



Learning with hidden variables

Generative and Deep Learning (GDL)

Daide Bacciu (davide.bacciu@unipi.it)



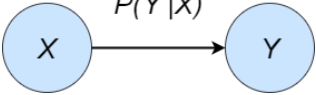
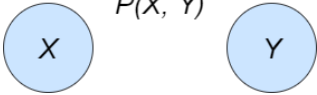
UNIVERSITÀ DI PISA



Lecture Outline

- ◇ Learning with non observed variables
- ◇ Latent/hidden variable models
- ◇ Maximum likelihood learning with latent variable
- ◇ Expectation-Maximization algorithm
- ◇ Exact maximum likelihood learning in mixture models

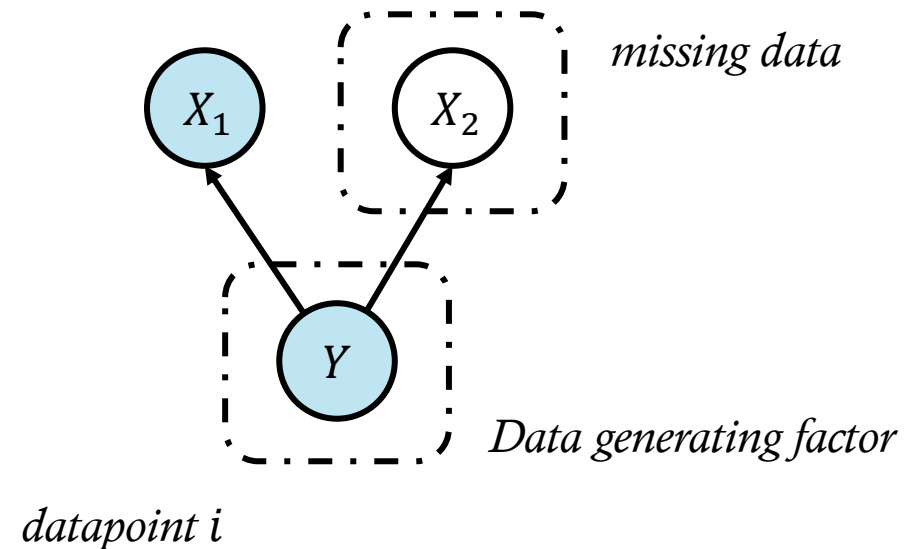
Learning with graphical models

		Structure	
		Fixed Structure 	Fixed Variables 
Data	Complete	Naive Bayes Calculate Frequencies (ML)	Discover dependencies from the data Structure Search Independence tests
	Incomplete	Latent variables EM Algorithm (ML) MCMC, VBEM (Bayesian)	Difficult Problem Structural EM

Latent/Hidden Variable Models

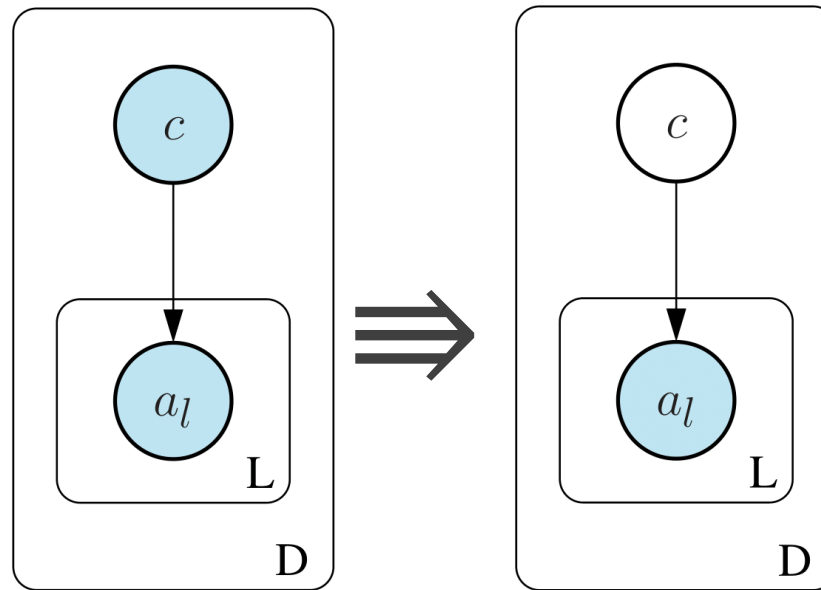
Dealing with Incomplete Data

- ◇ Incomplete data
 - ◇ Missing observations (**imputation**)
 - ◇ Unobserved random variables
- ◇ Key idea
 - ◇ Complete the data **making hypothesis on the unobserved variables**
 - ◇ Review the hypothesis
 - ◇ Iterate



What if...

Sample classification is **no longer observable**



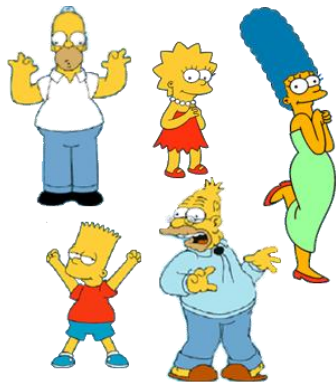
c is now a **Hidden/Latent Variable** within a **Mixture Model**

Mixture Model

No classification \Rightarrow Seek a natural grouping



What is a natural grouping?



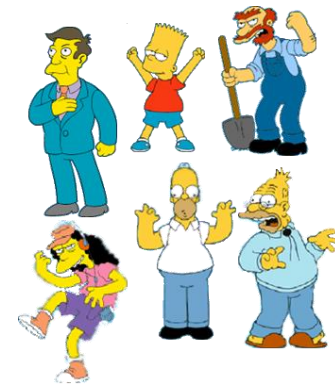
Family



School



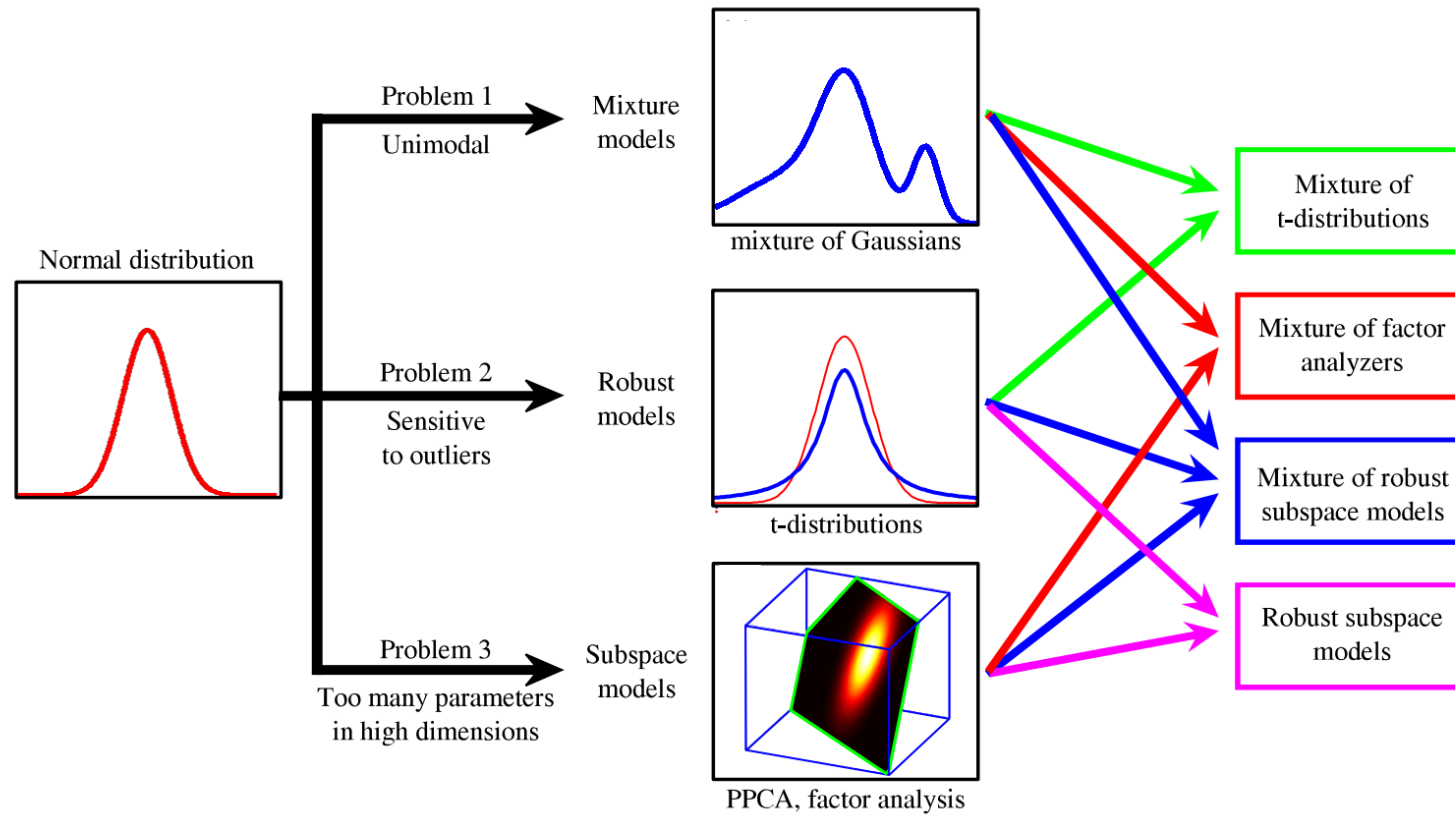
Females



Males

Credit goes to Eamonn Keogh @ UCR

Mixture Models

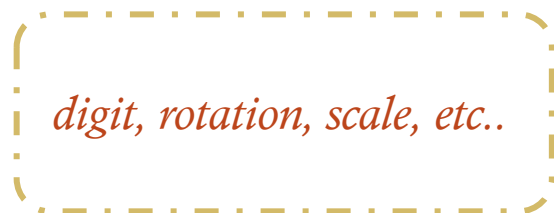


Latent Variable Models

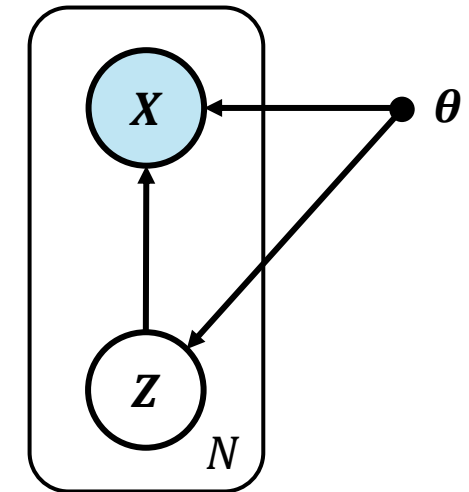
We want to model a probability distribution $P(\mathbf{X} = \mathbf{x}|\theta)$, parametrized by θ , with \mathbf{X} being a multivariate RV.



high-dimensional data \mathbf{x}



latent factors \mathbf{z} (low-dimensional)



N datapoints

These models contains both

- ◇ Observed random variables \mathbf{X} (i.e. for which we have training data)
- ◇ Unobserved ([hidden/latent](#)) variables \mathbf{Z} (e.g. data clusters)

Joint distribution with latents

We want to model a probability distribution $P(\mathbf{X} = \mathbf{x}|\theta)$, with \mathbf{X} being a multivariate RV

θ is a generic term to denote params, but these distributions have **different** params

Joint distribution **defined** as

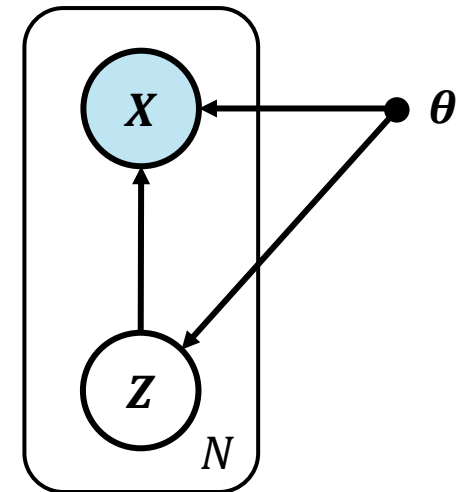
$$P(\mathbf{x}, \mathbf{z}|\theta) = P(\mathbf{x}|\mathbf{z}, \theta)P(\mathbf{z}|\theta)$$

Marginal likelihood (marginalize out the multivariate RV \mathbf{Z})

continuous $P(\mathbf{x}|\theta) = \int P(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int P(\mathbf{x}|\mathbf{z}, \theta)P(\mathbf{z}|\theta) d\mathbf{z}$

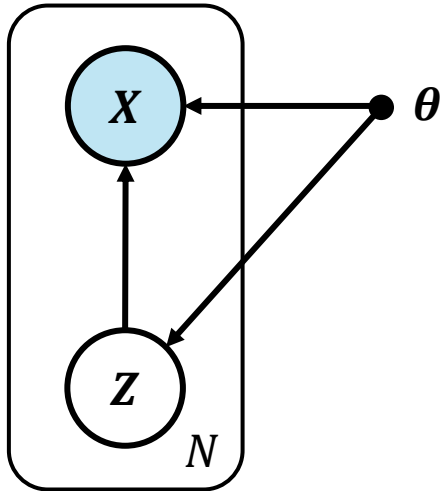
discrete $P(\mathbf{x}|\theta) = \sum_{\mathbf{z}} P(\mathbf{x}, \mathbf{z}|\theta) = \sum_{\mathbf{z}} P(\mathbf{x}|\mathbf{z}, \theta)P(\mathbf{z}|\theta)$

in both cases $P(\mathbf{x}|\theta) = \mathbb{E}_{P(\mathbf{z}|\theta)}[P(\mathbf{x}|\mathbf{z}, \theta)]$



Learning with latent variables

Learning in latent variable models



Learning: Obtain the value (ML/MAP) or a probability distribution (Bayesian) for the parameters θ of all the distributions involved
Here: $P(\mathbf{X}|\mathbf{Z}, \theta)$ and $P(\mathbf{Z}|\theta)$

We focus on **Maximum Likelihood estimation**

$$\theta_{ML} = \arg \max_{\theta} P(\mathbf{X}|\theta)$$

but we need to circumvent the issue that we don't know how to compute $P(\mathbf{X}|\theta)$ because it does not contain the \mathbf{Z} needed by the model

Completing the likelihood

- Given the likelihood maximization

$$\theta_{ML} = \arg \max_{\theta} \mathcal{L}(\theta|\mathbf{X}) = \arg \max_{\theta} P(\mathbf{X}|\theta)$$

- $\mathcal{L}(\theta|\mathbf{X})$ is the **incomplete likelihood** and \mathbf{X} is the **incomplete data** containing only observed random variables
- Assume **complete data** $\mathbf{d} = (\mathbf{X}, \mathbf{Z})$ exists where $z \in \mathbf{Z}$ are **hidden or latent** variables
- Obtain the **complete likelihood**

$$\mathcal{L}_c(\theta|\mathbf{d}) = P(\mathbf{d}|\theta) = P(\mathbf{X}, \mathbf{Z}|\theta) = P(\mathbf{Z}|\mathbf{X}, \theta)P(\mathbf{X}|\theta)$$

Key Idea

Maximize the complete likelihood \mathcal{L}_c instead of \mathcal{L}

Expectation Maximization

Algorithm solving an optimization problem involving maximization of **complete log-likelihood** $\mathcal{L}_c(\theta)$ w.r.t. model parameters θ

$$\theta^{(k+1)} = \arg \max_{\theta} \underbrace{\sum_{\mathbf{z}} P(\mathbf{Z} = \mathbf{z} | \mathbf{X}, \theta^{(k)}) \log P(\mathbf{X}, \mathbf{Z} = \mathbf{z} | \theta)}_{\mathbb{E}_{P(\mathbf{Z} | \mathbf{X}, \theta^{(k)})}[\cdot]}$$

Requires that posterior computation is feasible

A 2-step iterative algorithm

E-Step: Given the current estimate of the model parameters $\theta^{(k)}$, compute

$$Q^{(k+1)}(\theta | \theta^{(k)}) = \mathbb{E}_{\mathbf{z} \sim P(\mathbf{Z} | \mathbf{X}, \theta^{(k)})} [\log P(\mathbf{X}, \mathbf{Z} = \mathbf{z} | \theta^{(k)})]$$

M-Step: Find the new estimate of the model parameters

$$\theta^{(k+1)} = \arg \max_{\theta} Q^{(k+1)}(\theta | \theta^{(k)})$$

Iterate until

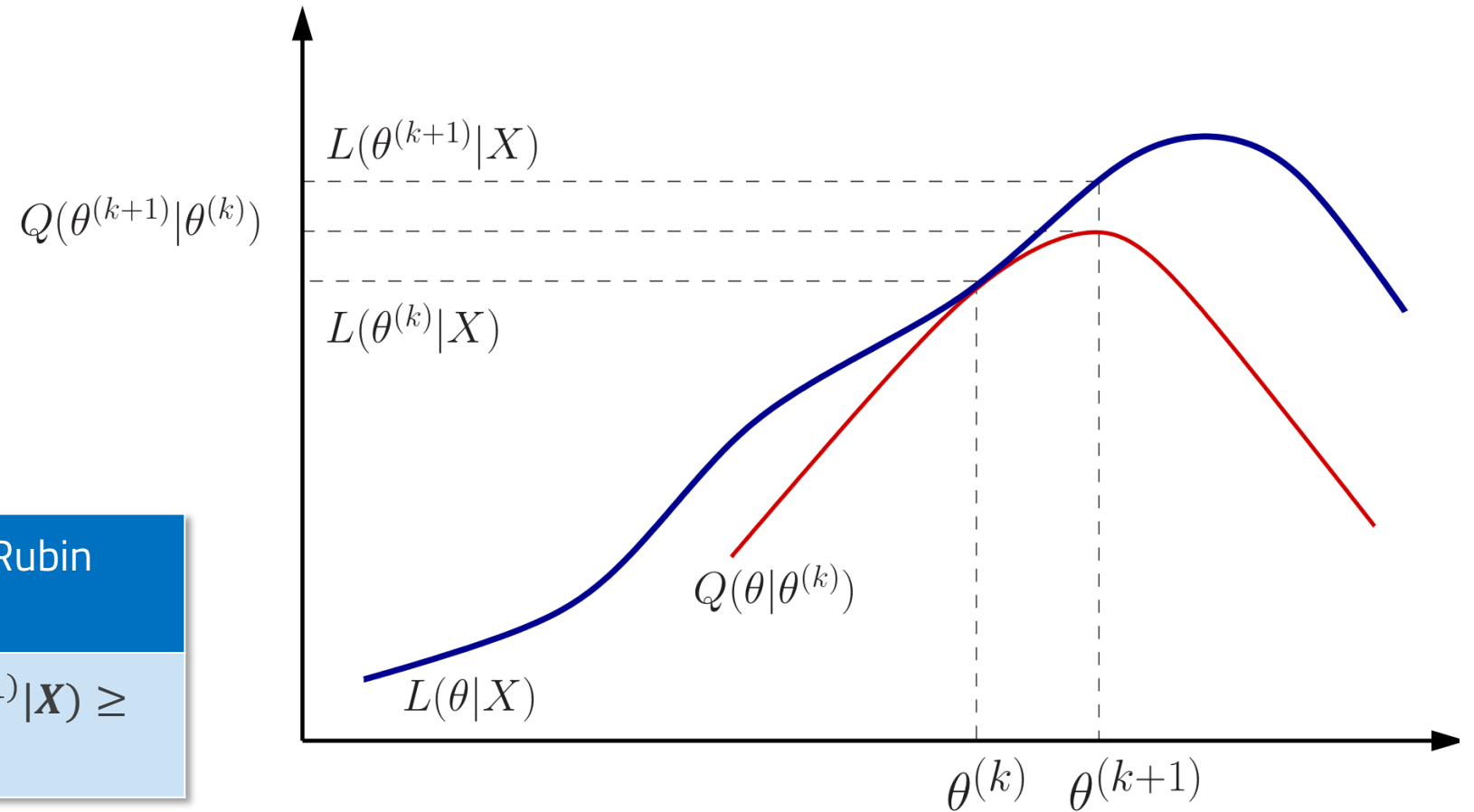
$$|\mathcal{L}_c(\theta)^{k+1} - \mathcal{L}_c(\theta)^k| < \epsilon$$

(or stop if maximum number of iterations is reached)

EM Graphically

Theorem (Dempster, Laird and Rubin (1977))

An EM process ensures $\mathcal{L}(\theta^{(k+1)} | \mathbf{X}) \geq \mathcal{L}(\theta^{(k)} | \mathbf{X})$



What if posterior is intractable?

Sampling methods

- ❖ Monte Carlo
 - ❖ Basic integration
 - ❖ Importance Sampling
 - ❖ Rejection Sampling
 - ❖ ...
- ❖ Markov Chain Monte Carlo (MCMC)
 - ❖ Metropolis Algorithm
 - ❖ Gibbs Sampling
 - ❖ Hamiltonian Monte Carlo (HMC)

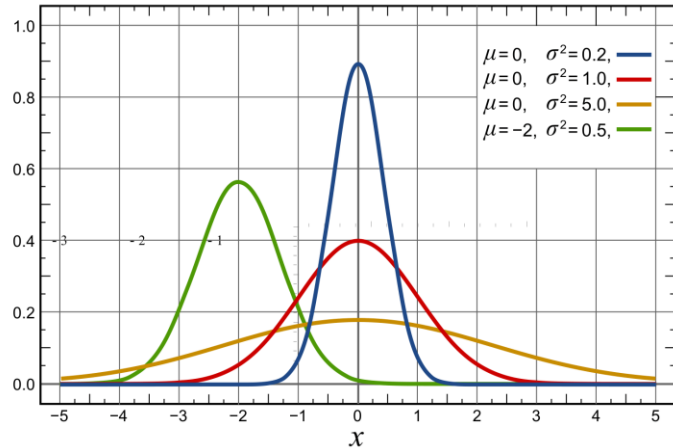
Variational methods

- ❖ Mean Field Variational Bayes
- ❖ Variational Expectation Maximization (VEM)
- ❖ Stochastic Variational Inference (SVI)
- ❖ Amortized Variational Inference
 - ❖ Variational Autoencoders

We shall see (some of) these in the next lectures!

Gaussian Mixture Models

Univariate Gaussian



Random Variable

Unidimensional $x \in \mathbb{R}$

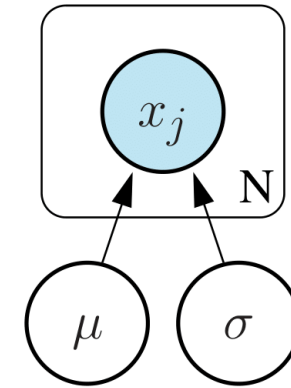
Parameterized Family

$$p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp - \left(\frac{(x-\mu)^2}{2\sigma^2} \right)$$

Model Parameters $\theta = (\mu, \sigma)$

Graphical model

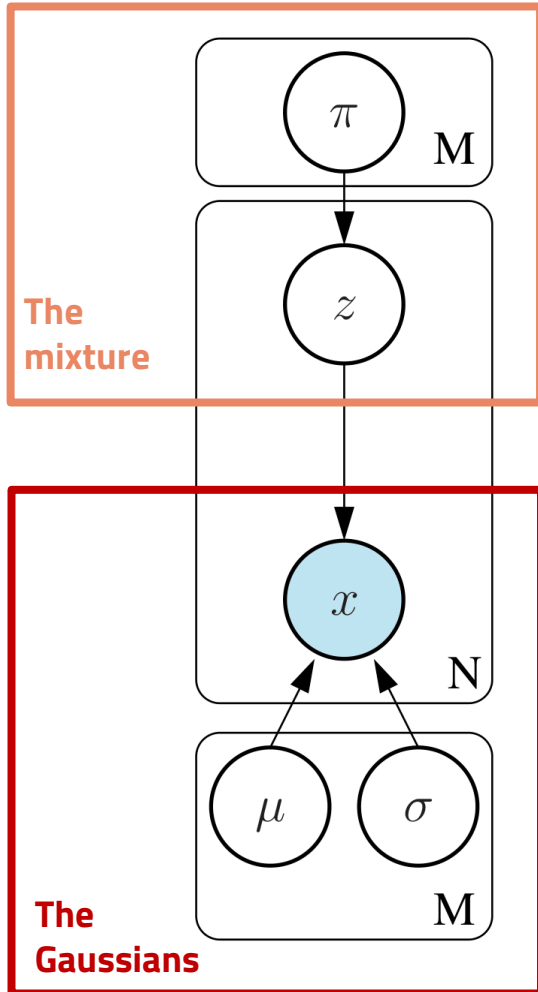
Fully observable with single RV
 $x_j \in \mathbb{R}$ and two parameters μ
and σ



Can be learned by maximum likelihood optimization

$$\begin{aligned} \langle \mu_{ML}, \sigma_{ML} \rangle &= \arg \max_{\mu, \sigma} \log P(\mathbf{d}|\mu, \sigma) = \arg \max_{\mu, \sigma} \sum_j^N \log P(x_j|\mu, \sigma) \\ &= \arg \max_{\mu, \sigma} \sum_j^N \left(\log \sqrt{2\pi}\sigma - \left(\frac{(x_j - \mu)^2}{2\sigma^2} \right) \right) \end{aligned}$$

Gaussian Mixture Model (GMM)



◇ A model for **continuous observable data** $\mathbf{x} \in \mathbb{R}^d$ involving **M Gaussians** that are mixed by **latent/hidden variables** $z \in \{1, \dots, M\}$

◇ Each observation \mathbf{x} is generated by a Gaussian with probability

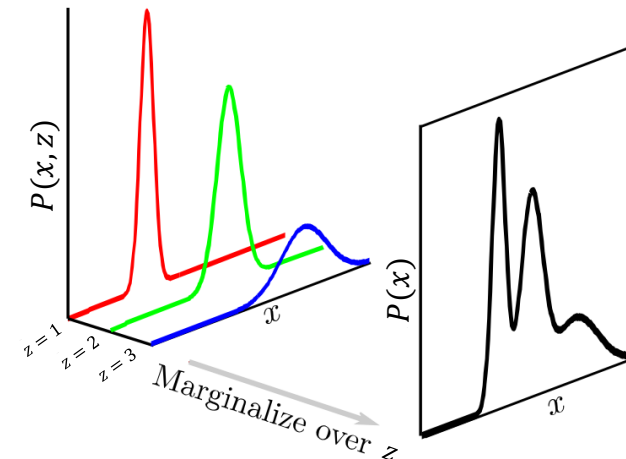
$$P(\mathbf{x}|z = m) \sim \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\sigma}_m)$$

◇ z is the **hidden mixture selection variable** with prior

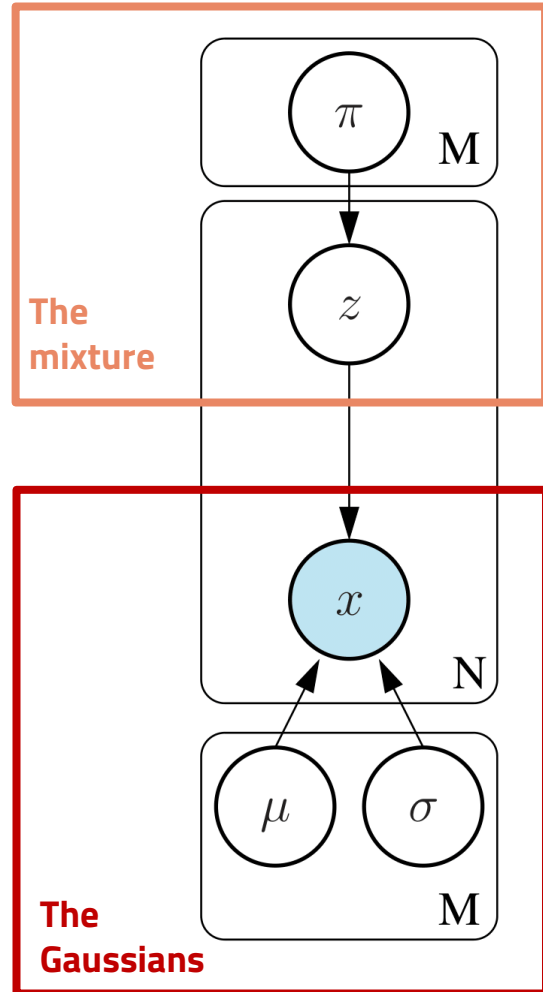
$$P(z = m) = \pi_m$$

◇ Model parameters $\theta = (\pi, \mu, \sigma)$ **estimated by the EM algorithm**

The mixture density is created by marginalizing the hidden z in the joint $P(\mathbf{x}, z)$



Gaussian Mixture Model (GMM)



◇ A model for **continuous observable data** $\mathbf{x} \in \mathbb{R}^d$ involving **M Gaussians** that are mixed by **latent/hidden variables** $z \in \{1, \dots, M\}$

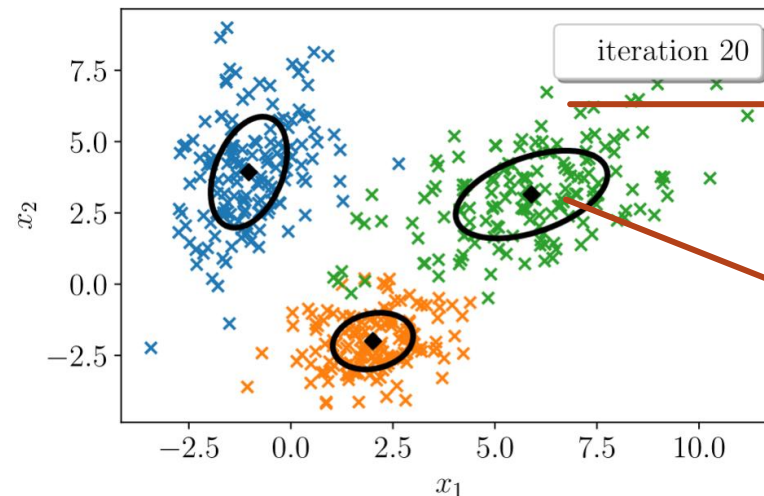
◇ Each observation \mathbf{x} is generated by a Gaussian with probability

$$P(\mathbf{x}|z = m) \sim \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\sigma}_m)$$

◇ z is the **hidden mixture selection variable** with prior

$$P(z = m) = \pi_m$$

◇ Model parameters $\theta = (\pi, \mu, \sigma)$ **estimated by the EM algorithm**



Can be used to cluster data into groups

The "group" corresponds to the value of the latent variable

Inference (feasible!)

$$P_{\theta}(z = m|\mathbf{x}) = \frac{P_{\theta}(\mathbf{x}|z = m)P_{\theta}(z = m)}{P_{\theta}(\mathbf{x})}$$

GMM Likelihood

A set of i.i.d observations $\mathbf{X} = \{x_1, \dots, x_n\}$ has incomplete likelihood

$$\mathcal{L}(\theta|\mathbf{X}) = P(\mathbf{X}|\theta) = \prod_{j=1}^N P(x_j|\theta)$$

Use **marginalization** to introduce the **hidden** variable $z_j = m$ (for all m)

$$\mathcal{L}(\theta|\mathbf{X}) = \prod_{j=1}^N \sum_{m=1}^M P(x_j, z_j = m|\theta)$$

Following the **independence relationships** in the graphical model this rewrites

$$\mathcal{L}(\theta|\mathbf{X}) = \prod_{j=1}^N \sum_{m=1}^M P(z_j = m|\pi) P(x_j|z_j = m, \mu, \sigma)$$

Categorical

univariate Gaussian
from slide 18

A technicality towards the complete likelihood

- ◆ We define the **complete** log-likelihood using **auxiliary indicator variables**

$$\bar{z}_{jm}(z_j) = \begin{cases} 1 & \text{if } z_j = m \\ 0 & \text{otherwise} \end{cases}$$

- ◆ In practice: it is like **assuming we know the hidden assignment** of each data point x_j to the mixture m for all j, m

Rationale: contribution of $P(x_j | z_j = m, \theta) \cdot P(z_j = m | \theta)$ factor to the likelihood should be 0 **if we know** that sample j belongs to mixture $m' \neq m$

Indicator variables are binary RVs for which, by definition, we have $P(\bar{z}_{jm} = 1) = P(z_j = m | x_j)$

GMM Complete Likelihood

Assume we know the hidden assignment $z_{jm} = 1$ for each sample j , we have the following complete likelihood

$$\mathcal{L}_c(\theta|\mathbf{X}, \mathbf{Z}) = \prod_{j=1}^N \sum_{m=1}^M z_{jm} \pi_m P(x_j | \mu_m, \sigma_m)$$

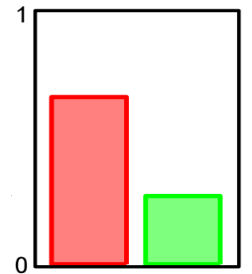
GMM Complete Likelihood

Assume we know the hidden assignment $z_{jm} = 1$ for each sample j , we have the following complete likelihood

$$\mathcal{L}_c(\theta|\mathbf{X}, \mathbf{Z}) = \prod_{j=1}^N \sum_{m=1}^M z_{jm} \pi_m P(x_j | \mu_m, \sigma_m) = \prod_{j=1}^N \prod_{m=1}^M \left(\pi_m P(x_j | \mu_m, \sigma_m) \right)^{z_{jm}}$$

Can be better handled by taking the log

$$\log \mathcal{L}_c(\theta|\mathbf{X}, \mathbf{Z}) = \sum_{j=1}^N \sum_{m=1}^M z_{jm} \log \left(\pi_m P(x_j | \mu_m, \sigma_m) \right)$$



Now we can compute

$$\begin{aligned} Q(\theta|\theta^{(k)}) &= E_{P(\mathbf{Z}|\mathbf{X}, \theta^{(k)})} [\log \mathcal{L}_c(\theta|\mathbf{X}, \mathbf{Z})] \\ &= \sum_{j=1}^N \sum_{m=1}^M P(z_j = m | x_j, \theta^{(k)}) \log \left(\pi_m^{(k)} P(x_j | \mu_m^{(k)}, \sigma_m^{(k)}) \right) \end{aligned}$$

$\xrightarrow{\hspace{10em}} P(\bar{z}_{jm} = 1) = P(z_j = m | x_j)$

Expectation Step

To compute the auxiliary function

$$Q(\theta|\theta^{(k)}) = \sum_{j=1}^N \sum_{m=1}^M P(z_j = m | x_j, \theta^{(k)}) \log \left(\pi_m^{(k)} P \left(x_j | \mu_m^{(k)}, \sigma_m^{(k)} \right) \right)$$

we need to estimate the **posterior** $P(z_j = m | x_j, \theta^{(k)})$ based on the **current** $\theta^{(k)}$ values

$$P(z_j = m | x_j, \theta^{(k)}) = \frac{\pi_m^{(k)} P \left(x_j | \mu_m^{(k)}, \sigma_m^{(k)} \right)}{\sum_{m'} \pi_{m'}^{(k)} P \left(x_j | \mu_{m'}^{(k)}, \sigma_{m'}^{(k)} \right)}$$

Maximization Step

Maximization of the auxiliary function with respect to the parameters $\theta = (\pi, \mu, \sigma)$

$$Q(\theta|\theta^{(k)}) = \sum_{j=1}^N \sum_{m=1}^M P(z_j = m|x_j, \theta^{(k)}) \log \pi_m \\ + \sum_{j=1}^N \sum_{m=1}^M P(z_j = m|x_j, \theta^{(k)}) \log P(x_j|\mu_m, \sigma_m)$$

As usual this is performed by solving

$$\frac{\partial Q(\theta|\theta^{(k)})}{\partial \theta} = 0$$

while taking into account the **sum-to-one constraint** $\sum_{m=1}^M \pi_m = 1$

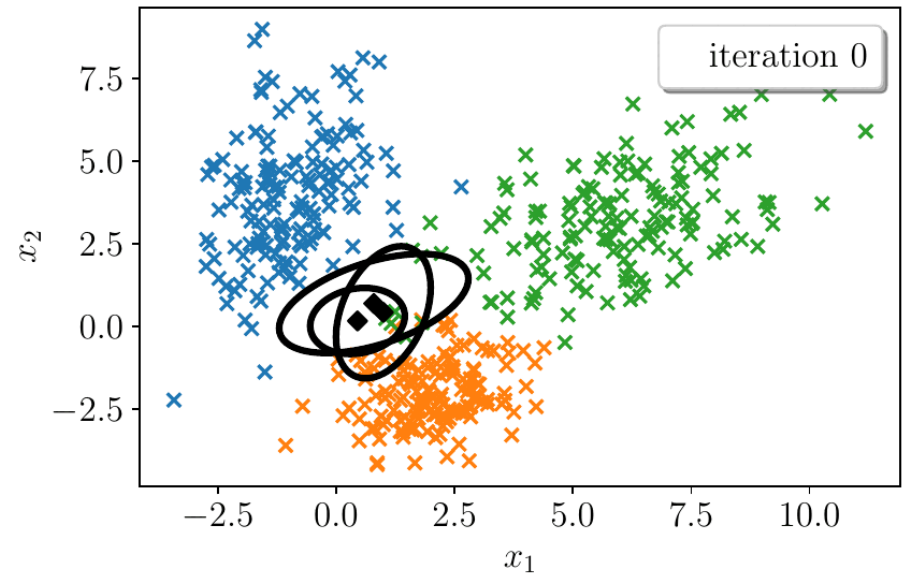
M-Step Learning Equations

Solving the **independent maximization problems** with respect to π_m , μ_m and σ_m yields

$$\pi_m^{(k+1)} = \frac{\sum_{j=1}^N P(z_j = m | x_j, \theta^{(k)})}{N}$$

$$\mu_m^{(k+1)} = \frac{\sum_{j=1}^N x_j P(z_j = m | x_j, \theta^{(k)})}{\sum_{j=1}^N P(z_j = m | x_j, \theta^{(k)})}$$

$$\sigma_m^{(k+1)} = \sqrt{\frac{\sum_{j=1}^N P(z_j = m | x_j, \theta^{(k)}) (x_j - \mu_m^{(k)})^2}{\sum_{j=1}^N P(z_j = m | x_j, \theta^{(k)})}}$$



Wrap-Up

The EM Algorithm in a Nut-Shell

The EM algorithm can be **summarized** by the optimization problem

$$\theta^{(k+1)} = \arg \max_{\theta} \sum_{\mathbf{z}} P(\mathbf{Z} = \mathbf{z} | \mathbf{X}, \theta^{(k)}) \log \mathcal{L}_c(\theta | \mathbf{X}, \mathbf{Z} = \mathbf{z})$$

The **posterior** $P(\mathbf{Z} = \mathbf{z} | \mathbf{X}, \theta^{(k)})$ provides an **expected count** for the **event** \mathbf{z} that is the EM counterpart of basic ML counting

Take Home Messages

- ◆ Latent variable models allow generation of visible information (data) to be steered by hidden knowledge about non-observable variables
 - ◆ Assuming conditional **independence relationships are known**
- ◆ Expectation-Maximization
 - ◆ Use it for **dealing with ML learning with unobserved variables**
 - ◆ **Complete observations** with hidden variables
 - ◆ Use **posterior** to estimate unobserved counts (at least for the countable case)
 - ◆ Compute observed terms using such **pseudo-counts**
- ◆ All of the above hinges on the fact that **we can compute the exact posterior**
 - ◆ At some point we will have to deal with how to approximate the posterior

Next Lecture(s)

- ◇ A probabilistic model for sequences: [Hidden Markov Models](#) (HMMs)
- ◇ Exact inference on a [chain](#) with [observed and unobserved](#) variables
- ◇ The Expectation-Maximization algorithm for HMMs
- ◇ Graphical models with [varying structure](#): Dynamic Bayesian Networks