

# Machine Learning

## Week 6: Unsupervised Learning: PCA, Mixture Models, Cluster Analysis

Maheesan Niranjana

School of Electronics and Computer Science  
University of Southampton

Autumn Semester 2016/17

## Week Six: Overview

- Review: Maximum Likelihood and Bayesian Estimation
- Deriving Principal Component Analysis
- Mixture Gaussian Model
- Expectation Maximization (EM) Algorithm
- K-Means Clustering

### Note:

You need not learn the derivations in estimating mixture model parameters by heart. But we need to go through the algebra to gain an insight into the formal basis of the algorithms we use in Machine Learning.

# Unsupervised Learning

- Given:  $\{\mathbf{x}_n\}_{n=1}^N$  (as opposed to  $\{\mathbf{x}_n, f_n\}_{n=1}^N$ )
- We might extract cluster structures
  - Notion of distance between points of data
  - Criterion to determine how many clusters (often from prior knowledge)
  - Underlying probabilistic model
- We might project data onto a subspace
$$\mathbf{x}_n \in \mathcal{R}^d \longrightarrow \mathbf{y}_n \in \mathcal{R}^q$$
  - $q = 2$  helps visualization
  - Subspace representation useful for
    - Data compression
    - Sometimes used to reduce features

Semi Supervised Learning:

$$\{\mathbf{x}_n, f_n\}_{n=1}^{N_1} \text{ and } \{\mathbf{x}_n\}_{n=N_1+1}^{N_2}$$

## Constrained Optimization: Lagrange Multipliers

Problem:

- Maximize  $f(\mathbf{x})$  (with respect to  $\mathbf{x}$ )
- Subject to  $g(\mathbf{x}) = c$

Method:

- Construct a function

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda (g(\mathbf{x}) - c)$$

- $L$  is called a Lagrangian;  $\lambda$  is called a Lagrange Multiplier
- The problem now is an unconstrained problem; we look for turning points by

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} - \lambda \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} = 0$$

# Example of Lagrange Multipliers:

## Principal Component Analysis

- $N$  data  $\mathbf{x}_n \in \mathcal{R}^d$  distributed with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$ .
- Project onto direction  $\mathbf{u}$ ; find the direction that maximizes projected variance.
- Projected variance is  $\mathbf{u}^t \mathbf{C} \mathbf{u}$
- We are only interested in the direction; not in increasing the projected variance by choosing  $\mathbf{u}$  with large magnitude.
- Set up a constrained optimization problem

$$\max_{\mathbf{u}} \mathbf{u}^t \mathbf{C} \mathbf{u} \quad \text{subject to } \mathbf{u}^t \mathbf{u} = 1$$

- Lagrangian

$$\mathcal{L} = \mathbf{u}^t \mathbf{C} \mathbf{u} - \lambda [\mathbf{u}^t \mathbf{u} - 1]$$

- $\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 0 \implies \mathbf{C} \mathbf{u} = \lambda \mathbf{u}$ ; i.e. principal directions are eigenvectors of covariance

## Mixture Model

We write a mixture of Gaussian densities:

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

- If the mixing proportions  $\pi_k$  satisfy

- $\pi_k \geq 0$
- $\sum_{k=1}^K \pi_k = 1$

$p(\mathbf{x})$  is a proper probability density.

- More powerful model – useful when data is multi-modal
- Parameters are: proportions, means and covariance matrices
- Parameter estimation ( $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ ) is not easy.
- $z_{nk}$ : association of  $n^{\text{th}}$  data to  $k^{\text{th}}$  mode unknown (latent)

## Log Likelihood:

( $\Delta$  represents all the means and covariances and  $\pi$  is a vector holding all the  $\pi_i$ 's)

$$\begin{aligned}\mathcal{L} &= \log p(\mathbf{X}|\Delta, \pi) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k)\end{aligned}$$

Note log of sums of variables; inconvenient to work with.

### Jensen's Inequality

$$\log E_{p(z)} \{f(z)\} \geq E_{p(z)} \{\log f(z)\}$$

- Introduce a new variable  $q_{nk}$ :  $q_{nk} \geq 0$  and  $\sum_{k=1}^K q_{nk} = 1$   
At every data  $n$ , we are defining a new probability distribution over the  $K$  components of the mixture model.
- We multiply and divide by the new variable:

$$\mathcal{L} = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k) \frac{q_{nk}}{q_{nk}}$$

We now treat the weighted sum as expectation over the newly introduced distribution:

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \log \sum_{k=1}^K q_{nk} \frac{\pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k)}{q_{nk}} \\ &= \sum_{n=1}^N \log \mathbf{E}_{q_{nk}} \left\{ \frac{\pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k)}{q_{nk}} \right\}\end{aligned}$$

That gives a form in which Jensen's inequality may be applied.

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \log \mathbf{E}_{q_{nk}} \left\{ \frac{\pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k)}{q_{nk}} \right\} \\ &\geq \sum_{n=1}^N \mathbf{E}_{q_{nk}} \left\{ \log \frac{\pi_k p(\mathbf{x}_n|\mu_k, \Sigma_k)}{q_{nk}} \right\}\end{aligned}$$

What we do is to optimize this lower bound, rather than the log likelihood itself, with respect to the unknowns  $\{q_{nk}, \pi_k, \mu_k, \Sigma_k\}$ .

$$\begin{aligned}
 \mathcal{B} &= \sum_{n=1}^N \mathbf{E}_{q_{nk}} \left\{ \log \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right\} \\
 &= \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left( \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right) \\
 &= \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}
 \end{aligned}$$

The task now is to maximize this ( $\mathcal{B}$ ) with respect to the unknowns.

Maximize with respect to  $\pi_k$

- Only the first term depends on  $\pi_k$
- But we need to constrain the solutions for  $\pi_k$  because  $\sum_{k=1}^K \pi_k = 1$ .
- Set up the Lagrangian:

$$B_1 = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k - \lambda \left( \sum_{k=1}^K \pi_k - 1 \right)$$

- Differentiate and equate to zero:

$$\frac{\partial B_1}{\partial \pi_k} = \frac{\sum_{n=1}^N q_{nk}}{\pi_k} - \lambda = 0$$

$$\sum_{n=1}^N q_{nk} = \lambda \pi_k$$

Sum both sides over  $k$

$$\begin{aligned}
 \sum_{k=1}^K \sum_{n=1}^N q_{nk} &= \lambda \sum_{k=1}^K \pi_k \\
 N &= \lambda
 \end{aligned}$$

Hence  $\pi_k = \frac{1}{N} \sum_{n=1}^N q_{nk}$

## Maximizing with respect to $\mu_k$

- Only the second term of the bound  $\mathcal{B}$  depends on  $\mu_k$

$$\begin{aligned} B_2 &= \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left( \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x}_n - \mu_k)^t \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \right) \right) \\ &= -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left( (2\pi)^{d/2} |\Sigma_k| \right) - \frac{1}{2} (\mathbf{x}_n - \mu_k)^t \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) \end{aligned}$$

- Differentiate:

$$\frac{\partial B_2}{\partial \mu_k} = \sum_{n=1}^N q_{nk} \Sigma_k^{-1} (\mathbf{x}_n - \mu_k).$$

- Equate to zero and re-arrange terms

$$\mu_k = \frac{\sum_{n=1}^N q_{nk} \mathbf{x}_n}{\sum_{n=1}^N q_{nk}}$$

## Maximizing with respect to $\Sigma_k$

- Again only the second term matters; differentiating with respect to  $\Sigma_k$  is tricky (we'll not do this)
- Answer

$$\Sigma_k = \frac{\sum_{n=1}^N q_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^t}{\sum_{n=1}^N q_{nk}}$$

Updating  $q_{nk}$  needs to recognize the constraints (sum to one)

$$B = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk} - \lambda \left( \sum_{k=1}^K q_{nk} - 1 \right)$$

$$\frac{\partial B}{\partial q_{nk}} = \log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - (1 + \log q_{nk}) - \lambda$$

$$\begin{aligned} 1 + \log q_{nk} + \lambda &= \log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ \exp(\log q_{nk} + (\lambda + 1)) &= \exp(\log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \\ q_{nk} \exp(\lambda + 1) &= \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Sum over mixture components to get the Lagrange multiplier.

$$\exp(\lambda + 1) \sum_{k=1}^K q_{nk} = \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Because  $q_{nk}$  should sum to one, we have

$$q_{nk} = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

## Summary of Algorithm

$$\begin{aligned} \pi_k &= \frac{1}{N} \sum_{n=1}^N q_{nk} \\ \boldsymbol{\mu}_k &= \frac{\sum_{n=1}^N q_{nk} \mathbf{x}_n}{\sum_{n=1}^N q_{nk}} \\ \boldsymbol{\Sigma}_k &= \frac{\sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^t}{\sum_{n=1}^N q_{nk}} \end{aligned}$$

$$q_{nk} = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Compare with maximum likelihood estimation of parameters of a single Gaussian and with posterior probabilities we studied in Bayesian classification.

# Expectation Maximization

## Auxilliary Variable as Posteriors

Interpret:

- Mixture model as a Gaussian classifier with  $K$  classes
- $\pi_k$  as prior probabilities
- Each of the  $\mathcal{N}(\mu_k, \Sigma_k)$  as class conditional densities / likelihoods.

$$\begin{aligned} p(z_{nk} = 1 | \mathbf{x}_n, \pi, \Delta) &= \frac{p(z_{nk} = 1 | \pi_k) p(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K p(z_{nk} = 1 | \pi_k) p(\mathbf{x}_n | \mu_j, \Sigma_j)} \\ &= q_{nk} \end{aligned}$$

- Each data item has a weighted contribution to the estimation of parameters.
- Unknown assignment  $z_{nk}$ ; **E**xpected value of this unknown assignment is  $q_{nk}$ , the posterior probability
- **M**aximize (the lower bound) to re-estimate parameters.

## K-Means Clustering Algorithm

```
Input:  $\mathbf{X} = \{ \mathbf{x}_n^t \}_{n=1}^N, K$   
Output:  $\mathbf{C}, \text{Idx}$   
initialize:  $\mathbf{C} = \{ \mathbf{c}_j^t \}_{j=1}^K$   
  
repeat  
  . assign  $n^{\text{th}}$  sample to nearest  $\mathbf{c}_j$   
  .  $\text{Idx}(n) = \min_j ||\mathbf{x}_n - \mathbf{c}_j||^2$   
  
  . recompute  $\mathbf{c}_j = \frac{1}{N_j} \sum_{n=j} \mathbf{x}_n$   
  
until no change in  $\mathbf{c}_1, \mathbf{c}_2, \dots \mathbf{c}_k$   
  
return  $\mathbf{C}, \text{Idx}$ 
```



# K-Means as Mixture Gaussian

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma_k)$$

- Set  $\Sigma_k = \sigma_k^2 \mathbf{I}$
- At every iteration, set largest  $q_{nk}$  (largest over  $k$ ) to one and others to zero. Winner take all at each datapoint.
- Computation of  $q_{nk}$  is expectation of latent variable  $z_{nk}$  – **E** step
- Re-estimation of  $\mu_k$  and  $\Sigma_k$  become maximum likelihood estimates from data assigned to each cluster (because  $q_{nk}$  is either one or zero) – **M** step