

Machine Learning

4771

Instructor: Tony Jebara

Topic 12

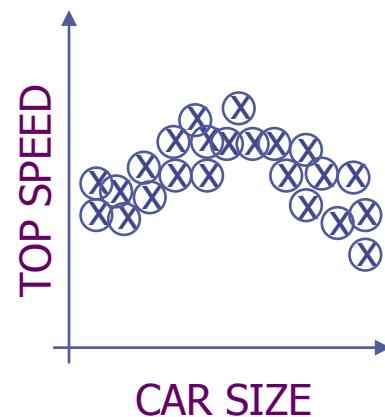
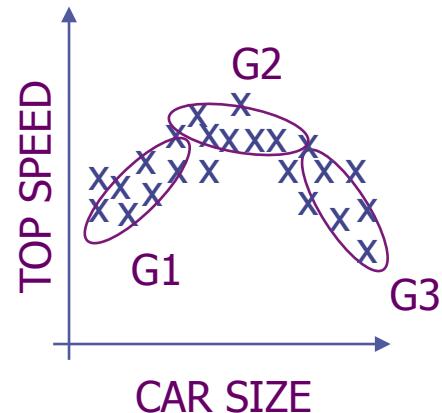
- Mixture Models and Hidden Variables
- Clustering
- K-Means
- Expectation Maximization

Mixtures for More Flexibility

- With mixtures (e.g. mixtures of Gaussians) we can handle more complicated (e.g. multi-bump, nonlinear) distributions.

subpopulations: G1=compact car
 G2=mid-size car
 G3=cadillac

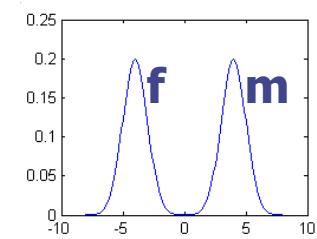
- In fact, if we have enough Gaussians (maybe infinite) we can approximate any distribution...



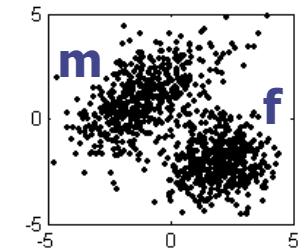
Mixtures as Hidden Variables

- Consider a dataset with K subpopulations but don't know which subpopulation each point belongs to

I.e. looking at height of adult people, we see $K=2$ subpopulations: males & females



I.e. looking at weight and height of people we see $K=2$ subpopulations: males & females



- Because of the 'hidden' variable (y can be 1 or 2), these distributions are not Gaussians but **Mixture of Gaussians**

Unsupervised

$$\begin{aligned}
 p(\vec{x}) &= \sum_y p(\vec{x}, y) = \sum_y p(y) p(\vec{x} | y) = \sum_y \pi_y N(\vec{x} | \vec{\mu}_y, \Sigma_y) \\
 \text{probability that } \vec{x}_i \text{ occurs} &= \sum_{y=1}^K \pi_y \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma_y|}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu}_y)^T \Sigma_y^{-1} (\vec{x} - \vec{\mu}_y)\right)
 \end{aligned}$$

pop.; probability of g; labels.

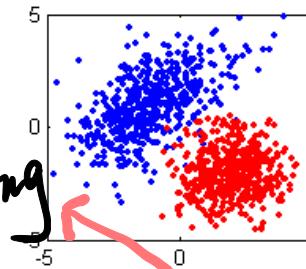
Unlabeled data → Clustering

- Recall classification problem:
maximize the log-likelihood of
data given models:

$$l = \sum_{n=1}^N \log p(\vec{x}_n, y_n | \pi, \mu, \Sigma)$$

$$= \sum_{n=1}^N \log \pi_{y_n} N(\vec{x}_n | \vec{\mu}_{y_n}, \Sigma_{y_n})$$

*Supervised learning
(with label)
origin*

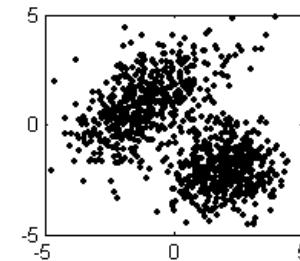


$$= \log \prod_{i=1}^N p(y_i | \pi, \mu, \Sigma)$$

- If we don't know the class
treat it as a hidden variable
maximize the log-likelihood with
unlabeled data: *not* $p(\vec{x}_n, y | \pi, \mu, \Sigma)$

$$l = \sum_{n=1}^N \log p(\vec{x}_n | \pi, \mu, \Sigma) = \sum_{n=1}^N \log \sum_{y=1}^K p(\vec{x}_n, y | \pi, \mu, \Sigma)$$

$$= \sum_{n=1}^N \log (\pi_1 N(\vec{x}_n | \vec{\mu}_1, \Sigma_1) + \dots + \pi_K N(\vec{x}_n | \vec{\mu}_K, \Sigma_K))$$

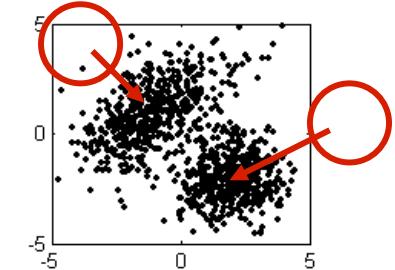


*Unsupervised
(without label)*

- Instead of classification, we now have a **clustering** problem
- original log likelihood = $\log \prod_{i=1}^N p(\vec{x}_i | \pi, \mu, \Sigma) = \sum_{i=1}^N \log p(\vec{x}_i | \pi, \mu, \Sigma) = \sum_{i=1}^N \log \sum_{y=1}^K p(\vec{x}_i, y | \pi, \mu, \Sigma)$

K-Means Clustering

- K-means solves a Chicken-and-Egg problem:
If knew classes, we can get model (max likelihood!)
If knew the model, we can predict the classes (classifier!)
- Kmeans: guess a model, use it to classify the data, use classified data as labeled data to update the model, repeat.
- Assumes each point x has a discrete multinomial vector z



- 0) Input dataset $\{\vec{x}_1, \dots, \vec{x}_N\}$
 - 1) Randomly initialize means $\vec{\mu}_1, \dots, \vec{\mu}_K$
 - 2) Find closest mean for each point $\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases}$
 - 3) Update means $\vec{\mu}_i = \sum_{n=1}^N \vec{x}_n \vec{z}_n(i) / \sum_{n=1}^N \vec{z}_n(i)$
 - 4) If any z has changed go to 2
- $k = \text{the number of classes}$
- $\vec{\mu}_i = \sum_{n \in \text{cluster } i} \vec{x}_n / \text{number of elements in cluster}$
- $\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases}$
- $\vec{z}_n^{(1)} = \begin{cases} 1 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(2)} = \begin{cases} 0 & \\ 1 & \end{cases}$
- $\vec{z}_n^{(10)} = \begin{cases} 1 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(11)} = \begin{cases} 0 & \\ 1 & \end{cases}$
- $\vec{z}_n^{(12)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(13)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(14)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(15)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(16)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(17)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(18)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(19)} = \begin{cases} 0 & \\ 0 & \end{cases}$
- $\vec{z}_n^{(20)} = \begin{cases} 0 & \\ 0 & \end{cases}$

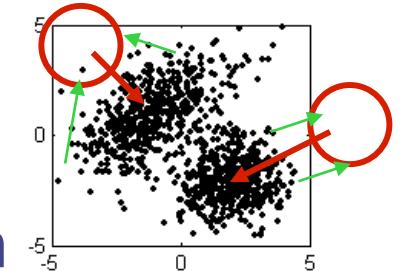
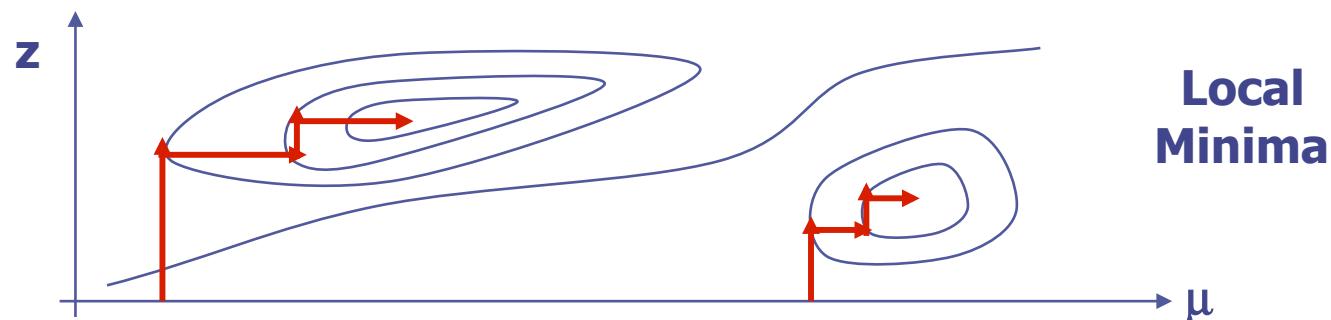
K-Means Clustering

- Geometric, each point goes to closest Gaussian
- Recompute the means by their assigned points
- Essentially minimizing the following cost function:

$$\min_{\mu} \min_z J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N) = \sum_{n=1}^N \sum_{i=1}^K \vec{z}_n(i) \|\vec{x}_n - \vec{\mu}_i\|^2$$

$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases} \quad \vec{\mu}_i = \frac{\sum_{n=1}^N \vec{x}_n \vec{z}_n(i)}{\sum_{n=1}^N \vec{z}_n(i)}$$

- Guaranteed to improve per iteration and converge
- Like Coordinate Descent (lock one var, maximize the other)
- A.k.a. Axis-Parallel Optimization or Alternating Minimization



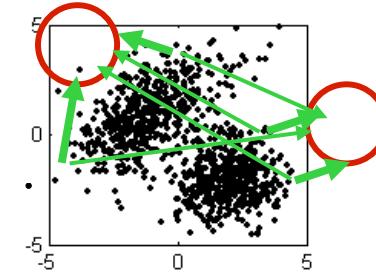
Expectation-Maximization (EM)

- EM is a soft/fuzzy version of K-Means (which does winner-takes-all, closest Gaussian Mean completely wins datapoint)

$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \left\| \vec{x}_n - \vec{\mu}_j \right\|^2 = \arg \max_j N(\vec{x}_n | \vec{\mu}_j, I) = \arg \max_j p(\vec{x}_n | \vec{\mu}_j) \\ 0 & \text{otherwise} \end{cases}$$

- Instead, consider soft percentage assignment of datapoint

$$\text{assign} \propto \pi_j \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} \left\| \vec{x}_n - \vec{\mu}_j \right\|^2\right) \theta = (\vec{\Sigma}, \vec{\mu})$$



$$\tau_{n,i} = \text{class } i \text{ responsibility for Data } n = p(\vec{z} = \vec{\delta}_i | \vec{x}_n, \theta)$$

$$\theta = (\vec{\Sigma}, \vec{\mu})$$

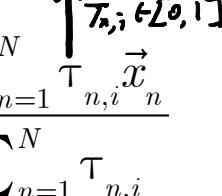
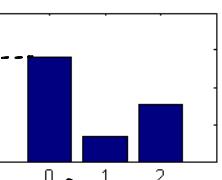
- EM is 'less greedy' than K-Means uses $\tau_{n,i} = p(\vec{z} = \vec{\delta}_i | \vec{x}_n, \theta)$ as shared responsibility for \vec{x}_n

$$\tau_{n,1}, \dots, \tau_{n,K} = T_{n,1} = p(\vec{z} = \vec{\delta}_1 | \vec{x}_n, \theta)$$

$$\sum_i \tau_{n,i} = 1$$

- Update for the means are then 'weighted' by responsibilities:

$$\mu_i = \frac{\sum_{n=1}^N \tau_{n,i} \vec{x}_n}{\sum_{n=1}^N \tau_{n,i}}$$



$$\begin{aligned}
 \tau_{n,i} &= p(z_n = \delta_i | x_n, \theta) = \frac{p(z_n = \delta_i, x_n | \theta)}{p(x_n | \theta)} = \frac{p(x_n | z_n = \delta_i, \theta) \cdot p(z_n = \delta_i | \theta)}{p(x_n | \theta)} = \frac{\pi(\mathbf{x}_n | u_i, \Sigma_i) \cdot \pi_i}{\sum_j \pi(\mathbf{x}_n | u_j, \Sigma_j) \cdot \pi_j} \\
 p(x_n | z_n = \delta_i, \theta) &= N(x_n | u_i, \Sigma_i) \\
 p(x_n | \theta) &= \sum_j p(x_n | z_n = \delta_j, \theta) = \sum_j p(x_n | \theta, z_n = \delta_j) \cdot p(z_n = \delta_j | \theta) \\
 &= \sum_j \pi(\mathbf{x}_n | u_j, \Sigma_j) \cdot \pi_j
 \end{aligned}$$

Tony Jebara, Columbia University

$$\begin{aligned}
 p(x_n | \theta) &= \sum_j p(x_n, z_n = \delta_j | \theta) \\
 &= \sum_j p(x_n | \theta, z_n = \delta_j) \cdot p(z_n = \delta_j | \theta)
 \end{aligned}$$

Expectation-Maximization

- EM uses expected value of $\vec{z}_n(i)$ rather than max

$$\tau_{n,i} = E\left\{\vec{z}_n(i) | \vec{x}_n\right\} = p\left(\vec{z}_n = \vec{\delta}_i | \vec{x}_n, \theta\right)$$

- EM updates covariances, mixing proportions AND means...
- The algorithm for Gaussian mixtures:

→ EXPECTATION:

$$\begin{aligned}
 \tau_{n,i}^{(t)} &= \frac{\pi_i N\left(\vec{x}_n | \vec{\mu}_i^{(t)}, \Sigma_i^{(t)}\right)}{\sum_j \pi_j N\left(\vec{x}_n | \vec{\mu}_j^{(t)}, \Sigma_j^{(t)}\right)} \\
 \vec{\mu}_i^{(t+1)} &= \frac{\sum_n \tau_{n,i}^{(t)} \vec{x}_n}{\sum_n \tau_{n,i}^{(t)}} \quad \pi_i^{(t+1)} = \frac{\sum_n \tau_{n,i}^{(t)}}{N} \\
 \Sigma_i^{(t+1)} &= \frac{\sum_n \tau_{n,i}^{(t)} \left(\vec{x}_n - \vec{\mu}_i^{(t+1)} \right) \left(\vec{x}_n - \vec{\mu}_i^{(t+1)} \right)^T}{\sum_n \tau_{n,i}^{(t)}}
 \end{aligned}$$

迭代模型
Iterative Model

-
- DEMO... like an iterative divide-and-conquer algorithm
 - But, divide&conquer is not a guarantee. Can we prove EM?

https://blog.csdn.net/qq_33369979/article/details/103913637

<https://blog.csdn.net/u010834867/article/details/90762296>