x_i) P(x_i)} = \phi_1(\theta) \\
 P(a_i, x_j | d_k) &= P(a_i | x_j, d_k) P(x_j | d_k) = \phi_2(\theta) \\
 P(b_l, x_j | d_k) &= P(b_l | x_j, d_k) P(x_j | d_k) = \phi_3(\theta)
 \end{aligned}$$

... to define $Q(\theta' | \theta)$

Summary

- Unsupervised learning useful to deal with unknown variables
- Clustering when labeled data are not available
- EM algorithm is a general method to estimate likelihood for mixed distributions including observed and latent variables
- Concepts to be extended to continuous latent variables