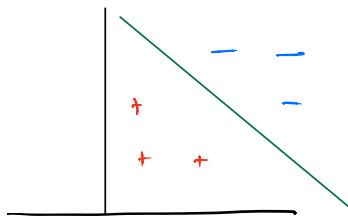
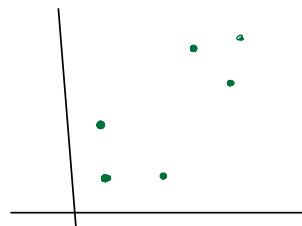


UNSUPERVISED LEARNING

TODAY: K-means, mixture of Gaussians, EM



Supervised Setting



Unsupervised, no labels!

Unsupervised is therefore
than Supervised

allow Stronger Assumptions
accept Weaker Guarantees

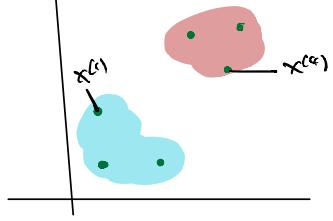
TECHNIQUES & IDEAS ARE VALUABLE
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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 $C^{(4)} = 1$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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 $C^{(i)} = \underset{j=1\dots K}{\operatorname{Argmin}} \| \mu^{(j)} - x^{(i)} \|^2$
3. Compute New cluster centers for $j=1\dots K$
$$\mu^{(j)} = \frac{1}{|\mathcal{D}_j|} \sum_{i \in \mathcal{D}_j} x^{(i)}$$

REPEAT until no points change $\mathcal{D}_j = \{i : C^{(i)} = j\}$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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.

Assignment Project Exam Help



<https://powcoder.com>

2 clusters 4 clusters

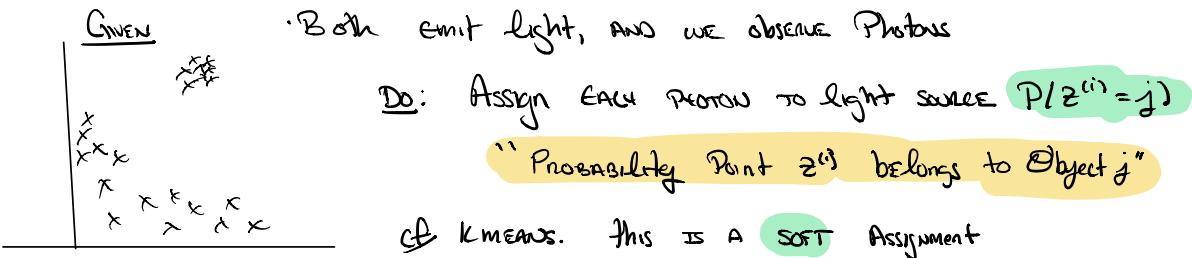
Add WeChat powcoder

Modeling Question¹

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 intensities & modes

Assignment Project Exam Help

Assume 1. Sources are well modeled by Gaussian (μ_j, σ^2)

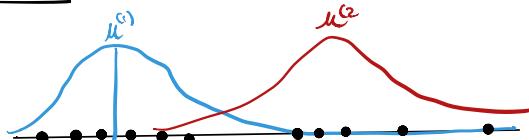
2. WE DO NOT ASSUME equal # of points per source
→ UNKNOWN MIXTURE

NB: Physics folks can care if known values make sense.

Add WeChat powcoder

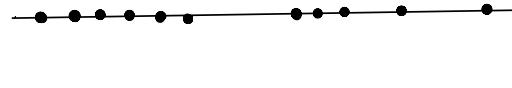
Mixture of Gaussians (MODEL & SETUP) - 1d for simplicity

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

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

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

According to the "Gmm model"

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

$$z^{(i)} \sim \text{Multinomial}(D) \quad \Phi_j \geq 0 \quad \sum_{j=1}^K \Phi_j = 1 \quad \text{"which source"}$$

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

The parameters to be found are highlighted

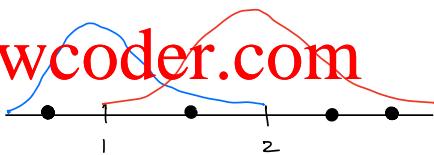
WE CALL $z^{(i)}$ A HIDDEN OR LATENT VARIABLE. $z^{(i)}$ IS NOT DIRECTLY OBSERVED

Assignment Project Exam Help

helpful to think in terms
of Sample.

<https://powcoder.com>

$\Phi_1 = 0.7 \quad \Phi_2 = 0.3$
 $\mu_1 = 1 \quad \mu_2 = 2 \quad \sigma_1^2 = \sigma_2^2 = \frac{1}{3}$ (equal)



Gmm Algorithm (Famous Algo \neq Class)

Mimics 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 $i=1\dots n$

$$w_j^{(i)} = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \sigma^2) \xrightarrow{\text{our goal}}$$

$$= \frac{P(z^{(i)} = j; x^{(i)}; \phi, \mu, \sigma^2)}{P(x^{(i)}; \phi, \mu, \sigma^2)} \quad \text{Bayes Rule}$$

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

* $\propto \exp \left\{ -\frac{(x^{(i)} - \mu_j)^2}{\sigma_j^2} \right\}$ "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

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

MLE. Let's make rigorous:

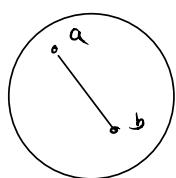
Assignment Project Exam Help

<https://powcoder.com>

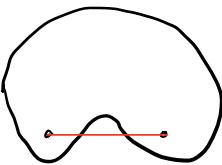
Add WeChat powcoder

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 $\subseteq \Omega$ as well.



Convex



NOT convex!

IN symbols,

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

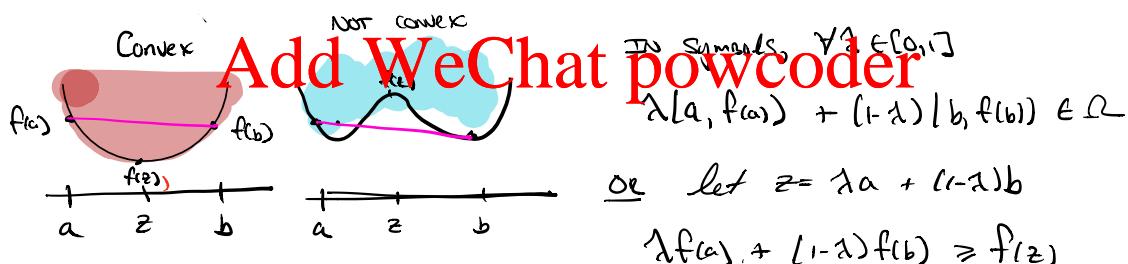
$$\lambda a + (1-\lambda)b \in \Omega$$

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

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)

<https://powcoder.com>



"Every curve is above function"

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

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

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

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

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

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

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

Assignment Project Exam Help

WE NEED CONCAVE FUNCTIONS

Add WeChat powcoder

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

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 DETAILS

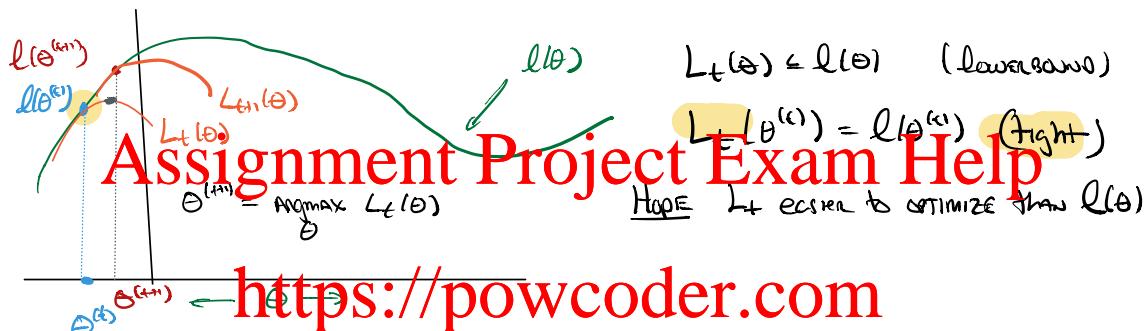
EM Algorithm as max likelihood

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

↑ DATA ↑ PARAMETERS

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

Picture of Algorithm



Rough Algo

Add WeChat powcoder

(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 $Q(1) = 0 \rightsquigarrow$

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

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

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

Key step holds for any such Q : (a)

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

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!

Assignment Project Exam Help

so, $Q(z) = P(z | x; \theta)$ then
 $\ell = \log \frac{P(x; \theta)}{Q(z)}$ does not depend on z , so constant!

NB: $Q(z)$ does depend on $\theta + x$ - we will select a $Q^{(1)}(z)$

Add WeChat powcoder

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

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

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? $L(\theta^{(t+1)}) \geq L(\theta^{(t)})$

IS IT GLOBALLY OPTIMAL? (MORE SEE PICTURE)

WE DERIVED HARD & SOFT CLUSTERING METHODS

EM Algorithm \Rightarrow FEAS of MLE.
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder