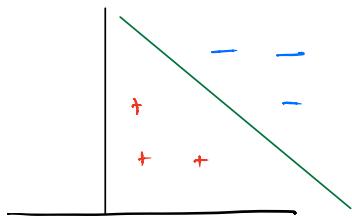


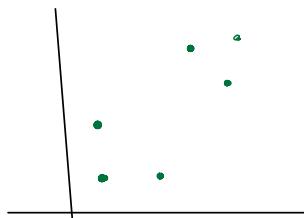
UNSUPERVISED LEARNING

TODAY: K-means, mixture of Gaussians, EM



Supervised Setting

Unsupervised is therefore
than Supervised



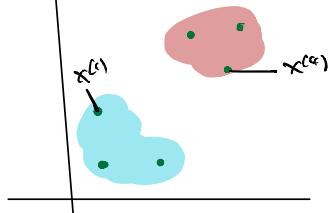
Unsupervised, no labels!

allow Stronger Assumptions
accept Weaker Guarantees

TECHNIQUES \neq IDEAS ARE VALUABLE

K-MEANS

GIVEN $k=2$



Do:



GIVEN $x^{(1)} \dots x^{(n)} \in \mathbb{R}^d$ \notin Integer k , # of clusters

DO find assignment of $x^{(i)}$ to ONE of k clusters

$C^{(i)} = j$ Point i in cluster j

e.g. $C^{(2)} = 2$ while $C^{(1)} = 1$

How do we find these clusters? Iterative Approach



1. Randomly init $\mu^{(1)}, \mu^{(2)}$

for each $i = 1 \dots n$

2. Assign each point to closest cluster $\leftarrow C^{(i)} = \underset{j=1 \dots K}{\operatorname{Argmin}} \| \mu^{(j)} - x^{(i)} \|^2$

3. Compute New cluster centers

REPEAT until no points change

for $j = 1 \dots K$

$$\mu^{(j)} = \frac{1}{|\Omega_j|} \sum_{i \in \Omega_j} x^{(i)} \text{ s.t. } \Omega_j = \{i : C^{(i)} = j\}$$

Comments

Does K-means terminate? Yes!

$$J(c, u) = \sum_{i=1}^n \|x^{(i)} - c^{(i)}\|^2 \text{ decreases monotonically}$$

(SEE NOTES)

Does it find a Global minimum? Not necessarily... NP-HARD

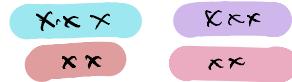
SIDE NOTE: K-means++ from GREAT Stanford Students

- + Improved Apx Ratio through Clever Init
- + DEFAULT IN SKLEARN

How do you choose k? No ONE right answer.



2 clusters



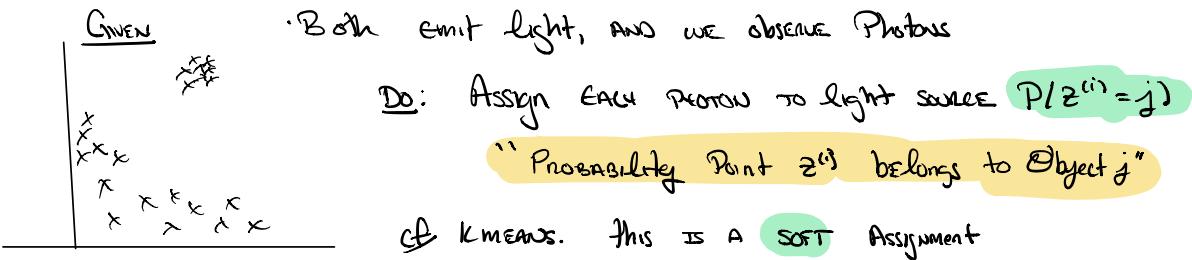
4 clusters

Modeling Question!

Mixture of Gaussians

Toy Astronomy Example (based on a paper from UW)

- QUASARS \neq STARS are sources of light



Challenges + Many Sources (say we know K , # of sources)

+ Sources have different intensity & modes

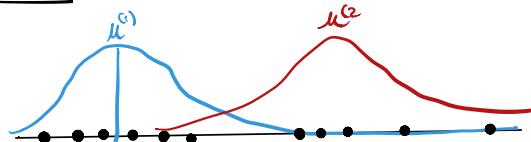
Assume

1. Sources are well modeled by GAUSSIAN (μ_j, σ^2)
2. WE DO NOT assume equal # of points per source
→ UNKNOWN MIXTURE

NB: Physics folks can check if recovered values make sense.

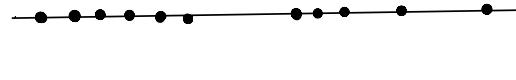
Mixture of Gaussians (MODEL & SETUP) - 1d for simplicity

MODEL:



WE OBSERVE POINTS w/o labels:

$$x^{(i)} \in \mathbb{R}$$



OBSERVATION 1 if we knew "Cluster labels" \rightarrow Solve w/ GFA.



Compute $\mu^{(1)}, \mu^{(2)}$ and be done.

CHALLENGE WE don't

Given $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}$ AND positive integer k

Do find P s.t. for $i=1 \dots n$ & $j=1 \dots k$ clusters

$$P(z^{(i)} = j) \text{ soft assignment}$$

According to the "Gmm model"

$$P(x^{(i)}, z^{(i)}) = P(x^{(i)} | z^{(i)}) P(z^{(i)}) \quad \text{Bayes Rule}$$

$$z^{(i)} \sim \text{Multinomial}(\vec{\phi}) \quad \vec{\phi} \geq 0 \quad \sum_{j=1}^k \phi_j = 1 \quad \text{"which source"}$$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \sigma_j^2) \quad \text{gaussian in each source}$$

The parameters to be found are highlighted

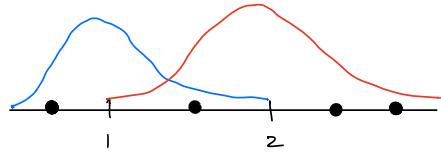
We call $z^{(i)}$ a hidden or latent variable. $z^{(i)}$ is not directly observed

helpful to think in terms

of Sample.

$$\phi_1 = 0.7 \quad \phi_2 = 0.3$$

$$\mu_1 = 1 \quad \mu_2 = 2 \quad \sigma_1^2 = \sigma_2^2 = \frac{1}{3} \text{ (equal)}$$



Gmm Algorithm (Famous Algo \neq Class)

Mixture K-means

1. (E-STEP) "Guess latent values" of $z^{(i)}$ FOR EACH POINT
2. (M-STEP) UPDATE PARAMETERS

ABSTRACTLY OUR FIRST EXAMPLE OF EM-ALGORITHM (Expectation Maximization)

(E-STEP) GIVEN Data \neq current guess at parameters $(\phi, \mu, \sigma^2, \dots)$
DO Predict latent variable $z^{(i)}$ for $i=1\dots n$

$$w_j^{(i)} = P(z^{(i)} = j | x^{(i)}, \phi, \mu, \sigma) \quad \text{our goal}$$

$$= \frac{P(z^{(i)} = j, x^{(i)}; \phi, \mu, \sigma)}{P(x^{(i)}; \phi, \mu, \sigma)} \quad \text{Bayes Rule}$$

$$= \frac{\underset{j}{\underbrace{P(x^{(i)} | z^{(i)} = j; \phi, \mu)}} \underset{\text{②}}{\underbrace{P(z^{(i)} = j; \phi, \mu)}}}{\sum_{l=1}^L \underset{\text{②}}{\underbrace{P(x^{(i)} | z^{(i)} = l; \phi, \mu, \sigma)}} \underset{\phi_l}{\underbrace{P(z^{(i)} = l; \phi)}}}$$

* $\propto \exp \left\{ -\frac{(x^{(i)} - \mu_j)^2}{\sigma_j^2} \right\}$ "How likely is $x^{(i)}$ according to Gaussian (μ_j, σ_j^2) "

● "How likely point from cluster"

Key Point WE CAN COMPUTE ALL TERMS! RETURN $w_j^{(i)}$

M-STEP

GIVEN $w_j^{(i)}$ our current estimate of $P(Z^{(i)} = j)$ for $i = 1, \dots, n$
 $j = 1, \dots, k$ clusters

DO Estimate Observed Parameters (using MLE)

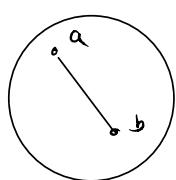
e.g. $\phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)} \approx$ fraction of elements in cluster j

$$m_j = \frac{\sum w_j^{(i)} x^{(i)}}{\sum w_j^{(i)}} \quad \dots \text{etc} \dots$$

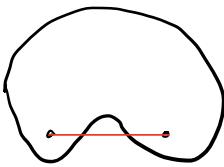
MLE. Let's make rigorous:

Detour Convexity \nRightarrow JENSEN (This is a key result, we'll go slowly)

A SET Ω IS CONVEX if for any $a, b \in \Omega$ the line joining a, b is in Ω as well.



Convex



NOT convex!

IN symbols,

$$\forall \lambda \in [0,1], a, b \in \Omega$$

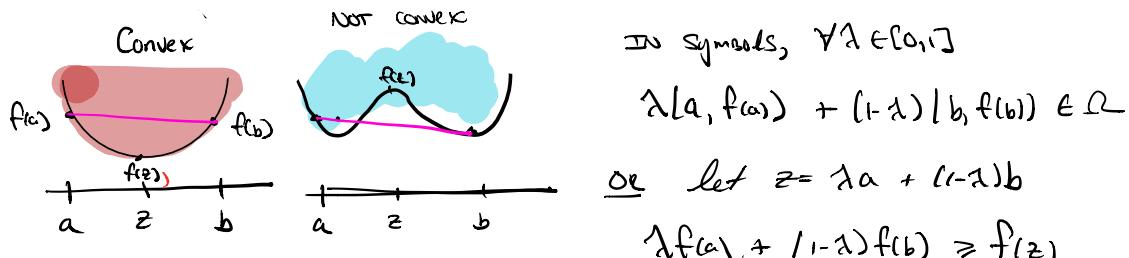
$$\lambda a + (1-\lambda)b \in \Omega$$

(NEED TO SHOW $a, b \in \Omega$)

GIVEN a function f , the graph of f G_f is defined as

$$G_f = \{ (x, y) : y \geq f(x) \}$$

A function is convex if its graph is convex (as a set)



"Every curve is above function"

If f twice differentiable, $\forall z$ $f'(z) \geq 0 \Rightarrow f$ is convex

$$\text{def } f(a) = f(z) + f'(z)(a-z) + f''(z_a)(a-z)^2 \quad \text{for } a \in [a, z]$$

$$f(b) = f(z) + f'(z)(b-z) + f''(z_b)(b-z)^2 \quad \text{for } z \in [z, b]$$

$$\lambda f(a) + (1-\lambda)f(b) = f(z) + f'(z)(\cancel{\lambda a + (1-\lambda)b} - z) + c \quad c \geq 0$$

$$\text{i.e. } \lambda f(a) + (1-\lambda)f(b) \leq f(z) \quad \square \quad \text{def of } z.$$

We say f is strongly convex if $\forall x \in \text{Dom}(f) \quad f''(x) > 0$.

Ex: $f(x) = x^2 \Rightarrow f''(x) = 2 \Rightarrow$ strongly convex

$f(x) = x^2(x-1)^2$: graph above (lower convex)

JENSEN'S INEQUALITY $\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$ for convex f .

Ex: x takes value a with prob λ

· takes value b with prob $1-\lambda$

$$\mathbb{E}[f(x)] = \lambda f(a) + (1-\lambda)f(b)$$

$$f(\mathbb{E}[x]) = f(z) \quad z = \lambda a + (1-\lambda)b$$

NS: can prove finitely
Supported Distribution
By induction

for convex f , definition implies this in this case!

stronger if f is strongly convex, and $\mathbb{E}[f(x)] = f(\mathbb{E}[x])$

$\Rightarrow x$ is a constant (elements: almost surely)

WE NEED CONCAVE FUNCTIONS g concave iff $-g$ is convex

$$\mathbb{E}[g(x)] \leq g(\mathbb{E}[x])$$

Ex: $g(x) = \log(x) \Rightarrow g''(x) = -x^{-2}$ on $(0, \infty)$ NEGATIVE



WHAT ABOUT $f(x) = ax + b$ CONVEX & CONCAVE since $f''(x) = 0$.

END DETAIL

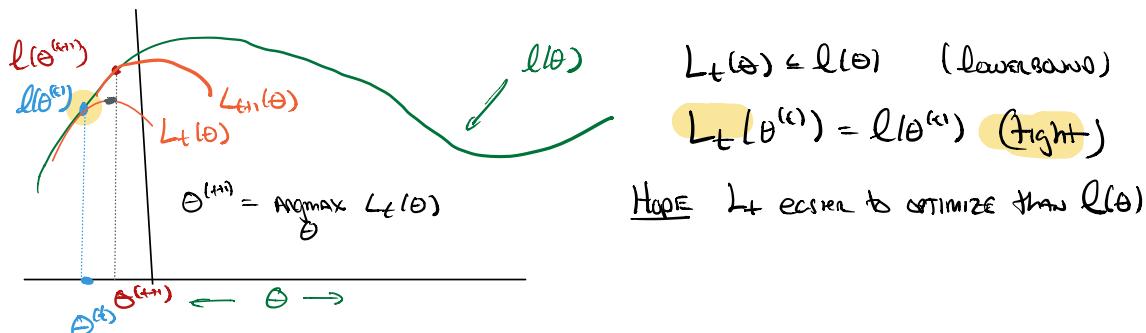
EM Algorithm as max likelihood

$$\ell(\theta) = \sum_{i=1}^n \log P(x^{(i)}; \theta)$$

↑ DATA ↑ PARAMETERS

WE ASSUME $P(x; \theta) = \sum_z P(x, z; \theta)$ of GMM LATENT VARIABLE

Picture of Algorithm



Rough Algo

(E-STEP) 1. GIVEN $\theta^{(t)}$ FIND L_t

(M-STEP) 2. GIVEN L_t , SET $\theta^{(t+1)} = \operatorname{argmax}_{\theta} L_t(\theta)$

How do we construct L_t ? (Let's look at single point)

$$\log \sum_z P(x, z; \theta) = \log \sum_z \frac{Q(z) P(x, z; \theta)}{Q(z)} \quad \text{for any } Q(z)$$

WE PICK $Q(z)$ S.T. $\sum_z Q(z) = 1$ AND $Q(z) = 0 \Leftrightarrow$

$$= \log \mathbb{E}_z \left[\frac{P(x, z; \theta)}{Q(z)} \right] \quad (\text{Symbol Rushing})$$

$$\geq \mathbb{E}_z \left[\log \frac{P(x, z; \theta)}{Q(z)} \right] \quad \text{JENSEN?} \quad (\log \text{ is concave})$$

$$= \sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)} \quad (\text{DEF of } \mathbb{E})$$

Key step holds for any such Q : (a)

This gives a family of lower bounds, one for each choice of Q ($P_\theta \leq l$)

How do we make it tight? Select Q to make inequality tight

What if... $\log \frac{P(x, z; \theta)}{Q(z)} = c$ for some constant, then JENSEN's is Equality!

$$P(x, z; \theta) = P(z|x; \theta) P(x; \theta)$$

so, $Q(z) = P(z|x; \theta)$ then

$c = \log P(x; \theta)$ does not depend on z , so constant!

NB: $Q(z)$ does depend on $\theta + x$ - we will select a $Q^{(i)}(z)$
for every point independently.

WE DEFINE Evidence-based Lower Bound (ELBO), sum over z

$$\text{ELBO}(x, Q, z) = \sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)}$$

WE'VE SHOWN $l(\theta) \geq \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta)$ for any $Q^{(i)}$ satisfying (a)
lower bound

$$l(\theta^{(i)}) = \sum_{z \in \mathcal{Z}} \text{ELBO}(x^{(i)}, Q^{(i)}, \theta^{(i)}) \quad \text{for choice of } Q^{(i)} \text{ above.}$$

WRAP UP

1. (E-STEP) $Q^{(t)}(z) = P(z^{(t)})x^{(t)}; \theta)$ for $i=1 \dots n$

2. (M-STEP) $\theta^{(t+1)} = \underset{\theta}{\operatorname{Argmax}} L_t(\theta)$

in which $L_t(\theta) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta)$

WHY DOES THIS TERMINATE? $L(\theta^{(t+1)}) \geq L(\theta^{(t)})$

IS IT GLOBALLY OPTIMAL? (MORE SEE PICTURE)

WE DERIVED HARD & SOFT CLUSTERING METHODS

EM ALGORITHM IN TERMS OF MLE.