# LECTURE 9: THE EXPECTATION-MAXIMIZATION ALGORITHM STAT 545: INTRO. TO COMPUTATIONAL STATISTICS

Vinayak Rao Purdue University

September 27, 2016

#### NORMAL

The Multivariate normal (MVN) density on  $\mathbb{R}^d$ :

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

Given N i.i.d. observations  $X \equiv \{x_1, \ldots, x_N\}$ , the likelihood is

$$\mathcal{L}(X|\mu, \Sigma) = \prod_{i=1}^{N} p(x_i|\mu, \Sigma)$$

#### NORMAL

The Multivariate normal (MVN) density on  $\mathbb{R}^d$ :

$$p(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$

Given N i.i.d. observations  $X \equiv \{x_1, \ldots, x_N\}$ , the likelihood is

$$\mathcal{L}(X|\mu, \Sigma) = \prod_{i=1}^{N} p(x_i|\mu, \Sigma)$$

Maximum likelihood estimation (MLE): learn parameters by maximizing  $\mathcal{L}(X|\mu, \Sigma)$  w.r.t  $\mu$  and  $\Sigma$ .

How? Calculate derivatives and set to 0.

# MLE FOR THE MVN

More convenient is the log-likelihood  $\ell(X|\mu, \Sigma) = \log \mathcal{L}(X|\mu, \Sigma)$ :

$$\ell(X|\mu, \Sigma) = \sum_{i=1}^{N} \log p(x_i|\mu, \Sigma)$$

For the Gaussian,

$$\ell(X|\mu, \Sigma) = -\frac{1}{2} \sum_{o=1}^{N} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{N}{2} \log |\Sigma| - \text{const}$$

# MLE FOR THE MVN

More convenient is the log-likelihood  $\ell(X|\mu, \Sigma) = \log \mathcal{L}(X|\mu, \Sigma)$ :

$$\ell(X|\mu, \Sigma) = \sum_{i=1}^{N} \log p(x_i|\mu, \Sigma)$$

For the Gaussian,

$$\ell(X|\mu, \Sigma) = -\frac{1}{2} \sum_{o=1}^{N} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{N}{2} \log |\Sigma| - \text{const}$$
$$\mu_{M_i} = \frac{1}{2} \sum_{o=1}^{N} x_i - \Sigma_{M_i} = \frac{1}{2} \sum_{o=1}^{N} (x_i - \mu_{M_i}) (x_i - \mu_{M_i})^T$$

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^{N} x_i, \quad \Sigma_{ML} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \mu_{ML})(x_i - \mu_{ML})^{T}$$

# MLE FOR THE MVN

More convenient is the log-likelihood  $\ell(X|\mu, \Sigma) = \log \mathcal{L}(X|\mu, \Sigma)$ :

$$\ell(X|\mu, \Sigma) = \sum_{i=1}^{N} \log p(x_i|\mu, \Sigma)$$

For the Gaussian,

$$\ell(X|\mu, \Sigma) = -\frac{1}{2} \sum_{o=1}^{N} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{N}{2} \log |\Sigma| - \text{const}$$
$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^{N} x_i, \quad \Sigma_{ML} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \mu_{ML}) (x_i - \mu_{ML})^T$$

MLE: moment matching (set mean/covariance to that of data) Holds for *exponential family distributions* (later) Consider a K-component discrete distribution  $\pi = (\pi_1, \ldots, \pi_K)$ 

• for 
$$X \sim \pi$$
,  $p(X = c) = \pi_c$ .

Consider a K-component discrete distribution  $\pi = (\pi_1, \ldots, \pi_K)$ 

• for 
$$X \sim \pi$$
,  $p(X = c) = \pi_c$ .

• Equivalently,

$$p(X) = \prod_{c=1}^{K} \pi_i^{\delta(X=c)} = \exp(\sum_{c=1}^{K} \delta(X=c) \log \pi_c)$$

Consider a K-component discrete distribution  $\pi = (\pi_1, \ldots, \pi_K)$ 

• for 
$$X \sim \pi$$
,  $p(X = c) = \pi_c$ .

• Equivalently,

$$p(X) = \prod_{c=1}^{K} \pi_i^{\delta(X=c)} = \exp(\sum_{c=1}^{K} \delta(X=c) \log \pi_c)$$

Given data, what is MLE of  $\pi$ ?

$$\pi_c = \frac{1}{N} \sum_{i=1}^N \delta(x_i = c)$$

Last week we saw a few clustering algorithms.

We also saw some limitations:

- Limited control on the cluster shapes (e.g. spherical clusters in k-means).
- Cannot capture variability across clusters.
- Cannot capture uncertainty in cluster assignments.
- Cannot capture information about relative cluster sizes.

We could adjust loss-function/optimization algorithm. Different approach: directly model data-generation process

- Can capture much richer structure more intuitively.
- Can make predictions about future data.
- Can deal with missing data naturally.

Like k-means, fix the number of clusters to K.

- · component c has parameter  $\theta_c$
- observations from cluster *c* distributed as  $p(x|\theta_c)$

Like k-means, fix the number of clusters to K.

- · component c has parameter  $\theta_c$
- observations from cluster *c* distributed as  $p(x|\theta_c)$

Draw cluster from  $\pi$ , a *K*-component probability vector

Like k-means, fix the number of clusters to K.

- · component c has parameter  $\theta_c$
- observations from cluster *c* distributed as  $p(x|\theta_c)$

Draw cluster from  $\pi$ , a *K*-component probability vector

Today we will consider the mixture of Gaussians (MoG)

- each component is a Gaussian
- $\theta_c = (\mu_c, \Sigma_c)$  is its mean and covariance

To generate the *i*th observation:

 $egin{aligned} \mathcal{C}_i &\sim \pi \ \mathbf{X}_i &\sim \mathcal{N}(\mathbf{X}_i | \mu_{\mathcal{C}_i}, \mathbf{\Sigma}_{\mathcal{C}_i}) \end{aligned}$ 

Sample it's cluster assignment Sample it's value To generate the *i*th observation:

 $egin{aligned} c_i &\sim \pi & & \text{Sample it's cluster assignment} \ x_i &\sim \mathcal{N}(x_i | \mu_{c_i}, \Sigma_{c_i}) & & \text{Sample it's value} \end{aligned}$ 

Joint probability:

$$P(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{c}_1, \dots, \mathbf{c}_N | \pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^N \pi_{c_i} \mathcal{N}(\mathbf{x}_i | \mu_{c_i}, \boldsymbol{\Sigma}_{c_i})$$
$$= \prod_{i=1}^N \prod_{i=1}^K \left[ \pi_j \mathcal{N}(\mathbf{x}_i | \mu_j, \boldsymbol{\Sigma}_j) \right]^{\mathbb{1}(c_i=j)}$$

# MODEL-BASED CLUSTERING



Given observations  $X = \{x_1, \ldots, x_N\}$ , we face three problems:

- What are the  $c_i$ ? (inference)
- What is  $\pi$  and  $\theta_c = (\mu_c, \Sigma_c)$ ? (learning)
- What is *K*? (model selection, not covered here)

#### LEARNING

Imagine we had the cluster assignments C. We saw:

$$P(X_1, \dots, X_N, C_1, \dots, C_N | \pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^N \prod_{j=1}^K \left[ \pi_j \mathcal{N}(X_i | \mu_j, \boldsymbol{\Sigma}_j) \right]^{\mathbb{1}(c_i = j)}$$
$$= \left( \prod_{j=1}^K \left( \pi_j \right)^{N_j} \right) \left( \prod_{j=1}^K \prod_{\{i \text{ s.t. } c_i = j\}} \mathcal{N}(X_i | \mu_j, \boldsymbol{\Sigma}_j) \right)$$

Conveniently separates out into  $\pi$  and component parameters.

$$\log P(X, C|\pi, \mu, \mathbf{\Sigma}) = \left(\sum_{j=1}^{K} N_j \log \pi_j\right) \left(\sum_{j=1}^{K} \sum_{\{i \text{ s.t. } c_i = j\}} \log \mathcal{N}(x_i|\mu_j, \mathbf{\Sigma}_j)\right)$$

$$\log P(X, C|\pi, \mu, \mathbf{\Sigma}) = \left(\sum_{j=1}^{K} N_j \log \pi_j\right) \left(\sum_{j=1}^{K} \sum_{\{i \text{ s.t. } c_i=j\}} \log \mathcal{N}(x_i|\mu_j, \mathbf{\Sigma}_j)\right)$$

MLE requires three sets of 'sufficient statistics':

- The number of observations assigned to each cluster  $(N_i)$ .
- · The empirical mean and mean-square of obs. in each cluster

$$\left(\frac{1}{N_j}\sum_{\{i \text{ s.t. } c_i=j\}} x_i, \frac{1}{N_j}\sum_{\{i \text{ s.t. } c_i=j\}} x_i x_i^{\top}\right)$$

k-means assigns obs. to clusters given parameters. Good idea?

k-means assigns obs. to clusters given parameters. Good idea? For an observation  $x_i$ , what is the posterior probability over  $c_K$ ?

$$\mathcal{P}(c_i|x_i, \pi, \mu, \mathbf{\Sigma}) \propto \mathcal{P}(x_i, c_i|\pi, \mu, \mathbf{\Sigma})$$
  
=  $\left(\prod_{j=1}^{K} \left[\pi_j \mathcal{N}(x_i|\mu_j, \mathbf{\Sigma}_j)\right]^{\mathbb{1}(c_i=j)}\right)$ 

k-means assigns obs. to clusters given parameters. Good idea? For an observation  $x_i$ , what is the posterior probability over  $c_k$ ?

$$\mathcal{P}(c_i|x_i, \pi, \mu, \mathbf{\Sigma}) \propto \mathcal{P}(x_i, c_i|\pi, \mu, \mathbf{\Sigma})$$
  
=  $\left(\prod_{j=1}^{K} \left[\pi_j \mathcal{N}(x_i|\mu_j, \mathbf{\Sigma}_j)\right]^{\mathbb{1}(c_i=j)}\right)$ 

- proportional to prior probability of cluster *j*,  $\pi_j$
- proportional to compatibility obs. *i* with parameters  $\theta_i$

k-means assigns obs. to clusters given parameters. Good idea? For an observation  $x_i$ , what is the posterior probability over  $c_k$ ?

$$\mathcal{P}(c_i|x_i, \pi, \mu, \mathbf{\Sigma}) \propto \mathcal{P}(x_i, c_i|\pi, \mu, \mathbf{\Sigma})$$
  
=  $\left(\prod_{j=1}^{K} \left[\pi_j \mathcal{N}(x_i|\mu_j, \mathbf{\Sigma}_j)\right]^{\mathbb{1}(c_i=j)}\right)$ 

- proportional to prior probability of cluster *j*,  $\pi_j$
- proportional to compatibility obs. *i* with parameters  $\theta_i$

Written as *r<sub>ic</sub>*: 'responsibility' of cluster *c* for obs. *i*.

k-means assigns obs. to clusters given parameters. Good idea? For an observation  $x_i$ , what is the posterior probability over  $c_k$ ?

$$egin{aligned} \mathsf{P}(c_i|x_i,\pi,oldsymbol{\mu},oldsymbol{\Sigma}) &\propto \mathsf{P}(x_i,c_i|\pi,oldsymbol{\mu},oldsymbol{\Sigma}) \ &= \left(\prod_{j=1}^{\mathcal{K}} \left[\pi_j \mathcal{N}(x_i|\mu_j,oldsymbol{\Sigma}_j)
ight]^{\mathbbm{1}(c_i=j)}
ight) \end{aligned}$$

- proportional to prior probability of cluster *j*,  $\pi_j$
- proportional to compatibility obs. *i* with parameters  $\theta_i$

Written as r<sub>ic</sub>: 'responsibility' of cluster c for obs. i.

```
rr <- rep(0,K)
for(i in 1:K) {
    rr[i] <- pi[i] * dmvnorm(x, mu[[i]],sigma[[i]]) }
rr = rr / sum(rr);</pre>
```

How do we update parameters given these probabilities?

$$\mu = \frac{\sum_{i=1}^{N} r_{ic} x_i}{\sum_{i=1}^{N} r_{ic}}$$
$$\Sigma = \frac{\sum_{i=1}^{N} r_{ic} x_i x_i^{\top}}{\sum_{i=1}^{N} r_{ic}}$$
$$\pi_c = \frac{1}{N} \sum_{i=1}^{N} r_{ic}$$

How do we update parameters given these probabilities?

$$\mu = \frac{\sum_{i=1}^{N} r_{ic} x_i}{\sum_{i=1}^{N} r_{ic}}$$
$$\Sigma = \frac{\sum_{i=1}^{N} r_{ic} x_i x_i^{\top}}{\sum_{i=1}^{N} r_{ic}}$$
$$\pi_c = \frac{1}{N} \sum_{i=1}^{N} r_{ic}$$

Compare with when we actually knew the cluster assignments.

- Initialize parameters  $\pi$ , ( $\mu_c$ ,  $\Sigma_c$ ) arbitrarily
- Calculate the observation responsibilities  $r_{ic}$  given parameters
- Update parameters given responsibilities
- Repeat till convergence

Suprising fact: EM converges to stationary point of the log-likelihood:

$$\log P(X|\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \sum_{C=\mathcal{C}} P(X, C|\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Suprising fact: EM converges to stationary point of the log-likelihood:

$$\log P(X|\pi, \mu, \mathbf{\Sigma}) = \log \sum_{C=C} P(X, C|\pi, \mu, \mathbf{\Sigma})$$

Can directly calculate gradients w.r.t. parameters and optimize. Doable but messy: Suprising fact: EM converges to stationary point of the log-likelihood:

$$\log P(X|\pi, \mu, \mathbf{\Sigma}) = \log \sum_{C=C} P(X, C|\pi, \mu, \mathbf{\Sigma})$$

Can directly calculate gradients w.r.t. parameters and optimize. Doable but messy:

- Sums inside logarthms is inconvenient.
- Need to calculate gradients w.r.t. covariance matrices.
- Need to choose step sizes.