*Ch 9 of Bishop*

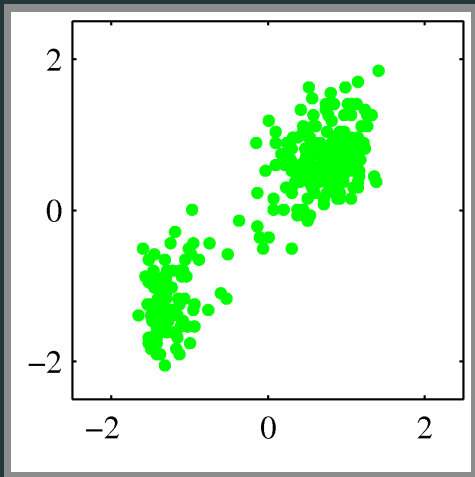# ORIE 4742 - Info Theory and Bayesian ML

Chapter 11: Mixture Models

May 5, 2021

Sid Banerjee, ORIE, Cornell

# Plan for next few classes

- More complex generative models / decision problems

  - latent variable models (mixture model)

  - optimizing complex fns  — simulated annealing
    - random walks
    - gradient descent
    - stochastic gradient descent

  — fitting unknown fns with GPs (Bayesian opt$^n$)

  - neural networks

# example: clustering points in $\mathbb{R}^2$



- Note - No 'training' data (no ground truth)

Idea - 'K-means'

- pick $K$ = number of clusters (Eg. $K=2$)

- pick cluster centers $\mu_1, \mu_2$

- For each point $x_n$, pick 'cluster membership' $r_n = \{r_{n1}, \ldots, r_{nk}\}$

  s.t. $\sum_{i=1}^{k} r_{ni} = 1$, $r_{ni} \in \{0,1\}$

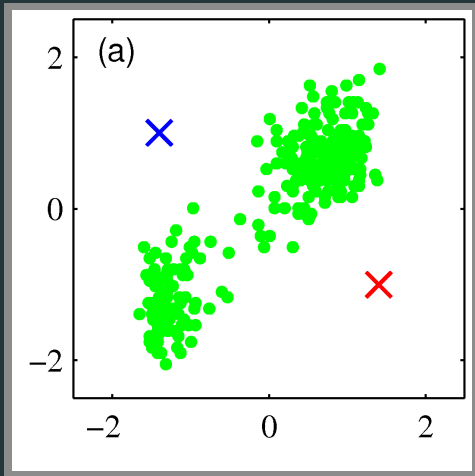- Aim - minimize 'distortion'

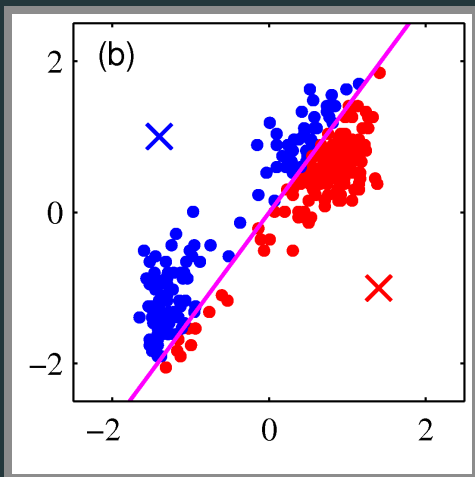$$\sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|x_n - \mu_k\|_2^2$$

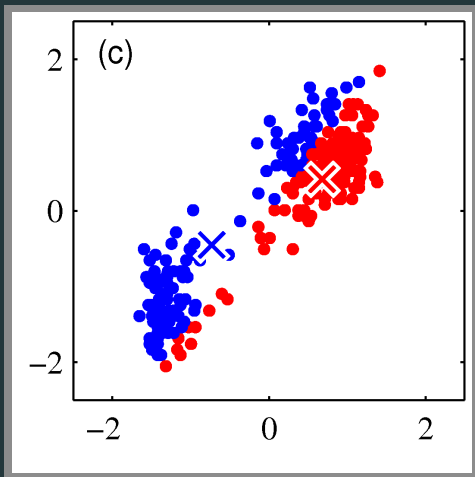(ie 'facility location')

Start by 'guessing' $\mu_1, \mu_2$



(a)

(b)

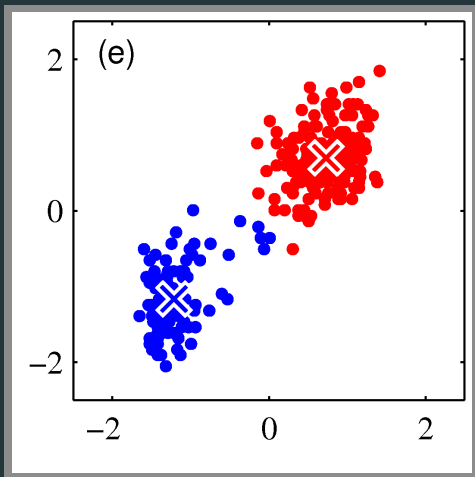∅ **next** **update** Cluster center $\mu_k$ to minimize $\sum \|x_n - \mu_k\|_2^2$ for $x_n$ st. $q_{nk} = 1$

(c)

(e)

(f)

ie - Iteratively set $\{r_{nk}\}$ and $\{\mu_a\}$

- Suppose we knew cluster labels ('supervised learning')
  - Can use 'standard' Bayesian ML classification models (Naive Bayes, Logistic regression, GP classification)

- The clustering problem - no cluster labels (unsupervised learning) - no examples of 'correct' answers

# latent variable generative models

exists, but is not in the data



'Common' latent variable

Eg - regression

## the Gaussian mixture model

- data $D = \{X_1, X_2 \ldots, X_N\} \in \mathbb{R}^d$
- each point $X_n$ has a latent cluster label in $\{1, 2, \ldots, K\}$
  denoted by $Z_n \in \{0,1\}^K$, $\sum_{k=1}^{K} z_{n,i} = 1$     (1-of-$K$ encoding)

  *indicator vectors*

- latent variable: $Z_n \sim \text{Mult}(\pi_1, \pi_2, \ldots, \pi_K)$ where $\sum_{i=1}^{K} \pi_i = 1$
  data: if latent cluster is $k \in [K]$, then $X_n \sim \mathcal{N}(\mu_k, \Sigma_k)$   ie. $Z_n = e_i$

  with prob $\frac{\pi_i}{\Sigma \pi_i}$

- joint likelihood:

$$p\left(X, Z \mid \mu, \Sigma, \pi\right) = \prod_{n=0}^{N-1} \prod_{k=0}^{K-1} \left[\pi_k \mathcal{N}\left(X_n \mid \mu_k, \Sigma_k\right)\right]^{z_{n,k}}$$

$(2\pi)^{-d/2} \Sigma_k^{-1/2} e^{-(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)}$
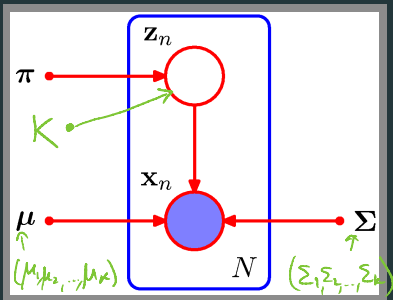
- log-likelihood of data:

$$\log p\left(X \mid \mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n \mid \mu_k, \Sigma_k\right)\right]$$

*not convex!*

## the Gaussian mixture model

log-likelihood of data:

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[ \sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right) \right]$$



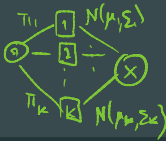Bayes Net for GMM

hyper params

· $K$
· $\Pi = \{\pi_1, \pi_2, \ldots, \pi_K\}$, $\sum_{i=1}^{K} \pi_i = 1$
· For each $k \in [k]$: $\mu_k \in \mathbb{R}^d$

$\Sigma_k = d \times d$, pos def

## the responsibility function

given a Gaussian mixture model with known $\{\mu_k, \Sigma_k\}_{k \in [K]}$, and any data point $X$, we can associate a responsibility parameter to each cluster for the point to be the probability of the underlying latent cluster

### responsibility



$$\gamma(z_k) = \mathbb{P}(z_k = 1 | X) = \frac{\pi_k \mathcal{N}(X | \mu_k, \Sigma_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(X | \mu_j, \Sigma_j)}$$

- Prior $\Pi$, data $X$ $\Rightarrow$ Posterior $\equiv$ $\gamma$

- 'Responsibility' each cluster has for 'explaining' the data $X$

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$

$$\propto \Sigma_k^{-1/2} \exp\left(-\frac{1}{2}\left(x_n - \mu_k\right)^T \Sigma_k^{-1}\left(x_n - \mu_k\right)\right)$$

$$s.t \quad \sum_{i=1}^{k} \pi_i = 1$$

- Set $\dfrac{\partial \log p\left(X|\mu, \Sigma, \pi\right)}{\partial \theta} = 0$ for $\begin{pmatrix} \text{first order} \\ \text{conditions} \end{pmatrix}$

$$\theta \in \left\{\mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K, \pi_1, \pi_2 \dots, \pi_k\right\}$$

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$

- $\underline{\text{FOC}}$ : $\quad -\sum_{n=1}^{N} \underbrace{\frac{\pi_k \mathcal{N}(X|\mu_n, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x|\mu_j, \Sigma_j)}}_{\gamma_{X_n}(z_k)} \underbrace{\Sigma_k}(X_n - \mu_n) = 0$

  assuming invertible, multiply by $\Sigma^{-1}$

$\Rightarrow \quad \mu_k^* = \underbrace{\frac{\sum_{n=1}^{N} \gamma_{X_n}(z_k) X_n}{\sum_{n=1}^{N} \gamma(z_k)}}_{} = \frac{1}{N_k} \underbrace{\sum_{n=1}^{N} \gamma_{X_n}(z_k) X_n}_{\text{weighted sum of } X_n}$

'effective' # of pts in cluster $k$

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$

Similarly (after some algebra)     $\left(N_k = \sum_{n=1}^{N} \gamma_{X_n}(z_k)\right)$

$$\sum_{k}^{*} = \frac{1}{N_k} \sum_{n=1}^{N} \gamma_{X_n}(z_k^{*}) \left(X_n - \mu_k^{*}\right)^{\top} \left(X_n - \mu_k^{*}\right)$$

weighted sum          empirical cov mat

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N} \log \left[ \sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$
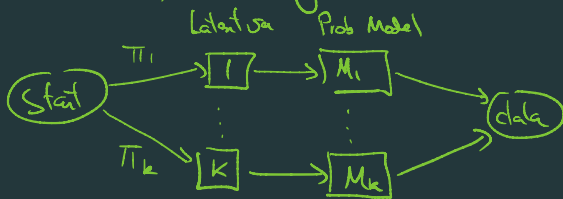
$$\text{s.t} \quad \sum_{k=1}^{K} \pi_k = 1$$

- $\min_{\lambda} \max_{\pi_k} \; \ln\left(p(X|\mu, \Sigma, \pi)\right) + \lambda\left(\sum_{k=1}^{k} \pi_k - 1\right)$

- Inner problem: $\sum_{n=1}^{N} \frac{1}{\pi_k} \gamma_{x_n}(z_k) + \lambda = 0 \quad, \quad \sum \pi_k^* = 1$
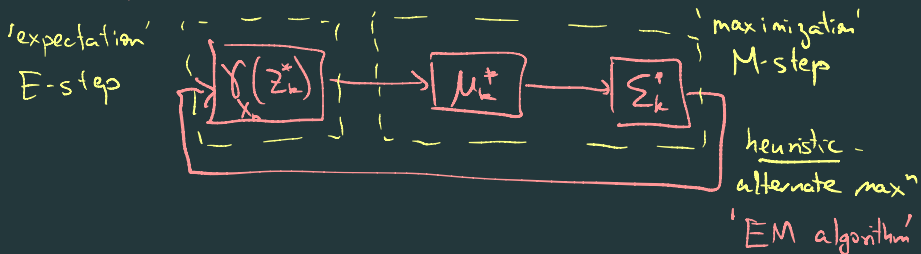  wrt $\pi_k$

$$\Rightarrow \quad \pi_k^* = \frac{N_k}{\sum_{k=1}^{K} N_k} \quad, \quad N_k = \sum_{n=1}^{N} \gamma_{x_n}(z_k^*)$$

# Notes

1) Works for any mixture model



2) Problem - The FOCs are 'circular'

## problems with MLE for GMMs

log-likelihood of data:

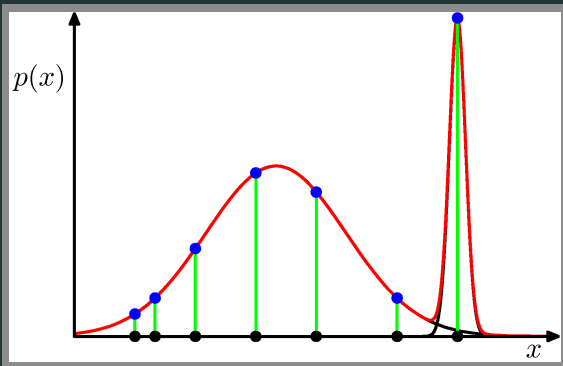$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$

1) log-likelihood  non-convex  ⇒ unclear if unique
maxima

2) What if we change 'labels' of clusters? likelihood
remains same! ⇒ K! alternate maxima

('benign' alternate maxima)

# problems with MLE for GMMs (Singularity Problem) 'kaboom'

log-likelihood of data:

$$\log p\left(X|\mu, \Sigma, \pi\right) = \sum_{n=0}^{N-1} \log \left[\sum_{k=0}^{K-1} \pi_k \mathcal{N}\left(X_n|\mu_k, \Sigma_k\right)\right]$$



$p(x)$

$x$

- Can partition points into 2 clusters
  $X \setminus X_1$ , $X_1$
- Fit $X_1$ with a very sharp dist$^n$
  $\Rightarrow$ likelihood $\nearrow \infty$
  (bad local minima)

# MLE for GMM: an alternate viewpoint

# the EM algorithm

(a)

(b)

(c)

# EM algorithm in action

(e)

(f)