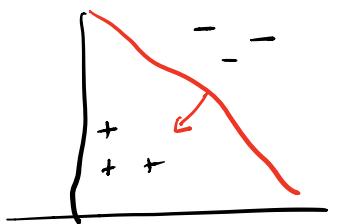
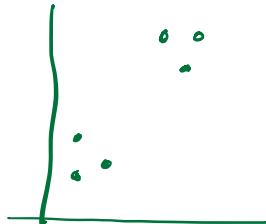


## Unsupervised Learning

TODAY: K-MEANS, MIXTURE OF GAUSSIANS, EM.



Supervised - Points And Labels



UNSUPERVISED - NO LABELS!

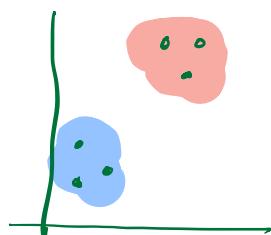
TECHNIQUES ARE VALUABLE (PEDAGOGICALLY & PRACTICALLY)

### K-MEANS

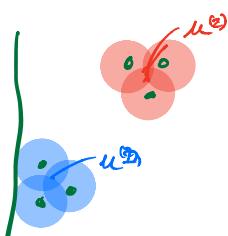
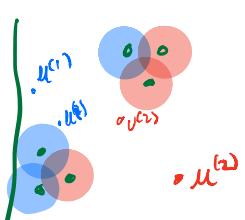
GIVEN  $\{x^{(1)}, \dots, x^{(n)}\} \in \mathbb{R}^d \neq k$  # of clusters

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

$c^{(i)} = j$  "Point  $i$  is cluster  $j$ "



How do we find these clusters?



1. Randomly Init  $\mu^{(1)}, \mu^{(2)}$
2. Assign EACH point to Cluster  $\longleftrightarrow C^{(i)} = \underset{j=1 \dots k}{\operatorname{Argmin}} \| \mu^{(j)} - x^{(i)} \|^2$
3. Compute NEW clusters CENTERS  
REPEAT UNTIL NO points CHANGE  $\mu^{(j)} = \frac{1}{|\Omega_j|} \sum_{i \in \Omega_j} x^{(i)}$   
where  $\Omega_j = \{i : C^{(i)} = j\}$  "Compute mean"

### Comments

+ Does K-MEANS TERMINATE? Yes!

$$J(C, \mu) = \sum_{i=1}^n \|x^{(i)} - \mu_{C(i)}\|^2 \text{ is decreasing monotonically}$$

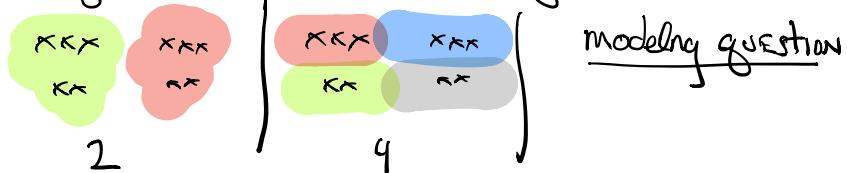
(Is nones)

+ Does it find a minimum? NOT NECESSARILY (NP-HARD)

SIDE NOTE: K-MEANS++ 2007 from GREAT Stanford Students

- IMPROVED APPROXIMATION BOUNDS through Clever (not
- DEFAULT IN SKLEARN

+ How do you choose k? NO ONE Right Answer

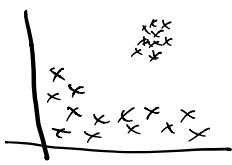


### Mixture of Gaussians

Toy Astronomy Example (BASED on real one example)

- QUASARS & STARS ARE SOURCES OF LIGHT

BOTH EMIT LIGHT, AND WE OBSERVE SHOTNOISE



Goal: Assign each photon to light source  $P(z^{(i)} = j)$   
 "Probability Point  $z^{(i)}$  belongs to Object  $j$ "

of K-MEANS. This is a **soft** Assignment

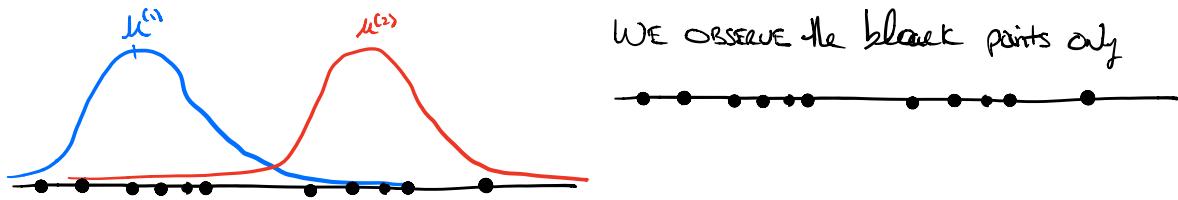
### Challenges

- Many Sources (say we know # of sources)
- Sources have different intensity  $\neq$  shape.

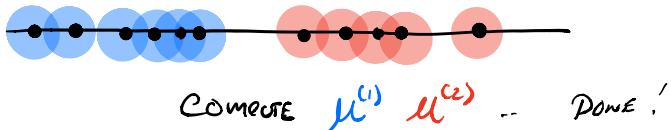
Assume: 1. Sources are Gaussian like ( $\mu_j, \sigma_j^2$ )  
 2. We do not assume equal # of points per source (mixture)

NS: In this example, physics can check how plausible it is.

### Mixture of Gaussians: Model & Setup



OBSERVATION 1: If we KNEW "cluster lasers"  $\rightarrow$  solve with QDA



Challenge: we don't!

GIVEN:  $x^{(1)} \dots x^{(n)} \in \mathbb{R}$  AND  $K > 0$

DO: FIND Prob  $z^{(i)}$  FOR  $i=1 \dots n$  TO ONE OF  $K$  CLUSTERS

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

Gmm model

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

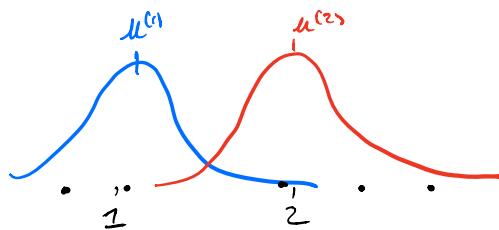
$$z^{(i)} \sim \text{multinomial}(\phi) \text{ s.t. } \sum_{j=1}^K \phi_j = 1, \phi_j \geq 0$$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \sigma_j^2)$$

THE PARAMETERS TO BE FOUND ARE IN BLUE

WE CALL  $z^{(i)}$  A HIDDEN OR LATENT VARIABLE

$z^{(i)}$  IS NOT DIRECTLY OBSERVED.



GENERATE DATA:

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

$$\mu^1 = 1 \quad \mu^2 = 2 \\ \sigma_1^2 = \sigma_2^2 = 1$$

Gmm Algorithm (FAMOUS ALGORITHM  $\neq$  type)

MIRRORS K-MEANS

1. (E-STEP) "GUESS" LATENT VALUES OF  $z^{(i)}$  (FOR EACH POINT)
2. (M-STEP) UPDATE PARAMETERS (MLE)

Why Abstract? VERY GENERAL EM ALGORITHM

E-STEP 1:

GIVEN: DATA  $\neq$  CURRENT PARAMETERS

Gmm

$$x \nmid (\phi, \mu, \sigma)$$

Do: Predict latent value  $z^{(i)}$  for  $i = 1 \dots n$

$$w_j^{(i)} = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \sigma) \quad \text{("How likely is data")}$$

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

$$= \frac{P(x^{(i)} | z^{(i)} = j; \phi, \mu) P(z^{(i)} = j | \phi)}{\sum_l P(x^{(i)} | z^{(i)} = l; \phi, \mu) P(z^{(i)} = l | \phi)}$$

○  $\propto \exp \left\{ -\frac{(x^{(i)} - \mu)^2}{\sigma^2} \right\}$  "How likely from this Gaussian"

○  $= \phi_j$  or  $\phi_0$  "How likely this point came from Cluster"

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

### M-STEP

GIVEN  $w_j^{(i)} = P(z^{(i)} = j)$  ESTIMATE for all  $i \in \{n\}$  of clusters

Do: Estimate the other observed PARAMETERS USY MLE

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

$$\mu_j = \frac{\frac{1}{n} \sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}} \approx$$
 "soft cluster center"

$$\begin{matrix} \vdots \\ \sigma_j^2 \end{matrix}$$

more generally. Let's make rigorous!

DETTOUR: CONVEXITY  $\nRightarrow$  JENSEN.

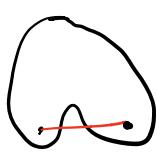
THIS IS A KEY RESULT, SO WANT TO GO SLOWLY! ASK!

A SET  $\Omega$  IS CONVEX if for any  $a, b \in \Omega$

the line joining  $a, b$  is in  $\Omega$



CONVEX



NOT CONVEX

IN SYMBOLS,

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

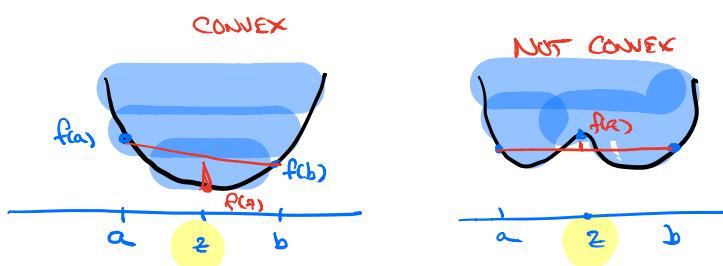
$$\lambda a + (1-\lambda)b \in \Omega$$

(NEED TO CHECK  $f(a), f(b) \in \Omega$ )

GIVEN A FUNCTION  $f$  the graph of  $f$  is  $G_f$  defined

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

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



CONVEX MEANS

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

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

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

"Every chord is above function"

If  $f$  is twice differentiable,  $f''(x) > 0 \Rightarrow f$  is convex

$$\text{Pf: } f(a) = f(z) + f'(z)(a-z) + f''(z)(a-z)^2 \quad z \in (a, b)$$

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

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

--- ≤ --- □

WE say  $f$  is strictly convex if  $\forall x \in \text{dom}(f)$   $f''(x) > 0$

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

$f(x) = x^2(x-1)^2$  is the graph above that is not convex.

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

Ex:  $X$  takes value  $a$  with prob. 1

takes value  $b$  with prob.  $1-1$

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

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

for convex  $f$ , our definition above implies JENSEN'S thm.

NB: for finitely supported distributions, prove Jensen's by induction

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

then  $X$  is a constant (experiments almost surely)

WE NEED CONCAVE FUNCTIONS  $g$  is concave if  $-g$  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

curve below function.

WHAT ABOUT  $f(x) = ax + b$ ?  $f''(x) = 0 \Rightarrow$  convex AND concave!

$$\begin{aligned} \mathbb{E}[f(x)] &\geq f(\mathbb{E}[x]) \quad \nmid \quad \mathbb{E}[f(x)] \leq f(\mathbb{E}[x]) \\ \Rightarrow \mathbb{E}[f(x)] &= f(\mathbb{E}[x]) \quad \text{linear!} \end{aligned}$$

### END DETOUR

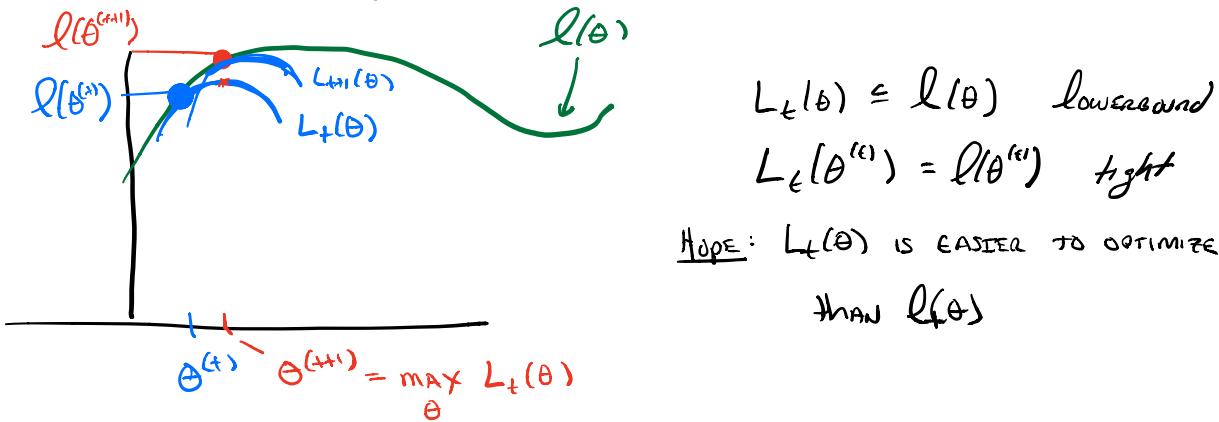
### EM Algorithms as max likelihood

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

PARAMETERS  
DATA

WE ASSUME  $P(x; \theta) = \sum_z P(x, z; \theta)$  cf. GMM  
latent variable

### Picture of Algorithm



### Rough Algorithm

(E-STEP) 1. Find  $f_t(\theta)$  given  $\theta^{(t)}$

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

How do we construct  $L_t(\theta)$

Let's EXAMINE A single point

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

Pick  $Q(z)$  s.t.  $\sum_z Q(z) = 1, Q(z) \geq 0$  ( $\star$ )

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

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

$$= \sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)} \quad (\text{symbol pushing, defl E})$$

① this holds for any  $Q$  satisfying  $(\star)$

② this gives a family of lower bounds (PROPERTY 1 ABOVE)

Pick  $Q(z)$  to satisfy PROPERTY 2. that is

$$\log \sum_z P(x, z; \theta) = \sum_z \log P(x, z; \theta)$$

1.e. when is JENSEN'S tight?

$$\text{if we pick } \frac{P(x, z; \theta)}{Q(z)} = c \quad \text{i.e. } Q(z) = P(z|x; \theta)$$

NB:  $Q(z)$  depends on  $\theta \nmid x$  - different  $Q(z)$  for every point.

WE DEFINE Evidence-Based Lower bound (ELBO), sum over  $z$

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

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

$$\ell(\theta^{(t)}) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta^{(t)}) \text{ for choice of } Q^{(i)} \text{ above}$$

WRAP UP:

1. (E-STEP)  $Q^{(t)}(z) = P(z^{(i)} | x^{(i)}, \theta)$
2. (M-STEP)  $\theta^{(t+1)} = \underset{\theta}{\text{ARGMAX}} \ell_t(\theta)$   
in which  $\ell_t(\theta) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}, \theta)$

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

IS IT GLOBALLY OPTIMAL? NO! (SEE PICTURE)

IN THIS LECTURE, WE SAY HARD & SOFT CLUSTERING METHODS  
WE DERIVED GENERAL ALGORITHM (EM) IN TERMS  
OF MLE.