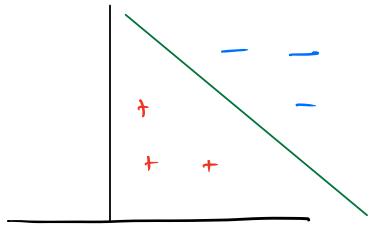
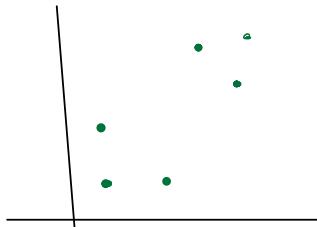


UNSUPERVISED LEARNING

TODAY: K-means, mixture of Gaussians, EM



Supervised Setting



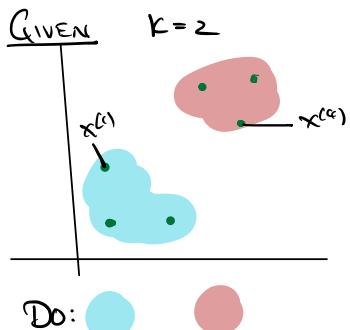
Unsupervised, no labels!

Unsupervised is harder
than Supervised

allow Stronger Assumptions
accept Weaker Guarantees

TECHNIQUES & IDEAS ARE VALUABLE

K-MEANS



GIVEN $x^{(1)} \dots x^{(n)} \in \mathbb{R}^d$ \nsubseteq 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 $c^{(4)} = 1$

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 $\longleftrightarrow C^{(i)} = \underset{j=1 \dots k}{\operatorname{Argmin}} \| \mu^{(j)} - x^{(i)} \|^2$

3. Compute New cluster centers

for $j = 1 \dots k$

REPEAT until no points change

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

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

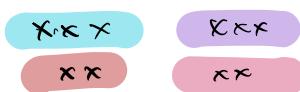
SIDE NOTE: K-means++ from great Stanford Students

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

How do you choose k? No one right answer.



2 clusters

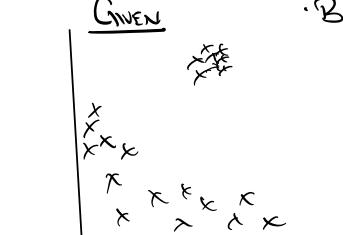


4 clusters

Modeling Question!

Mixture of Gaussians

Toy Astronomy Example (based on a paper from UW)

- Quasars & stars are sources of light
 - Both emit light, and we observe photons
- GIVEN 
- Do: Assign each photon to light source $P(z^{(i)} = j)$
- "Probability Point $z^{(i)}$ belongs to Object j "
- ↳ kmeans. This is a **soft** assignment

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

+ Sources have different intensity & shapes

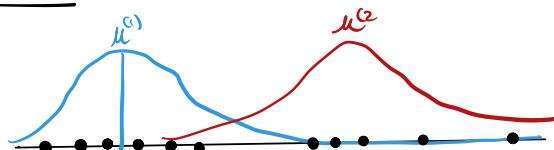
Assume

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.

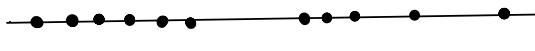
Mixture of Gaussians (MODEL & SETUP) - 1d for simplicity

model:



WE OBSERVE POINTS w/o source!

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



OBSERVATION 1 if we knew "Cluster labels" \rightarrow solve w/ GA.



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 $\hat{\gamma}$ s.t. for $i=1 \dots n$ & $j=1 \dots k$ clusters

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

According to the "Gmm model"

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

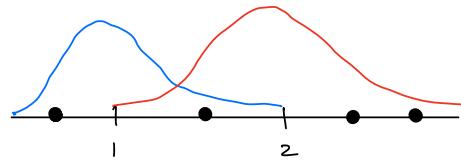
$$z^{(i)} \sim \text{Multinomial}(\vec{\phi}) \quad \vec{\phi} \geq 0 \quad \sum_{j=1}^k \phi_j = 1 \quad \text{"which source"}$$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \sigma_j^2) \quad \text{GAUSSIAN IN EACH SOURCE}$$

The parameters to be found are highlighted

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

Example "think Sampling"



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

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

1. Pick Blue (w/ Prob 0.7) or Red (w prob 0.3)

2. Use appropriate mean μ_1 (red) or μ_2 (blue)

Repeats

Gmm Algorithm (Famous Algo \notin Class)

Minimizing χ^2 -meaus

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 $i=1..n$

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

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

$$= \frac{P(x^{(i)} | z^{(i)} = j; \phi, \mu, \sigma) P(z^{(i)} = j; \phi, \mu, \sigma)}{\sum_{l=1}^L P(x^{(i)} | z^{(i)} = l; \phi, \mu, \sigma) P(z^{(i)} = l; \phi, \mu, \sigma)} \quad \begin{matrix} \phi_j \\ \phi_l \end{matrix}$$

* $\propto \exp \left\{ -\frac{(x^{(i)} - \mu_j)^2}{\sigma_j^2} \right\}$ "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 $i = 1, \dots$
 $j = 1, \dots, k$ clusters

DO Estimate Observed Parameters (using MLE)

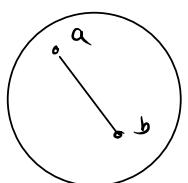
e.g. $\phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)} \approx$ fraction of elements in cluster j

$$\mu_j = \frac{\sum w_j^{(i)} x^{(i)}}{\sum w_j^{(i)}} \quad \text{-- etc...}$$

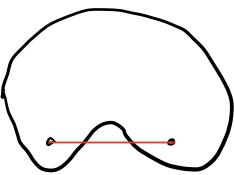
MLE. Let's make rigorous!

Defn Convex $\not\equiv$ 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,

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

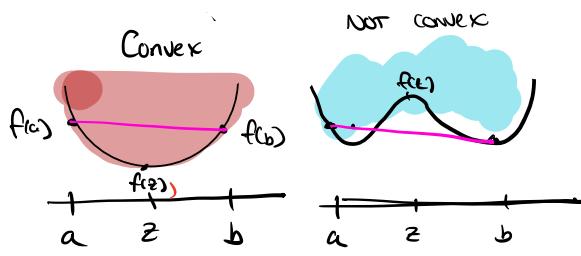
$$\lambda a + (1-\lambda)b \in \Omega$$

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

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)



In symbols, $\forall \lambda \in [0,1]$

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

or let $z = \lambda a + (1-\lambda)b$

$$\lambda f(a) + (1-\lambda)f(b) \geq f(z)$$

"Every curve is above function"

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

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

$$f(b) = f(z) + f'(z)(b-z) + f''(z_b)(b-z)^2 \quad \exists z \in [z, b]$$

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

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

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

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

$f(x) = x^2(x-1)^2$: graph below (not 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 $1-\lambda$

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

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

NB: can prove finitely supported distribution by induction

for convex f , definition implies this in this case!

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

$\Rightarrow x$ is a constant (except: 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 DETAILED

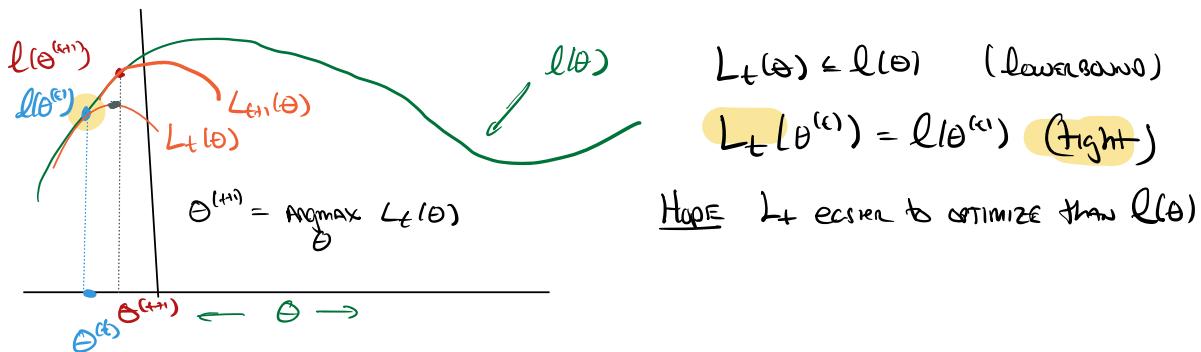
EM Algorithm as max likelihood

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

PARAMETERS
DATA

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 $Q(z) = 0 \iff$

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

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

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

Key step holds for any such Q : (x)

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

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 (x)
lower bound

$$l(\theta^{(i)}) = \sum_{z=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}, \theta^{(i)}) \quad \text{for choice of } Q^{(i)} \text{ above.}$$

WAP UP

1. (E-STEP) $Q^{(i)}(z) = P(z^{(i)})x^{(i)}; \theta)$ 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? $\ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)})$

IS IT GLOBALLY OPTIMAL? (NOPE SEE PICTURE)

WE DERIVED HARD & SOFT CLUSTERING METHODS

EM ALGORITHM IN TERMS OF MLE.