

Lecture 19

Mixtures + EM algo

Q3 3/26 4pm
LVIII PCA posted
HW7 due 3/30

Lecture VII: Clustering: K-means and Mixtures of Gaussians

Marina Meilă
mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath
Cheriton School of Computer Science
University of Waterloo

March, 2026

Paradigms for clustering ✓

K-means clustering ✓

Mixtures of Gaussians and the EM algorithm ✓ ←

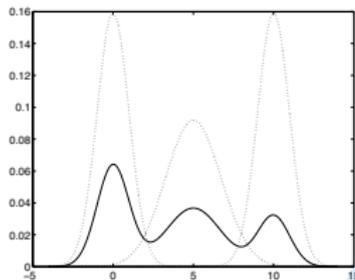
Initialization ←

Special topics in clustering

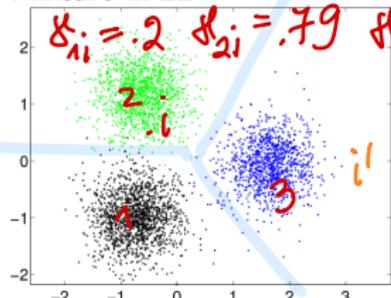
Reading HTF Ch.: 14.3, Murphy Ch.: Ch 11.[1], 11.2.1-3, 11.3, Ch 25, Bach Ch.:

Model based clustering: Mixture models

Mixture in 1D



Mixture in 2D



- ▶ The **mixture density**

$$f(x) = \sum_{k=1}^K \pi_k f_k(x)$$

density on \mathbb{R}^D

$\mathcal{N}(\mu_k, \Sigma_k)$

- ▶ $f_k(x)$ = the **components** of the mixture
 - ▶ each is a density
 - ▶ f called **mixture of Gaussians** if $f_k = \text{Normal}_{\mu_k, \Sigma_k}$
- ▶ π_k = the **mixing proportions**,
 $\sum_k = 1^K \pi_k = 1$, $\pi_k \geq 0$.
- ▶ **model parameters** $\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$
- ▶ The **degree of membership** of point i to cluster k

$$\gamma_{ki} \stackrel{\text{def}}{=} P[x_i \in C_k] = \frac{\pi_k f_k(x)}{f(x)} \text{ for } i = 1:n, k = 1:K \quad (8)$$

- ▶ depends on x_i and on the model parameters

$\delta_{3i} \approx 1$

Given params θ : where did x^i come from?

Bayes rule

$$\Pr[x^i \text{ sampled from } k] = \frac{\overbrace{P[k]}^{\pi_k} \cdot p(x^i | k)}{f(x)} = \delta_{ki} \quad k=1:K$$

the soft clustering

$$\text{Ex} \rightarrow \sum_{k=1}^K \delta_{ki} = 1$$

Criterion for clustering: Max likelihood

- ▶ denote $\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$ (the parameters of the mixture model)
- ▶ Define **likelihood** $P[\mathcal{D}|\theta] = \prod_{i=1}^n f(x_i)$
- ▶ Typically, we use the **log likelihood**

$$l(\theta) = \ln \prod_{i=1}^n f(x_i) = \sum_{i=1}^n \left(\ln \underbrace{\sum_k \pi_k f_k(x_i)}_{f(x_i)} \right) \quad (9)$$

has local maxima

- ▶ denote $\theta^{ML} = \underset{\theta}{\operatorname{argmax}} l(\theta)$
- ▶ θ^{ML} determines a soft clustering γ by (8)
- ▶ a soft clustering γ determines a θ (see later)
- ▶ Therefore we can write

$$\mathcal{L}(\gamma) = -l(\theta(\gamma))$$

Algorithms for model-based clustering

Maximize the (log-)likelihood w.r.t θ

- ▶ directly - (e.g by gradient ascent in θ) 
- ▶ by the EM algorithm (very popular!) 
- ▶ indirectly, w.h.p. by "computer science" algorithms 

w.h.p = with high probability (over data sets)

The Expectation-Maximization (EM) Algorithm

Algorithm Expectation-Maximization (EM)

Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, number clusters K
Initialize parameters $\pi_{1:K} \in \mathbb{R}$, $\mu_{1:K} \in \mathbb{R}^d$, $\Sigma_{1:K} \in \mathbb{R}^{d \times d}$ at random¹
Iterate until convergence

E step (Optimize clustering) for $i = 1 : n$, $k = 1 : K$

cluster data

$$\gamma_{ki} = \frac{\pi_k f_k(x)}{f(x)}$$

M step (Optimize parameters) set $\Gamma_k = \sum_{i=1}^n \gamma_{ki}$, $k = 1 : K$ (number of points in cluster k)

re-estimate Θ

\Leftrightarrow ML with "weighted" x_i 's

$$\pi_k = \frac{\Gamma_k}{n}, \quad k = 1 : K$$

$$\mu_k = \sum_{i=1}^n \frac{\gamma_{ki}}{\Gamma_k} x_i \quad \leftarrow \text{weighted avg}$$

$$\Sigma_k = \frac{\sum_{i=1}^n \gamma_{ki} (x_i - \mu_k)(x_i - \mu_k)^T}{\Gamma_k}$$

► $\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K}$ are the maximizers of $l_c(\theta)$ in (13)

► $\sum_k \Gamma_k = n$

Intuition: $\gamma_{ki} = \text{indicator} = \begin{cases} 1 & i \in C_k \\ 0 & i \notin C_k \end{cases} \Rightarrow \Gamma_k = \sum_i 1_{i \in C_k} = n_k$

¹ Σ_k need to be symmetric, positive definite matrices

The EM Algorithm – Motivation

- Define the **indicator variables**

$$z_{ik} = \begin{cases} 1 & \text{if } i \in C_k \\ 0 & \text{if } i \notin C_k \end{cases} \quad (10)$$

denote $\bar{z} = \{z_{ki}\}_{k=1:K}^{i=1:n}$

- Define the **complete log-likelihood**

$$l_c(\theta, \bar{z}) = \sum_{i=1}^n \sum_{k=1}^K z_{ki} \ln \pi_k f_k(x_i) \quad (11)$$

- $E[z_{ki}] = \gamma_{ki}$
- Then

$$E[l_c(\theta, \bar{z})] = \sum_{i=1}^n \sum_{k=1}^K E[z_{ki}] [\ln \pi_k + \ln f_k(x_i)] \quad (12)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \gamma_{ki} \ln \pi_k + \sum_{i=1}^n \sum_{k=1}^K \gamma_{ki} \ln f_k(x_i) \quad (13)$$

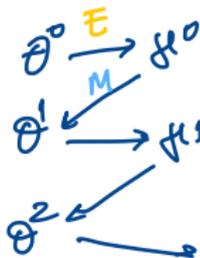
- ▶ If θ known, γ_{ki} can be obtained by (8)
(Expectation)
- ▶ If γ_{ki} known, π_k, μ_k, Σ_k can be obtained by separately maximizing the terms of $E[l_c]$
(Maximization)

Brief analysis of EM

$$Q(\theta, \gamma) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ki} \ln \underbrace{\pi_k f_k(x_i)}_{\theta}$$

Handwritten notes: γ over $i:K$

- ▶ each step of EM increases $Q(\theta, \gamma)$
 - ▶ Q converges to a local maximum
 - ▶ at every local maxi of Q , $\theta \leftrightarrow \gamma$ are fixed point
 - ▶ $Q(\theta^*, \gamma^*)$ local max for $Q \Rightarrow l(\theta^*)$ local max for $l(\theta)$
 - ▶ under certain regularity conditions $\theta \rightarrow \theta^{ML}$
 - ▶ the E and M steps can be seen as projections
-
- ▶ Exact maximization in **M step** is not essential.
Sufficient to increase Q .
This is called **Generalized EM**



θ^* γ^*

Converge?

$$E: \gamma^*(\theta^*) = \theta^*(\gamma^*) : M$$

Probabilistic alternate projection view of EM

- ▶ let z_i = which gaussian generated i ? (random variable), $X = (x_{1:n})$, $Z = (z_{1:n})$
- ▶ Redefine Q

Ex: $\longrightarrow Q(\tilde{P}, \theta) = \mathcal{L}(\theta) - KL(\tilde{P} || P(Z|X, \theta))$ $Q = \text{ELBO}$

where $P(X, Z|\theta) = \prod_i \prod_k P[z_i = k] P[x_i | \theta_k]$

$\tilde{P}(Z)$ is any distribution over Z ,

$KL(P(w) || Q(w)) = \sum_w P(w) \ln \frac{P(w)}{Q(w)}$ the **Kullback-Leibler divergence**

Then,

- ▶ **E step** $\max_{\tilde{P}} Q \Leftrightarrow KL(\tilde{P} || P(Z|X, \theta)) \leftarrow \min_{\tilde{P}}$
 - ▶ **M step** $\max_{\theta} Q \Leftrightarrow KL(P(X|Z, \theta^{old}) || P(X|\theta)) \leftarrow \min_{\theta}$
- ▶ Interpretation: KL is "distance", "shortest distance" = projection

θ_i

The M step in special cases

$$\phi = \dim x$$

$$x \in \mathbb{R}^{\phi}$$

- ▶ Note that the expressions for μ_k, Σ_k = expressions for μ, Σ in the normal distribution, with data points x_i weighted by $\frac{\gamma_{ki}}{\Gamma_k}$

M step

general case

$$\Sigma_k = \sum_{i=1}^n \frac{\gamma_{ki}}{\Gamma_k} (x_i - \mu_k)(x_i - \mu_k)^T$$

$$\Sigma_k = \Sigma$$

$$\Sigma \leftarrow \frac{\sum_{i=1}^n \sum_{k=1}^K \gamma_{ki} (x_i - \mu_k)(x_i - \mu_k)^T}{n}$$

"same shape & size" clusters

$$\Sigma_k = \sigma_k^2 I_d$$

$$\sigma_k^2 \leftarrow \frac{\sum_{i=1}^n \gamma_{ki} \|x_i - \mu_k\|^2}{d \Gamma_k}$$

"round" clusters

$$\Sigma_k = \sigma^2 I_d$$

$$\sigma^2 \leftarrow \frac{\sum_{i=1}^n \sum_{k=1}^K \gamma_{ki} \|x_i - \mu_k\|^2}{nd}$$

"round, same size" clusters



EX:

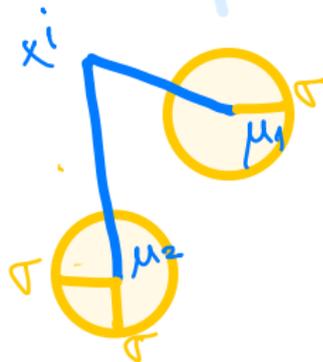
"shared parameters"

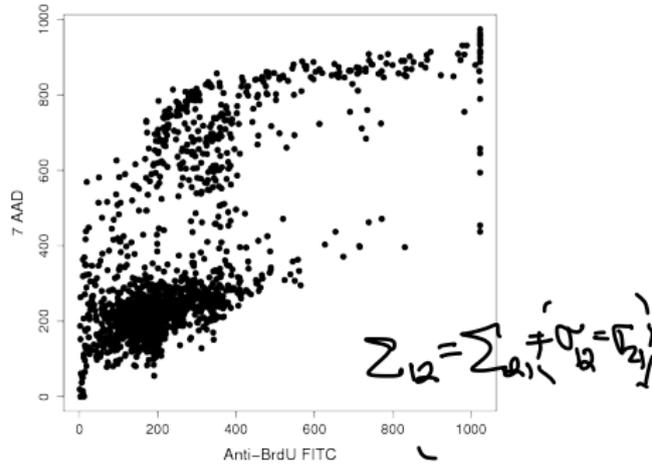
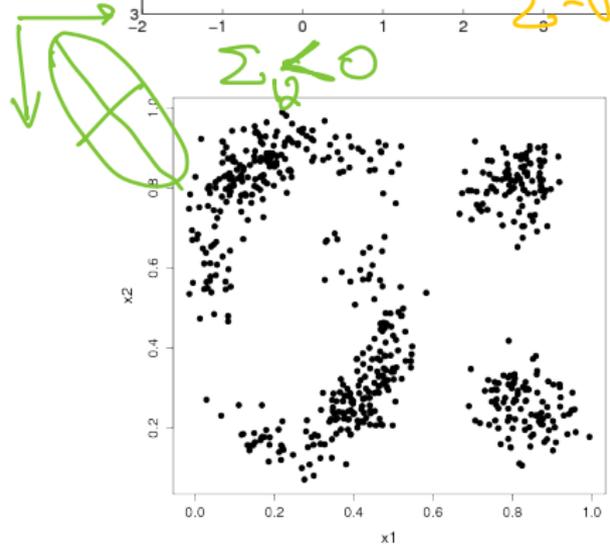
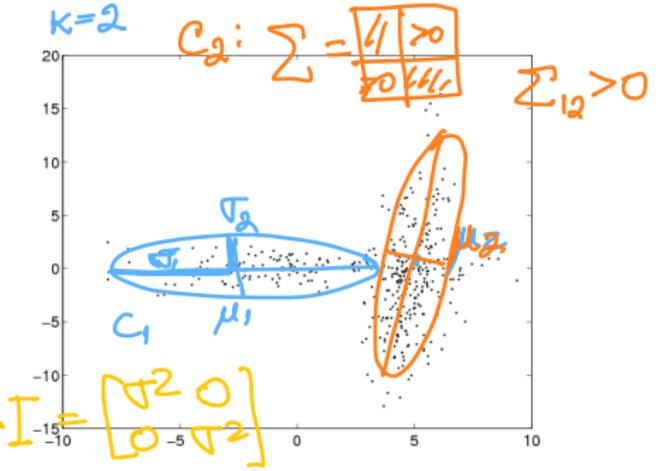
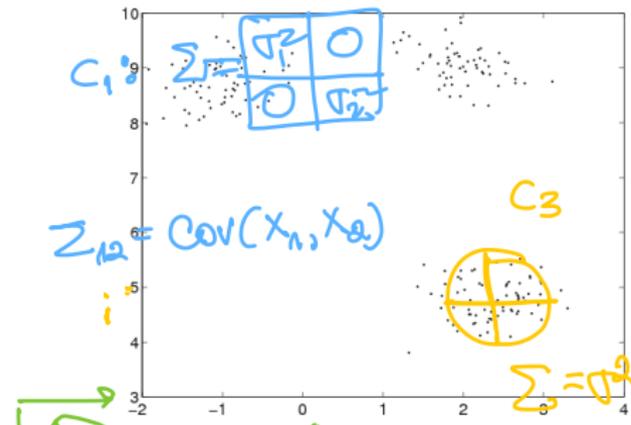
Exercise Prove the formulas above

- ▶ Note also that **K-means** is **EM** with $\Sigma_k = \sigma^2 I_d, \sigma^2 \rightarrow 0$ Exercise Prove it



If clusters in data not gaussian







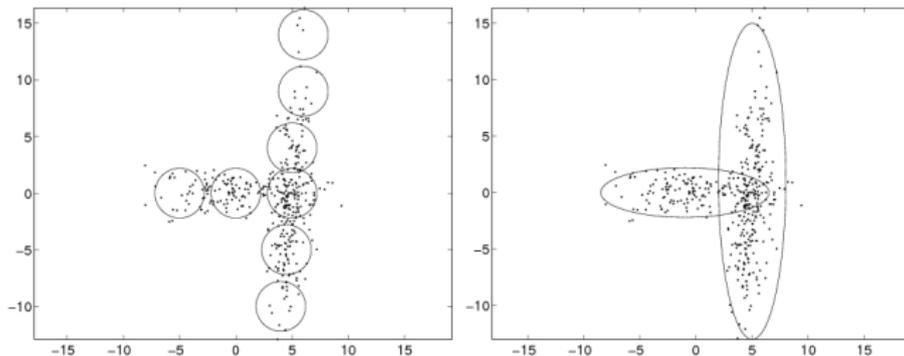
More special cases introduce the following description for a covariance matrix in terms of *volume*, *shape*, *alignment with axes* (=determinant, trace, e-vectors). The letters below mean: I=unitary (shape, axes), E=equal (for all k), V=unequal

- ▶ EII: equal volume, round shape (spherical covariance)
- ▶ VII: varying volume, round shape (spherical covariance)
- ▶ EEI: equal volume, equal shape, axis parallel orientation (diagonal covariance)
- ▶ VEI: varying volume, equal shape, axis parallel orientation (diagonal covariance)
- ▶ EVI: equal volume, varying shape, axis parallel orientation (diagonal covariance)
- ▶ VVI: varying volume, varying shape, equal orientation (diagonal covariance)
- ▶ EEE: equal volume, equal shape, equal orientation (ellipsoidal covariance)
- ▶ EEV: equal volume, equal shape, varying orientation (ellipsoidal covariance)
- ▶ VEV: varying volume, equal shape, varying orientation (ellipsoidal covariance)
- ▶ VVV: varying volume, varying shape, varying orientation (ellipsoidal covariance)

(from)

EM versus K-means

- ▶ Alternates between cluster assignments and parameter estimation
- ▶ Cluster assignments γ_{ki} are probabilistic
- ▶ Cluster parametrization more flexible



- ▶ Converges to local optimum of **log-likelihood**
Initialization recommended by **K-logK** method
- ▶ **Modern algorithms with guarantees** (for e.g. mixtures of Gaussians)
 - ▶ Random projections
 - ▶ Projection on principal subspace
 - ▶ **Two step EM** (=K-logK initialization + one more EM iteration)

A fundamental result

The Johnson-Lindenstrauss Lemma For any $\epsilon \in (0, 1]$ and any integer n , let d' be a positive integer such that $d' \geq 4(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln n$. Then for any set \mathcal{D} of n points in \mathbb{R}^d , there is a map $f: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ such that for all $u, v \in V$,

$$(1 - \epsilon) \|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon) \|u - v\|^2 \quad (14)$$

Furthermore, this map can be found in randomized polynomial time.

- ▶ note that the **embedding dimension** d' does **not** depend on the original dimension d , but depends on n, ϵ
- ▶ show that: the mapping f is linear and that w.p. $1 - \frac{1}{n}$ a **random projection (rescaled)** has this property
- ▶ their proof is elementary Projecting a fixed vector v on a random subspace is the same as projecting a random vector v on a fixed subspace. Assume $v = [v_1, \dots, v_d]$ with $v \sim$ i.i.d. and let \tilde{v} = projection of v on axes $1 : d'$. Then $E[\|\tilde{v}\|^2] = d' E[v_j^2] = \frac{d'}{d} E[\|v\|^2]$. The next step is to show that the variance of $\|\tilde{v}\|^2$ is very small when d' is sufficiently large.

Clustering in high d

1) Random Proj.

2) PCA

easier!

• choose $d' < d$
project $x^{1:n}$ on random subspace

$d' \sim c \ln K$
 $> K$

- faster, less memory
- d' large enough: K 's separate clusters denser
- clusters more Gaussian

$$V \in \mathbb{R}^{d \times d'} \quad V^T V = I_{d'} \quad v = \begin{bmatrix} | \\ | \\ | \end{bmatrix} \text{ basis} \quad \text{Proj } x = V V^T x$$

A two-step EM algorithm \iff K -log k init for K -means

Assumes K spherical gaussians, separation $\|\mu_k^{true} - \mu_{k'}^{true}\| \geq C\sqrt{d}\sigma_k$

1. Pick $K' = O(K \ln K)$ centers μ_k^0 at random from the data
2. Set $\sigma_k^0 = \frac{d}{2} \min_{k \neq k'} \|\mu_k^0 - \mu_{k'}^0\|^2$, $\pi_k^0 = 1/K'$
3. Run one E step and one M step $\implies \{\pi_k^1, \mu_k^1, \sigma_k^1\}_{k=1:K'}$
4. Compute "distances" $d(\mu_k^1, \mu_{k'}^1) = \frac{\|\mu_k^1 - \mu_{k'}^1\|}{\sigma_k^1 - \sigma_{k'}^1}$
5. Prune all clusters with $\pi_k^1 \leq 1/4K'$
6. Run **Fastest First Traversal** with distances $d(\mu_k^1, \mu_{k'}^1)$ to select K of the remaining centers. Set $\pi_k^1 = 1/K$.
7. Run one E step and one M step $\implies \{\pi_k^2, \mu_k^2, \sigma_k^2\}_{k=1:K}$

theorem For any $\delta, \varepsilon > 0$ if d large, n large enough, separation $C \geq d^{1/4}$ the **Two step EM** algorithm obtains centers μ_k so that

$$\|\mu_k - \mu_k^{true}\| \leq \|\text{mean}(C_k^{true}) - \mu_k^{true}\| + \varepsilon \sigma_k \sqrt{d}$$

In practice : • 1 center \leftrightarrow 1 cluster : k -log k

• β 's $\neq 0, 1$ • $\Sigma_k^0 \leftarrow \sigma_0^2 \text{Id}$

• $\pi_k = \frac{1}{K}$

σ_0 not too small

Selecting K for mixture models

The **BIC** (Bayesian Information) **HEURISTIC** criterion

- ▶ let θ_K = parameters for γ_K
- ▶ let $\#\theta_K$ = number independent parameters in θ_K
 - ▶ e.g. for mixture of Gaussians with full Σ_k 's in d dimensions

$$\#\theta_K = \underbrace{K-1}_{\pi_{1:K}} + \underbrace{Kd}_{\mu_{1:K}} + \underbrace{Kd(d-1)/2}_{\Sigma_{1:K}}$$

- ▶ define

$$BIC(\theta_K) = l(\theta_K) - \frac{\#\theta_K}{2} \ln n$$

- ▶ Select K that maximizes $BIC(\theta_K)$
- ▶ selects true K for $n \rightarrow \infty$ and other technical conditions (e.g. parameters in compact set)
- ▶ but theoretically not justified (and overpenalizing) for finite n

1. estimate Δ_K (or θ_K^{ML})
for $K=2,3,\dots,K_{max}$

2. select K^* by comparing

$\Delta_{1:K_{max}}$

$\theta_{1:K_{max}}$

Δ - hard clustering