Last updated: November 6, 2012

MIXTURE MODELS AND EM

J. Elder

- Some of these slides were sourced and/or modified from:
 - Christopher Bishop, Microsoft UK
 - Simon Prince, University College London
 - Sergios Theodoridis, University of Athens & Konstantinos Koutroumbas, National Observatory of Athens



Mixture Models and EM: Topics

Probability & Bayesian Inference

- 1. Intuition
- 2. Equations
- 3. Examples
- 4. Applications



Mixture Models and EM: Topics

Probability & Bayesian Inference

1. Intuition

- 2. Equations
- 3. Examples
- 4. Applications



What do we do if a distribution is not wellapproximated by a standard parametric model?







Mixtures of Gaussians

Probability & Bayesian Inference

Combine simple models p(x) into a complex model:

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
Component
Mixing coefficient
$$\forall k : \pi_k \ge 0 \qquad \sum_{k=1}^{K} \pi_k = 1$$

k=1





J. Elder

Mixtures of Gaussians

Probability & Bayesian Inference





Determining parameters μ , σ and π using maximum log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Log of a sum; no closed form maximum.

Solution: use standard, iterative, numeric optimization methods or the expectation maximization algorithm.



8

Probability & Bayesian Inference

$$\boldsymbol{p}(\mathbf{x}_{n}) = \sum_{k=1}^{k} \pi_{k} \mathcal{N}(\mathbf{x}_{n}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$$

- Assumptions
 - **o** for each training observation \mathbf{x}_n there is a hidden variable z_n .
 - $z_n = 1,...,K$ represents which Gaussian x_n came from
- With this interpretation, we have:

$$p(\mathbf{x}_{n}) = \sum_{k=1}^{K} p(\mathbf{x}_{n} \mid \mathbf{z}_{n} = k) p(\mathbf{z}_{n} = k)$$
where
$$p(\mathbf{z}_{n} = k) = \pi_{k}$$
and
$$p(\mathbf{x}_{n} \mid \mathbf{z}_{n} = k) \sim N(\mathbf{x}_{n}; \mu_{k}, \Sigma_{k})$$
Observed Data x
Observed



Probability & Bayesian Inference

$$p\left(\mathbf{x}_{n} \mid \pi_{1}, \dots, \pi_{k}, \mu_{1}, \dots, \mu_{k}, \Sigma_{1}, \dots, \Sigma_{k}\right) = \sum_{k=1}^{k} \pi_{k} \mathcal{N}\left(\mathbf{x}_{n}; \mu_{k}, \Sigma_{k}\right) = \sum_{k=1}^{k} p\left(\mathbf{z}_{n} = k\right) p\left(\mathbf{x}_{n} \mid \mathbf{z}_{n} = k\right)$$





10

Probability & Bayesian Inference $p\left(\mathbf{x}_{n} \mid \pi_{1}, \dots, \pi_{k}, \mu_{1}, \dots, \mu_{k}, \Sigma_{1}, \dots, \Sigma_{k}\right) = \sum_{k=1}^{k} \pi_{k} N\left(\mathbf{x}_{n}; \mu_{k}, \Sigma_{k}\right) = \sum_{k=1}^{k} p\left(\mathbf{z}_{n} = k\right) p\left(\mathbf{x}_{n} \mid \mathbf{z}_{n} = k\right)$

OUR GOAL

11

To estimate the parameters θ :

The means μ_{k}

The covariances Σ_{μ}

The weights (mixing coefficients) π_{μ} for all K components of the model.



THING TO NOTICE

If we knew the hidden variables z_n for the training data it would be easy to estimate parameters θ - just estimate individual Gaussians separately.



Probability & Bayesian Inference

THING TO NOTICE #2:

If we knew the parameters θ it would very easy to estimate the posterior distribution over each hidden variable z_n using Bayes' rule:





Expectation Maximization

Probability & Bayesian Inference

Chicken and egg problem:

- could find $z_{1...N}$ if we knew θ
- could find θ if we knew $z_{1...N}$

Solution: Expectation Maximization (EM) algorithm (Dempster, Laird and Rubin 1977)

EM for Gaussian mixtures can be viewed as alternation between 2 steps:

- 1. Expectation Step (E-Step)
 - For fixed θ find posterior distribution over responsibilities $z_{1...N}$
- 2. Maximization Step (M-Step)
 - Now use these posterior probabilities to re-estimate heta



K-Means: A Poor-Man's Approximation

Probability & Bayesian Inference

- The K-means algorithm is very similar to EM, except that
 - We assume the mixing coefficients are the same for each component.
 - We assume each component is isotropic with common variance.
 - We do not admit any uncertainty about the responsibilities at each iteration.



14

K-Means: A Poor-Man's Approximation

Probability & Bayesian Inference

- □ K-Means consists of the following steps:
 - 1. Randomly select the mean for each component.
 - 2. Associate each input with the nearest component mean.
 - 3. Recompute the means based upon the new association.
- □ Steps 2 and 3 are alternated until convergence.



15

Example

(a) (b) (C) 2 2 2 х 0 0 0 × × -2-2-22 -2 0 2 -20 2 -2 0 (e) (d) 2 2 2 0 0 0 -2 -2-2-2 0 2 -20 2 -2 2 0 (g) (h) (i) 2 2 2 0 0 0 -2 -20 2 -20 2 -2 0 2 -2



Mixture Models and EM: Topics

Probability & Bayesian Inference

- 1. Intuition
- 2. Equations
- 3. Examples
- 4. Applications



Responsibilities

Probability & Bayesian Inference

The **responsibility** r_{nk} of component k for observation x_n is the posterior probability that component k generated x_n .

Responsibility
$$r_{nk} \triangleq P(z_n = k \mid \mathbf{x}_n; \theta)$$

Responsibilities Update Equation:

Responsibility
$$r_{nk} = \frac{p(\mathbf{x}_n \mid \mathbf{z}_n = k; \theta) \pi_k(t)}{\sum_{k=1}^{K} p(\mathbf{x}_n \mid \mathbf{z}_n = k; \theta) \pi_k(t)}$$



□ Let N_k = effective number of observations explained by component k:

$$N_{k} \triangleq \sum_{n=1}^{N} r_{nk}$$



Example: Mixture of Gaussians

Probability & Bayesian Inference

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- Differentiating does not allow us to isolate the parameters analytically.
- However, it does generate equations that can be used for iterative estimation of these parameters by EM.



19

Example: Mixture of Gaussians

Probability & Bayesian Inference

Parameter Update Equations:

$$\mu_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} r_{nk} \mathbf{x}_{n}$$

$$\Sigma_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} r_{nk} \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k} \right) \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k} \right)^{t}$$

$$\pi_{k} = \frac{1}{N} \sum_{k=1}^{N} r_{nk}$$



Mixture of (Multivariate) Gaussians

Probability & Bayesian Inference

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- □ Here we will derive the formula for the mixing coefficents π_k .
- Textbook Problem 11.2: Derive the EM formulae for the mean μ_k and covariance Σ_k .
 - Please do this at home!
 - **u** When solving for $\Sigma_{k'}$ recall that:

$$\frac{d}{dA}|A| = |A|A^{-t}$$
 and $\frac{d}{dA}x^tAx = xx^t$



21

- EM typically takes longer to converge than K-Means, and the computations are more elaborate.
- It is thus common to use the solution to K-means to initialize the parameters for EM.



- EM is guaranteed to monotonically increase the likelihood.
- However, since in general the likelihood is nonconvex, we are not guaranteed to find the globally optimal parameters.



End of Lecture

October 24, 2012

Mixture Models and EM: Topics

Probability & Bayesian Inference

1. Intuition

25

- 2. Equations
- 3. Examples
- 4. Applications



Bivariate Gaussian Mixture Example

Probability & Bayesian Inference





26

2-Component Bivariate MATLAB Example

Probability & Bayesian Inference





27

2-Component Bivariate MATLAB Example

Probability & Bayesian Inference

%update responsibilities

for i=1:k
 p(:,i)=alphas(i).*mvnpdf(x,mu(i,:),squeeze(S(i,:,:)));
end
p=p./repmat((sum(p,2)),1,k);

%update parameters

```
for i=1:k
    Nk=sum(p(:,i));
    mu(i,:)=p(:,i)'*x/Nk;
    dev=x-repmat(mu(i,:),N,1);
    S(i,:,:)=(repmat(p(:,i),1,D).*dev)'*dev/Nk;
    alphas(i)=Nk/N;
```







Old Faithful Example

Probability & Bayesian Inference



Duration of eruption (min)



Face Detection Example: 2 Components

Probability & Bayesian Inference

0.4999 Prior 0.5001 Mean Face Model **Parameters** Standard deviation 0.5325 0.4675 Prior Mean Non-Face Model **Parameters** Standard deviation

30

Each component is still assumed to have diagonal covariance.

The face model and non-face model have divided the data into two clusters. In each case, these clusters have roughly equal weights.

The primary thing that these seem to have captured is the photometric (luminance) variation.

Note that the standard deviations have become smaller than for the single Gaussian model as any given data point is likely to be close to one mean or the other.

Machine Learning and Pattern Recognition

Results for MOG 2 Model

Probability & Bayesian Inference



Performance improves relative to a single Gaussian model, although it is not a dramatic improvement.

We have a better description of the data likelihood.

5327 Introduction to Machine Learning and Pattern Recognition

MOG 5 Components



MOG 10 Components

Probability & Bayesian Inference











CSE 4404/5327 Introduction to Machine Learning and Pattern Recognition

33

Results for Mog 10 Model

Probability & Bayesian Inference



Performance improves slightly more, particularly at low false alarm rates.

J. Elder

Background Subtraction

Probability & Bayesian Inference



Test Image



GOAL : (i) Learn background model (ii) use this to segment regions where the background has been occluded



What if the scene isn't static?

Probability & Bayesian Inference



Gaussian is no longer a good fit to the data.

Not obvious exactly what probability model would fit better.



Background Mixture Model



Background Subtraction Example

Probability & Bayesian Inference





