

# Mixture models

## IAML: Mixture models and EM

Victor Lavrenko and Charles Sutton  
School of Informatics

Semester 1

- Recall types of clustering methods
  - hard clustering: clusters do not overlap
    - element either belongs to cluster or it does not
  - soft clustering: clusters may overlap
    - strength of association between clusters and instances
- Mixture models
  - probabilistically-grounded way of doing soft clustering
  - each source: a generative model (Gaussian or multinomial)
  - parameters (e.g. mean/covariance are unknown)
- Expectation Maximization (EM) algorithm
  - automatically discover all parameters for the K “sources”

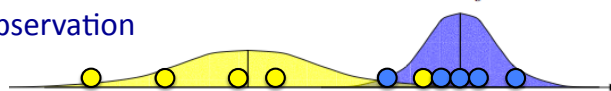
Copyright © 2011 Victor Lavrenko

## Mixture models in 1-d

- Observations  $x_1 \dots x_n$ 
  - K=2 Gaussians with unknown  $\mu, \sigma^2$
  - estimation trivial if we know the source of each observation

$$\mu_b = \frac{x_1 + x_2 + \dots + x_{n_b}}{n_b}$$

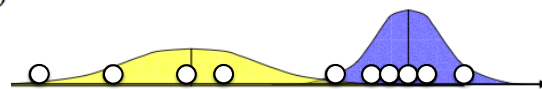
$$\sigma_b^2 = \frac{(x_1 - \mu_b)^2 + \dots + (x_{n_b} - \mu_b)^2}{n_b}$$



- What if we don't know the source?
- If we knew parameters of the Gaussians ( $\mu, \sigma^2$ )
  - can guess whether point is more likely to be a or b

$$P(b | x_i) = \frac{P(x_i | b)P(b)}{P(x_i | b)P(b) + P(x_i | a)P(a)}$$

$$P(x_i | b) = \frac{1}{\sqrt{2\pi\sigma_b^2}} \exp\left(-\frac{(x_i - \mu_b)^2}{2\sigma_b^2}\right)$$



Copyright © 2011 Victor Lavrenko

## Expectation Maximization (EM)

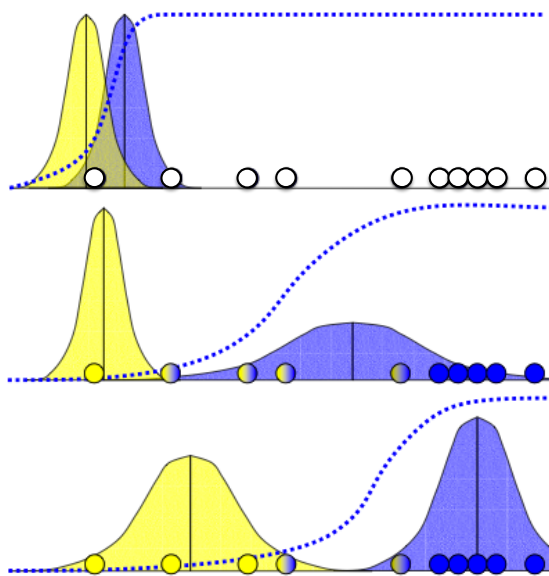
- Chicken and egg problem
  - need  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$  to guess source of points
  - need to know source to estimate  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$
- EM algorithm
  - start with two randomly placed Gaussians  $(\mu_a, \sigma_a^2), (\mu_b, \sigma_b^2)$

E-step: – for each point:  $P(b | x_i)$  = does it look like it came from b?

M-step: – adjust  $(\mu_a, \sigma_a^2)$  and  $(\mu_b, \sigma_b^2)$  to fit points assigned to them  
– iterate until convergence

Copyright © 2011 Victor Lavrenko

## EM: 1-d example



Copyright © 2011 Victor Lavrenko

$$P(x_i | b) = \frac{1}{\sqrt{2\pi\sigma_b^2}} \exp\left(-\frac{(x_i - \mu_b)^2}{2\sigma_b^2}\right)$$

$$b_i = P(b | x_i) = \frac{P(x_i | b)P(b)}{P(x_i | b)P(b) + P(x_i | a)P(a)}$$

$$a_i = P(a | x_i) = 1 - b_i$$

$$\mu_b = \frac{b_1 x_1 + b_2 x_2 + \dots + b_n x_n}{b_1 + b_2 + \dots + b_n}$$

$$\sigma_b^2 = \frac{b_1 (x_1 - \mu_b)^2 + \dots + b_n (x_n - \mu_b)^2}{b_1 + b_2 + \dots + b_n}$$

$$\mu_a = \frac{a_1 x_1 + a_2 x_2 + \dots + a_n x_n}{a_1 + a_2 + \dots + a_n}$$

$$\sigma_a^2 = \frac{a_1 (x_1 - \mu_a)^2 + \dots + a_n (x_n - \mu_a)^2}{a_1 + a_2 + \dots + a_n}$$

could also estimate priors:

$$P(b) = (b_1 + b_2 + \dots + b_n) / n$$

$$P(a) = 1 - P(b)$$

## Gaussian mixture models: d>1

- Data with d attributes, from k sources

- Each source c is a Gaussian

- Iteratively estimate parameters:

- prior: what % of instances came from source c?

$$P(c) = \frac{1}{n} \sum_{i=1}^n P(c | \vec{x}_i)$$

- mean: expected value of attribute j from source c

$$\mu_{c,j} = \sum_{i=1}^n \left( \frac{P(c | \vec{x}_i)}{nP(c)} \right) x_{i,j}$$

- covariance: how correlated are attributes j and k in source c?

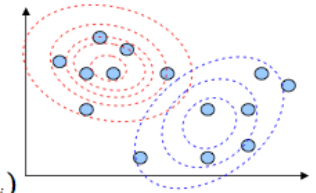
$$(\Sigma_c)_{j,k} = \sum_{i=1}^n \left( \frac{P(c | \vec{x}_i)}{nP(c)} \right) (x_{i,j} - \mu_{c,j})(x_{i,k} - \mu_{c,k})$$

- based on: our guess of the source for each instance

$$P(c | \vec{x}_i) = \frac{P(\vec{x}_i | c)P(c)}{\sum_{c'=1}^k P(\vec{x}_i | c')P(c')} \quad P(\vec{x}_i | c) = \frac{1}{\sqrt{2\pi|\Sigma_c|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu}_c)^T \Sigma_c^{-1} (\vec{x}_i - \vec{\mu}_c)\right)$$

$$\sum_{j=1}^d \sum_{k=1}^d (x_{i,j} - \mu_{c,j})(\Sigma_c^{-1})_{j,k} (x_{i,k} - \mu_{c,k})$$

Copyright © 2011



## How to pick K?

- Probabilistic model  $L = \log P(x_1 \dots x_n) = \sum_{i=1}^n \log \sum_{k=1}^K P(x_i | k)P(k)$ 
  - tries to “fit” the data (maximize likelihood)
- Pick K that makes L as large as possible?
  - $K = n$ : each data point has its own “source”
  - may not work well for new data points
- Split points into training set T and validation set V
  - for each K: fit parameters of T, measure likelihood of V
  - sometimes still best when  $K = n$
- Occam’s razor: pick “simplest” of all models that fit
  - Bayes Inf. Criterion (BIC):  $\max_p \{ L - \frac{1}{2} p \log n \}$
  - Akaike Inf. Criterion (AIC):  $\min_p \{ 2 p - L \}$

L ... likelihood, how well our model fits the data  
p ... number of parameters  
how “simple” is the model

Copyright © 2011 Victor Lavrenko

## Summary

- Walked through 1-d version
  - works for higher dimensions
    - d-dimensional Gaussians, can be non-spherical
  - works for discrete data (text)
    - d-dimensional multinomial distributions (pLSI)
- Maximizes likelihood of the data:  $P(x_1 \dots x_n) = \prod_{i=1}^n \sum_{k=1}^K P(x_i | k)P(k)$
- Very similar to K-means
  - sensitive to starting point, converges to a local maximum
  - convergence: when change in  $P(x_1 \dots x_n)$  is sufficiently small
  - cannot discover K (likelihood keeps growing with K)

Copyright © 2011 Victor Lavrenko