

Thursday, 11/12/2020

EE 660

MACHINE LEARNING
FROM SIGNALS:
FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

Lecture 24

Announcements

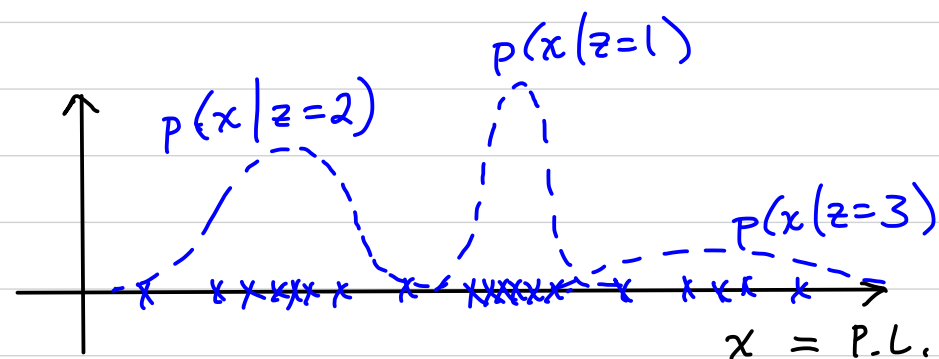
- Homework 9 will be posted

Today's topics

- Unsupervised learning (USL) (part 2)
 - Expectation Maximization (EM)
 - Mixture models in USL
 - EM algorithm and equations
 - Example
 - Similarity and dissimilarity measures for clustering

Mixture Models for USL

- Model each cluster as a pdf with unknown parameters
 - usually for continuous features $x_j \in \mathbb{R}$



$$p(\underline{x} | \underline{\theta}) = \sum_{k=1}^K p(\underline{x} | z=k, \underline{\theta}) \underbrace{p(z=k | \underline{\theta})}_{\pi_k}, \quad \begin{array}{l} z = \text{cluster index (Label)} \\ K = \text{total \# of clusters} \end{array}$$

↓
mixture

↓
use $\underline{\theta}_k =$ parameters for cluster k

$$(1) \quad p(\underline{x} | \underline{\theta}) = \sum_{k=1}^K \pi_k p(\underline{x} | z=k, \underline{\theta}_k), \quad \text{with } \sum_{k=1}^K \pi_k = 1$$

↑
our model for cluster k
prior or mixing parameter (weight) for cluster k .

- Goal: find MLE of $\underline{\theta}$, or $\underline{\theta}_k$ for $k=1, 2, \dots, K$

- Likelihood: $p(\mathcal{D} | \underline{\theta}) = \prod_{i=1}^N p(\underline{x}_i | \underline{\theta})$

◦ Generally not solvable analytically \Rightarrow Use EM

EM for Clustering Using Mixture Models

- Same basic algorithm as EM for SSL

$$\text{Let } \mathcal{D} = \{\underline{x}_i\}_{i=1}^N, \quad \mathcal{H} = \{z_i\}_{i=1}^N.$$

z_i is cluster label for \underline{x}_i , $z_i \in \{1, 2, \dots, K\}$.

- Algorithm (EM for estimating $\underline{\theta}$ and \mathcal{H})

1. Initialize $t=0$ and $\underline{\theta}^{(0)}$

2. Iterate (index $t: 0, 1, 2, \dots, T-1$)

2.1 E step: Compute $p(\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)})$

2.2 M step: Find $\underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)}} \{ \ln p(\mathcal{D}, \mathcal{H} | \underline{\theta}) \}$

2.3 Halt when $p(\mathcal{D} | \underline{\theta}^{(t+1)})$ converges

3. Output $\hat{\underline{\theta}} = \underline{\theta}^{(T)}$

Equations for E step

$$p(\mathcal{H} | \mathcal{D}, \underline{\theta}) = \prod_{i=1}^N p(z_i | \underline{x}_i, \underline{\theta})$$

$$p(z_i | \underline{x}_i, \underline{\theta}) = \frac{p(\underline{x} | z=k, \underline{\theta}_k) p(z=k | \underline{\theta}_k)}{\sum_{k'=1}^K p(\underline{x} | z=k', \underline{\theta}_{k'}) p(z=k' | \underline{\theta}_{k'})}$$

our model
 π_k

→ soft label (responsibility) $\gamma_{ik}^{(t)} = p(z_i=k | \underline{x}_i, \underline{\theta}_k^{(t)})$

↑ cluster index
data pt. index

Equations for M step

$$p(\mathcal{D}, \mathcal{H} | \underline{\theta}) = \prod_{i=1}^N p(\underline{x}_i, z_i | \underline{\theta}) = \prod_{i=1}^N \underbrace{p(\underline{x}_i | z_i, \underline{\theta}) p(z_i | \underline{\theta})}_{\text{our model}}$$

for $z_i=k$: $p(\underline{x}_i | z_i=k, \underline{\theta}_k) \pi_k$

$$\rightarrow \underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ \sum_{\mathcal{H}} p(\mathcal{H} | \mathcal{D}, \underline{\theta}^{(t)}) \ln p(\mathcal{D}, \mathcal{H} | \underline{\theta}) \right\}$$

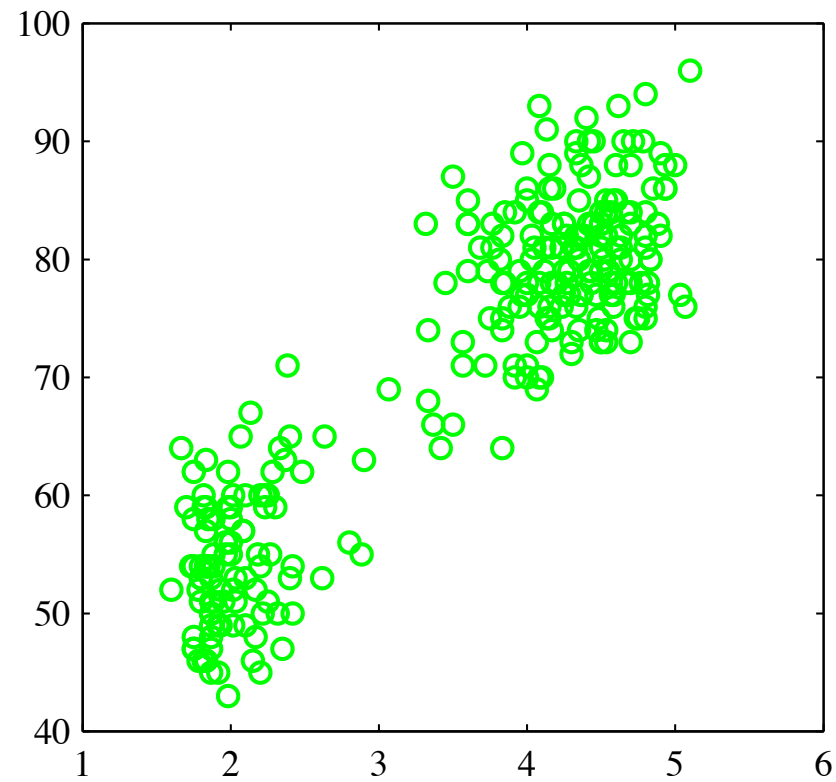
Comments

1. Algorithm characteristics - same as for SSL EM.
2. Choice of $\underline{\theta}^{(0)}$ if no prior knowledge?
 - Use a simpler clustering algorithm (e.g., Kmeans), then use its resulting clusters to calculate $\underline{\theta}^{(0)}$ for EM.
 - or ◦ Try many different (random) $\underline{\theta}^{(0)}$; compare the results using $p(\mathcal{D} | \underline{\theta})$.
3. Quality of clustering result depends on how appropriate the model is for the data.

Example

"Old Faithful" geyser.

Old Faithful Data



Time to next eruption (minutes) vs. duration of eruption (minutes)

From Bishop Fig. A.5

Model

Gaussian mixture model (GMM) using EM, $K=2$.

$$p(\underline{x} | \underline{\theta}) = \sum_{k=1}^2 p(\underline{x} | z=k, \underline{\theta}_k) p(z=k | \underline{\theta}_k) = \sum_{k=1}^2 \underbrace{N(\underline{x} | \underline{\mu}_k, \underline{\Sigma}_k)}_{\underline{\theta}_k} \pi_k$$

Algorithm steps

(a) $\underline{\theta}^{(0)}$: $\underline{\Sigma}_{z=1}^{(0)} = \underline{\Sigma}_{z=2}^{(0)} = \underline{I}$. $\underline{\mu}_k^{(0)}$ is shown.

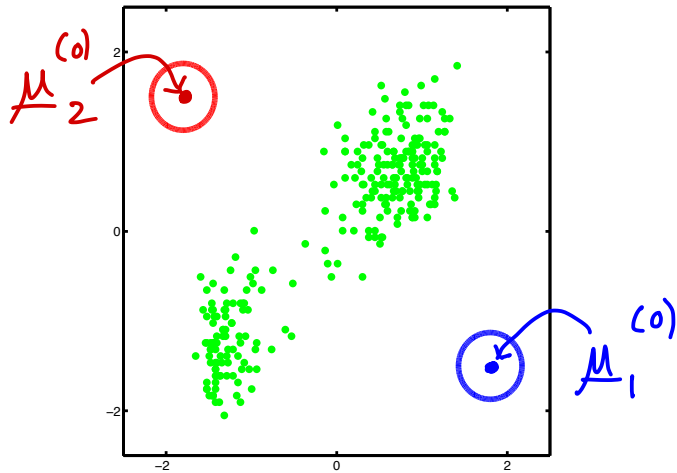
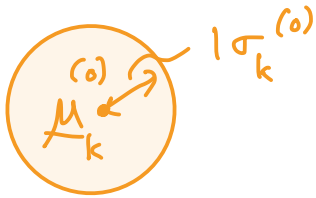
(b) First E step result: $\gamma_{ik}^{(0)} = p(z_i=k | \underline{x}_i, \underline{\theta}_k^{(0)})$

blue ink $\propto \gamma_{i,k=1}$; red ink $\propto \gamma_{i,k=2}$

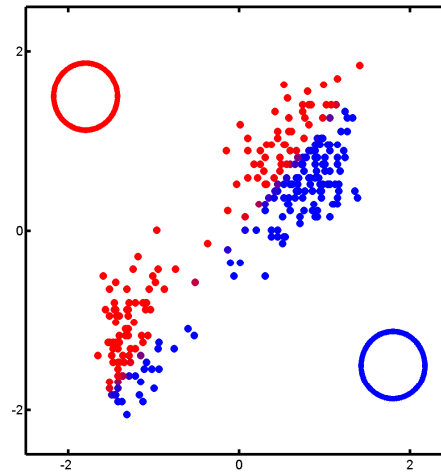
(c) After first M step: $\underline{\theta}^{(1)}$ gives new $\underline{\mu}_k^{(1)}$, $\underline{\Sigma}_k^{(1)}$

new Gaussians $p(\underline{x} | z=k, \underline{\theta}_k)$ are shown.

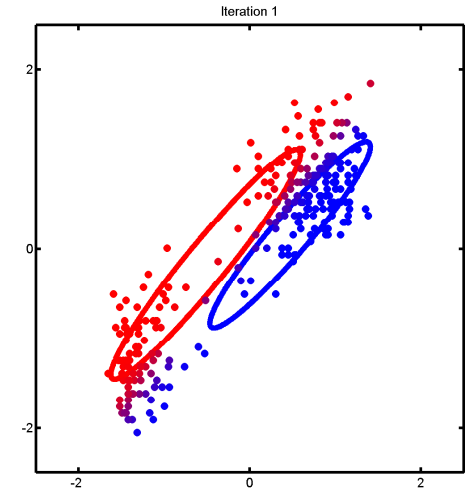
EM Example: Old Faithful Data



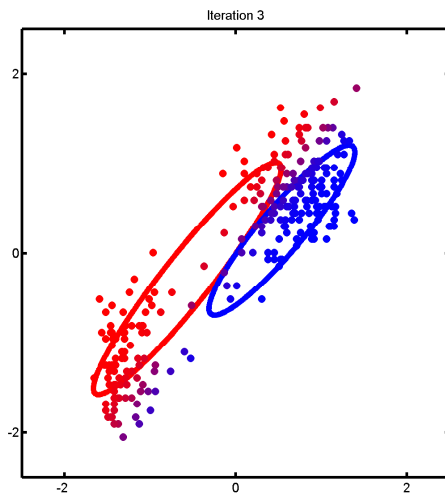
(a) Initial condition



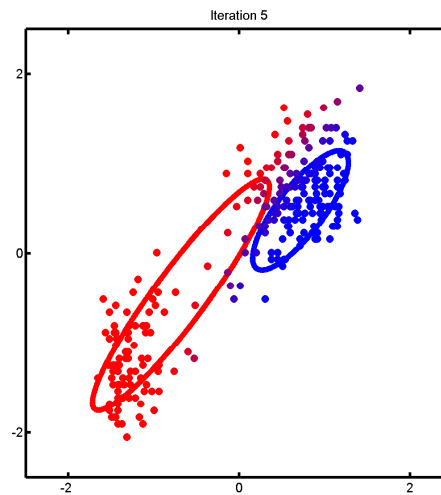
(b) After first E step



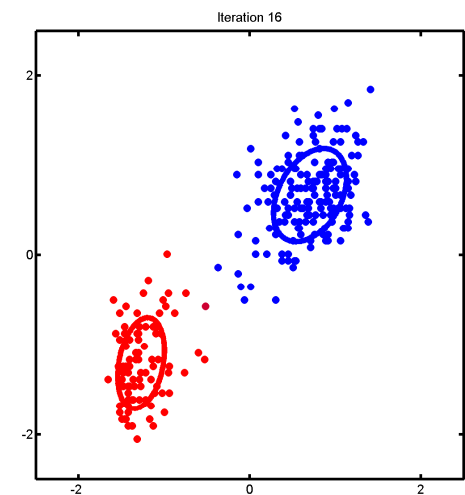
(c) After first M step



(d) After 3 iterations



(e) After 5 iterations



(f) After 16 iterations

Similarity and Dissimilarity Measures

Need measures for similarity or dissimilarity between:

-
- 2 points $\underline{x}_i, \underline{x}_j$.
 - A point and a cluster
 - 2 clusters

Need measure for quality of a partitioning [difficult problem]

Let $\Delta(\underline{x}_i, \underline{x}_{i'}) = d_{ii'}$ denote a dissimilarity function

Can use a distance function, e.g.:

$$\Delta(\underline{x}_i, \underline{x}_{i'}) = d_{ii'} = \sum_{j=1}^D \Delta_j(x_{ij}, x_{i'j})$$

$$\Delta_j(x_{ij}, x_{i'j}) \text{ can be: } (x_{ij} - x_{i'j})^2 \quad (\text{Eucl. dist.})^2$$

$$|x_{ij} - x_{i'j}| \quad (l_1 \text{ norm}$$

or city block dist.
or Manhattan dist.)

For nominal features (symbolic, categorical, or labels):

$$\begin{aligned}\Delta(x_{ij}, x_{i'j}) &= \# \text{ of features that are different} \\ &= \sum_{j=1}^D \mathbb{I}(x_{ij} \neq x_{i'j}) \\ &= \text{Hamming distance.}\end{aligned}$$

Can let $\Delta(\underline{x}_i, \underline{x}_{i'}) = \text{this or some other } d_{ii'}, \text{ s.t. } d_{ii'} \geq 0 \ \forall i, i'$
 and $d_{ii} = 0 \ \forall i$.
 or $f(d_{ii'})$, $f = \text{any monotonically increasing fcn.}$