

Latent Variable View of EM

Sargur Srihari

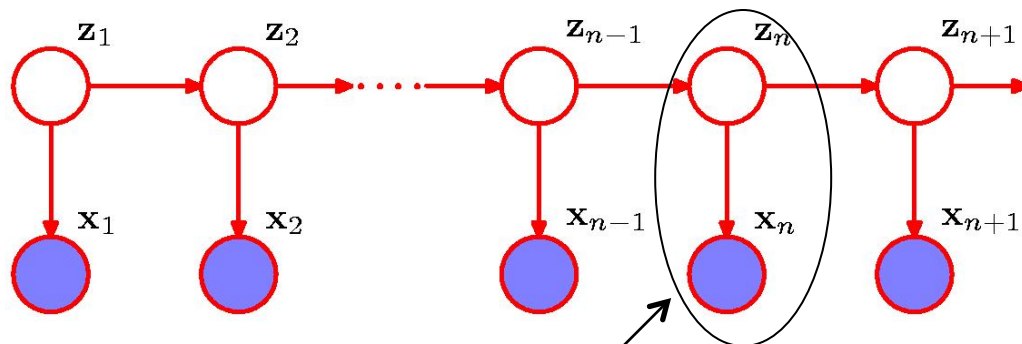
srihari@cedar.buffalo.edu

Examples of latent variables



- Joint distribution is $p(\mathbf{x}, \mathbf{z})$
- We don't have values for z

2. Hidden Markov Model



- A single time slice is a mixture with components $p(\mathbf{x}|\mathbf{z})$
- An extension of mixture model
 - Choice of mixture component depends on choice of mixture component for previous distribution
- Latent variables are multinomial variables z_n
 - That describe component responsible for generating x_n

Another example of latent variables

3. Topic Models (Latent Dirichlet Allocation)

- In NLP unobserved groups explain why some observed data are similar
- Each document is a mixture of various topics (latent variables)
- Topics generate words
 - CAT-related: milk, meow, kitten
 - DOG-related: puppy, bark, bone
- Multinomial distributions over words with Dirichlet priors

Main Idea of EM

- Goal of EM is:
 - find maximum likelihood models for distributions $p(\mathbf{x})$ that have latent (or missing) data
 - E.g., GMMs, HMMs
 - In case of Gaussian mixture models $p(\mathbf{x}) = \sum_{k=1}^K \pi_k N(\mathbf{x} | \mu_k, \Sigma_k)$
 - We have a complex distribution of observed variables \mathbf{x}
 - We wish to estimate its parameters μ_k, Σ_k, π_k
- Introduce latent variables \mathbf{z} so that $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$
 - joint distribution $p(\mathbf{x}, \mathbf{z})$ is more tractable (since we know forms of components) $p(\mathbf{x} | z_k = 1) = N(\mathbf{x} | \mu_k, \Sigma_k)$
 - Complicated form from simpler components
- The original distribution is obtained by marginalizing the joint distribution

Alternative View of EM

- This view recognizes key role of latent variables

- Observed data - matrix $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$
- Latent Variables - matrix $Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix}$

- where n^{th} row is the sample vector $x_n^T = [x_{n1} \ x_{n2} \ \dots \ x_{nD}]$
- with corresponding row $z_n^T = [z_{n1} \ z_{n2} \ \dots \ z_{nK}]$

- Goal of EM algorithm is to find maximum likelihood solution for $p(X)$ given some X
- When we do not have Z

Likelihood with Latent Variables

- Likelihood function (from sum rule) is

$$p(X | \theta) = \sum_Z p(X, Z | \theta)$$

– where θ is the set of all model parameters

- E.g., means, covariances, responsibilities

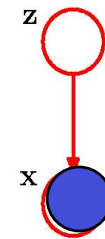
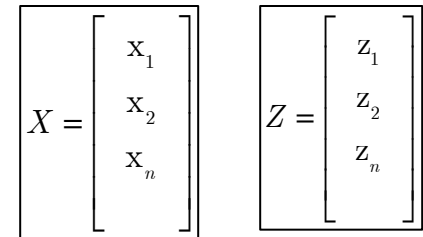
- Log-likelihood function is

$$\ln p(X | \theta) = \ln \left\{ \sum_Z p(X, Z | \theta) \right\}$$

- We wish to find θ that maximizes this
- Joint likelihood function can be written as

$$p(X, Z | \theta) = p(Z | X, \theta) p(X | \theta)$$

- We choose this form from the graph since we know X and not Z



Complication due to Latent Variables

- Log likelihood function is

$$\ln p(X | \theta) = \ln \left\{ \sum_Z p(X, Z | \theta) \right\}$$

Summation inside brackets
due to marginalization
Not due to log-likelihood

- Key Observation:

– Summation over latent variables appears inside logarithm

- Even if joint distribution $p(X, Z | \theta)$ belongs to exponential family the marginal distribution $p(X | \theta)$ does not
 - Taking log of Sum of Gaussians does not give simple quadratic
- Results in complicated expressions for maximum likelihood solution, i.e., what value of θ maximizes the likelihood

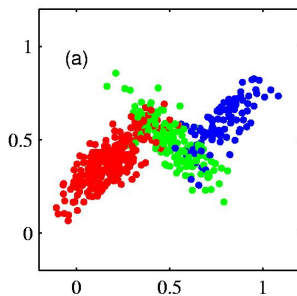
Complete and Incomplete Data Sets

Log-likelihood function is:

$$\ln p(X | \theta) = \ln \left\{ \sum_Z p(X, Z | \theta) \right\}$$

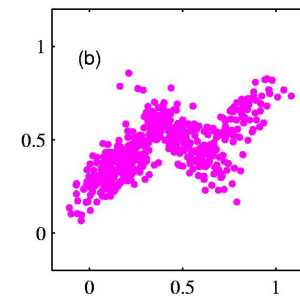
Complete Data $\{X, Z\}$

- For each observation in X we know corresponding value of latent variable in Z
- Since we can evaluate $p(X, Z | \theta) = p(Z | X, \theta) p(X | \theta)$
 - maximization over θ is straightforward



Incomplete Data $\{X\}$

- Actual data set
- Since we do not know Z
 - We cannot evaluate $p(X, Z | \theta)$ to maximize over θ



Maximizing Expectation of $p(X, Z|\theta)$

- Since we don't have the complete data set $\{X, Z\}$ to evaluate $p(X, Z|\theta)$ we instead evaluate its expectation
- Since we are given X , for a given θ we first determine the distribution of the latent variables $p(Z|X, \theta)$
- Then the expected log-likelihood of complete data is

$$E [\ln p(X, Z | \theta)] = \sum_Z p(Z | X, \theta) \ln p(X, Z | \theta)$$

Summation is due to expectation not sum rule!

- We maximize this by considering every value of θ
 - And for each θ by summing over each value of Z
 - Since logarithm acts directly on the joint $p(X, Z|\theta)$ and not on a summation it is tractable

E and M Steps

- *E Step*: Estimate the missing values
 - Use current parameter value θ^{old} to find the posterior distribution of the latent variables given by

$$p(Z \mid X, \theta^{\text{old}})$$

- *M Step*: Determine revised parameter estimate θ^{new} by *maximizing*

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$$

– *where*

$$Q(\theta, \theta^{\text{old}}) = \sum_Z p(Z \mid X, \theta^{\text{old}}) \ln p(X, Z \mid \theta)$$

Summation due to expectation

- is the *expectation* of $p(X, Z \mid \theta)$ for some general parameter value θ

General EM Algorithm

- Given joint distribution $p(X, Z | \theta)$ over observed variables X and latent variables Z governed by parameters θ goal is to maximize likelihood function $p(X | \theta)$

- Step 1: Choose an initial setting for the parameters θ^{old}

- Step 2: E Step: Evaluate $p(Z | X, \theta^{old})$

- Step 3: M Step: Evaluate θ^{new} given by

$$\theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$$

where $Q(\theta, \theta^{old}) = \sum_Z p(Z | X, \theta^{old}) \ln p(X, Z | \theta)$

- Check for convergence
 - of either log-likelihood or parameter values
- If not satisfied then let $\theta^{old} \leftarrow \theta^{new}$
- Return to Step 2

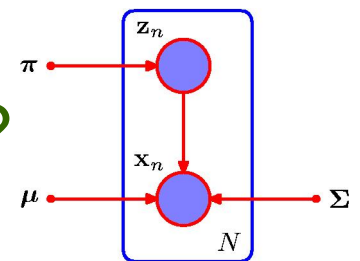
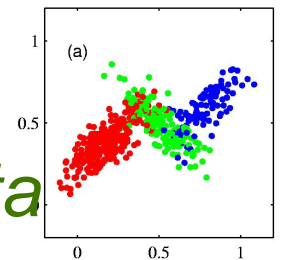
Missing Variables

- EM has been described for maximum likelihood function when there are discrete latent variables
- It can also be applied when there are unobserved variables corresponding to missing values in data set
 - Take the joint distribution of all variables and then marginalize over missing ones
 - EM is then used to maximize corresponding likelihood function
- Method is valid when data is missing at random
 - Not if missing value depends on unobserved values
 - E.g., if quantity exceeds some threshold

Gaussian Mixtures Revisited

- Apply EM (latent variable view) to GMM
- In the E-step we compute
 - Expectation of *log-likelihood of complete data* $\{X, Z\}$ wrt posterior of latent Variables Z
- In the M-step we maximize $Q(\theta, \theta^{old})$ wrt θ
 - Will show that this leads to the same m.l estimates for GMM parameters π, μ, Σ as before

$$Q(\theta, \theta^{old}) = \sum_Z p(Z | X, \theta^{old}) \ln p(X, Z | \theta)$$



Likelihood for Complete Data

- Likelihood function for the complete data set is

$$p(X, Z | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} N(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}}$$

- Log-likelihood is

$$\ln p(X, Z | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \ln \pi_k + \ln N(x_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- Much simpler than log-likelihood for incomplete data:

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

- Maximum likelihood solution for complete data can be obtained in closed form
- Since we don't have values for latent variables, we obtain its expectation wrt the posterior distribution of latent variables

Posterior Distribution of Latent Variables

- From $p(z) = \prod_{k=1}^K \pi_k^{z_k}$ and $p(x|z) = \prod_{k=1}^K N(x|\mu_k, \Sigma_k)^{z_k}$ we have

$$p(Z|X, \mu, \Sigma) \propto \prod_{n=1}^N \prod_{k=1}^K (\pi_k N(x_n|\mu_k, \Sigma_k))^{z_{nk}}$$

- From which we can get the expected value for the indicator variable as

$$E[z_{nk}] = \frac{\pi_k N(x_n|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n|\mu_j, \Sigma_j)} = \gamma(z_{nk})$$

- Substituting into complete log-likelihood:

$$E_Z [\ln p(X, Z|\pi, \mu, \Sigma)] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln N(x_n|\mu_k, \Sigma_k) \}$$

- Final procedure:** choose initial values for $\pi^{old}, \mu^{old}, \Sigma^{old}$
 - Evaluate the responsibilities (E-step)
 - Keep responsibilities fixed and use closed-form solutions for $\pi^{new}, \mu^{new}, \Sigma^{new}$

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

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

$$\pi_k = \frac{N_k}{N}$$

Relation to K-means

- EM for Gaussian mixtures has close similarity to K-means
- K-means performs a hard assignment of data points to clusters
 - Each data point is associated uniquely with one cluster
- EM makes a soft assignment based on posterior probabilities
- K-means does not estimate the covariances of the clusters but only the cluster means