

today:

- k-means
- EM

K-mean alg.: → can be derived as block-coordinate minimization alg. of obj. fct.:

$$\text{"distortion measure"} \rightarrow J(z, \mu) \triangleq \sum_{i=1}^n \|x_i - \mu_{z_i}\|_2^2 = \sum_{i=1}^n \left(\sum_{j=1}^k z_{i,j} \|x_i - \mu_j\|^2 \right)$$

cluster assignment
 z_1, \dots, z_n
 ← corners of Δ
 ("one hot encoding")

cluster index represented by z_i
 $\mu_1, \dots, \mu_K \in \mathbb{R}^d$
 cluster centroids

alg.:

- 1) initialize $\mu^{(1)}$

2) iterate until convergence:

$$\text{"E step": } z^{(t+1)} = \underset{z \in \text{valid assign.}}{\operatorname{argmin}} J(z, \mu^{(t)})$$

$$\Rightarrow z_{i,j}^{(t+1)} = 1 \text{ for } j^* = \operatorname{argmin}_{j \in \mathbb{R}^d} \|x_i - \mu_j^{(t)}\|$$

$$\text{"M step": } \mu^{(t+1)} = \underset{\mu \in \mathbb{R}^d}{\operatorname{argmin}} J(z^{(t+1)}, \mu)$$

$$\Rightarrow \boxed{\mu_j^{(t+1)} = \frac{\sum_i z_{i,j} x_i}{\sum_i z_{i,j}}} \quad \text{empirical mean of cluster}$$

Visualization:

<https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

Properties of K-means:

1) converge in finite # of iterations to a local min

2) not hard in general to compute global min in z

K-means++: clever initialization scheme which guarantees that alg. is within $\log k$ of global opt. (w.h.p.) + run K-means

→ idea: spread as much as possible the initial means

to avoid



3) choice of K ?

• one heuristic is $J(z, \mu, k) = \sum_i \sum_j z_{i,j} \|x_i - \mu_j\|^2 + \underbrace{\text{C} \cdot k}_{\text{hyperparameter}}$

→ we'll see later in class
 "non-parametric" models

Side reference I mentioned:
see <https://icml.cc/2012/papers/291.pdf>
for interpreting regularized K-means as
approximate inference in a Dirichlet process mixture
model... [by Kulis & Jordan]

→ we'll see later in class
"non-parametric" models
where "k" is basically infinite
and can get $p(\mathbf{x} | \text{data})$

hyperparameter

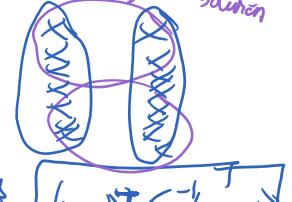
e.g., Dirichlet process
mixture model

K-means
solution

4) K-mean is very sensitive in distance measure: it assumes
spherical clusters

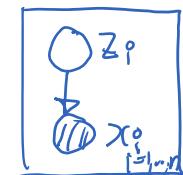
GMM fixes
that problem

$$\text{Mahalanobis distance } d_{\Sigma}(\mathbf{x}, \mathbf{x}') \triangleq \sqrt{(\mathbf{x}-\mathbf{x}')^T \Sigma^{-1} (\mathbf{x}-\mathbf{x}')}$$



EM - maximum likelihood in latent variable model

Setup:



z latent variable

x observed variable

$$\text{log-likelihood } \log p(\mathbf{x}; \boldsymbol{\theta}) = \log \left(\prod_{i=1}^n p(x_i; \boldsymbol{\theta}) \right)$$

$$= \sum_{i=1}^n \log p(x_i; \boldsymbol{\theta})$$

$$= \sum_{i=1}^n \log \left[\sum_{z_i} p(x_i; z_i; \boldsymbol{\theta}) \right]$$

problem? → yields multi-modal
opt. problems (non-convex)

options MLE in latent variable model

1) do gradient ascent on a non-convex obj.

2) EM alg. → block-coordinate ascent on auxiliary fn., which lower bounds $\log p(\mathbf{x}; \boldsymbol{\theta})$
nice interpretation in terms of filling "missing data"

i.e. E step → fill z with "soft-values"

M step → max w.r.t. $\boldsymbol{\theta}$ for fully observed model



trick overview:

$$\log \sum_z p(x_i, z) = \log \sum_z q(z) \frac{p(x_i, z)}{q(z)}$$

$$= \log \mathbb{E}_q \left[\frac{p(x_i, z)}{q(z)} \right]$$

$$\geq \mathbb{E}_q \left[\log \frac{p(x_i, z)}{q(z)} \right]$$

Jensen's
inequality
trick?

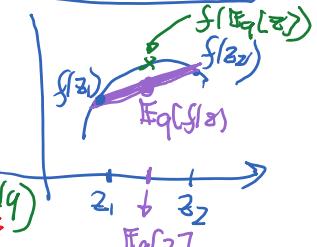
$$= \sum_z q(z) \log \frac{p(x_i, z)}{q(z)} - \sum_z q(z) \log q(z)$$

$$\triangleq \mathcal{L}(q, \boldsymbol{\theta}) \triangleq \mathbb{E}_q \left[\log p(x_i, z; \boldsymbol{\theta}) \right] + H(q)$$

Jensen's inequality

$$\mathbb{E}_q[f(g(z))] \leq f(\mathbb{E}_q[g(z)])$$

when f is concave



$$\text{true: } \hat{\mathcal{L}}(q, \theta) \triangleq \mathbb{E}_q [\log p(x, z; \theta)] + H(q) \quad \begin{array}{c} z_1 \\ \vdots \\ z_2 \end{array} \quad \mathbb{E}_q[z] \quad \begin{array}{c} \text{"expected complete} \\ \text{"entropy of"} \\ \text{"log-likelihood"} \end{array}$$

we have $\log p(x; \theta) \geq \hat{\mathcal{L}}(q, \theta) \quad \forall q \in \mathcal{Q}$

EM algorithm: E step: $q_{t+1} \triangleq \arg \max_{q \in \text{dist. over } z} \mathcal{J}(q, \theta_t) \Rightarrow q_{t+1}(z) = \frac{p(z|x; \theta_t)}{\mathbb{P}_q[z]}$

M Step: $\theta_{t+1} \triangleq \arg \max_f \mathcal{L}(q_{t+1}, \theta)$
 $= \arg \max_{\theta} \mathbb{E}_{q_{t+1}} [\log p(x, z; \theta)]$

this is another MLE problem, but for complete information

(often, replace Z with $\mathbb{E}_q[Z]$ in this expression)

16 hours

* we had $\log p(x; \theta) \geq \hat{\mathcal{L}}(q, \theta) \quad \log(\mathbb{E}_q[g(z)]) \geq \mathbb{E}_q[\log(g(z))]$

in Jensen's ineq., you get a strict inequality unless the dist. for $g(z)$ is degenerate (i.e. $g(z)$ takes only one value) i.e. when $g(z) = \text{constant}$
 \log (when f is strictly concave)

above $g(z) = \frac{p(x|z)}{q(z)} = \text{constant } \forall z \Rightarrow q(z) \propto p(x|z)$

$\mathcal{J}(q_{t+1}, \theta_t) = \log p(x; \theta_t) \geq \hat{\mathcal{L}}(q, \theta_t) \quad \forall q$ i.e. $\boxed{q(z) = p(z|x; \theta_t)}$

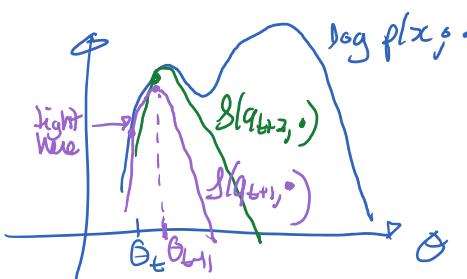
$\Rightarrow q_{t+1}$ maximizes $\mathcal{J}(q, \theta_t)$ w.r.t. q i.e. $\arg \max_q \mathcal{J}(q, \theta_t) = p(z|x; \theta_t) = q_{t+1}$

$\hat{\mathcal{J}}(q_{t+1}, \theta_t) = \log p(x; \theta_t)$

properties of EM algorithm

a) $\log p(x; \theta_{t+1}) \geq \log p(x; \theta_t)$

proof: $\log p(x; \theta_{t+1}) \geq \hat{\mathcal{L}}(q_{t+1}, \theta_{t+1})$



$\geq \hat{\mathcal{L}}(q_{t+1}, \theta_t)$
 [because q_{t+1} max. $\hat{\mathcal{L}}(q, \theta_t)$]

$\hat{\mathcal{L}}(q_{t+1}, \theta_t) = \log p(x; \theta_t) //$

b) θ_t in EM converges to a stationary pt. of $\log p(x; \theta)$

i.e. $\nabla_\theta \log p(x; \theta) \Big|_{\hat{\theta}} = 0$

Life means, initialization is crucial

→ usually do random restarts

• for GMM, could use k-means to initialize the μ 's

$$c) \mathcal{J}(q, \Theta) = \mathbb{E}_q [\log \frac{p(x, z; \Theta)}{q(z)}]$$

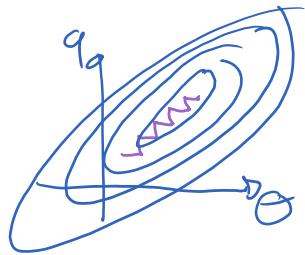
$$\log p(x; \Theta) - \mathcal{J}(q, \Theta) = -\mathbb{E}_q [\log \frac{p(z, x; \Theta)}{q(z)p(x; \Theta)}]$$

$$\log p(x; \Theta) \leftarrow \begin{cases} \mathbb{E}_q [\log q(\cdot)] \\ \mathbb{E}_q [\log \frac{q(z)}{p(z|x; \Theta)}] \end{cases} \stackrel{\Delta}{=} \text{KL}(q(\cdot) \parallel p(\cdot|x; \Theta))$$

KL divergence

We will revisit this
for variational inference $q \in \mathcal{Q}$

"Simple distributions"



Block-coordinate method sometimes slow

For GMM model:



$$z_i \sim \text{Mult}(\pi)$$

$$x_i | z_i = j \sim N(\mu_j, \Sigma_j)$$

shorthand to say $z_{i,j} = 1$

$$\Theta = (\pi, (\mu_j)_{j=1}^K, (\Sigma_j)_{j=1}^K)$$

notation: $x = x_{1:n}$
 $z = z_{1:n}$

$$\text{exercise: } p(z|x) = \prod_{i=1}^n p(z_i|x_i) = \prod_{i=1}^n p(z_i|x_i)$$

complete log-likelihood:

$$\log p(x, z; \Theta) = \sum_{i=1}^n [\underbrace{\log p(x_i|z_i; \Theta)}_{\text{Gaussian}} + \underbrace{\log p(z_i; \Theta)}_{\text{Multinomial}}]$$

$$= \sum_{i=1}^n \left[\sum_{j=1}^K z_{i,j} \log N(x_i | \mu_j, \Sigma_j) + \sum_{j=1}^K z_{i,j} \log \pi_j \right]$$

$$\mathbb{E}_q [\log p(x, z; \Theta)] = \sum_{i=1}^n \sum_{j=1}^K \mathbb{E}_q [z_{i,j}] (\log N(x_i | \mu_j, \Sigma_j) + \log \pi_j)$$

$$\mathbb{E}_q [1 \{z_{i,j}=1\}] = q(z_{i,j}=1)$$

$\mathbb{E}_q [z_{i,j}] = q(z_{i,j}=1)$ [marginal prob.]
dimin FIM $n_r(x) = n \{x \mid x \in \mathcal{B}_w\}$

$$\mathbb{E}_q[\mathbb{1}_{\{z_{ij}=1\}}] = q(z_{ij}=1)$$

↑ $q(z_{ij}) = q(z_{ij}=1) \quad \text{(marginal prob)}$

during EM, $q_{t+1}(z) = p(z|x; \theta_t)$

Weight $\gamma_{ij}^t \triangleq p(z_{ij}=1|x_i; \theta_t) = q_{t+1}(z_{ij}=1)$

E-step is computing $q_{t+1}(z) \triangleq p(z|x; \theta_t)$

$$= \prod_i p(z_i|x_i; \theta_t)$$

$$\Rightarrow q_{t+1}(z_i) = p(z_i|x_i; \theta_t)$$

$$\propto p(x_i|z_i; \theta_t) p(z_i; \theta_t)$$

\downarrow Gaussian \downarrow $\pi_{z_i}^t$

$$\gamma_{ij}^t = q_{t+1}(z_{ij}=1) = \frac{\pi_j^{(t)} N(x_i | \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{k=1}^K \pi_k^{(t)} N(x_i | \mu_k^{(t)}, \Sigma_k^{(t)})}$$

$\left\{ \begin{array}{l} p(x_i, z_i | \theta^{(t)}) \\ p(x_i | \theta^{(t)}) \end{array} \right.$

E step for GMM: computing $\gamma_{ij}^{(t)}$ for $i=1, \dots, n$ using $\theta^{(t)}$

M step : $\max_{\{\pi_j, \mu_j, \Sigma_j\}} \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij}^{(t)} [\log p(x_i|\mu_j, \Sigma_j) + \log \pi_j]$

~~exercise 9~~

$$\hat{\pi}_j^{(t+1)} = \frac{\sum_i \gamma_{ij}^{(t)}}{n} \quad \text{"soft counts"}$$

M step for EM for GMM

$$\hat{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} x_i}{\sum_{i=1}^n \gamma_{ij}^{(t)}}$$

$$\hat{\Sigma}_j^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ij}^{(t)} (x_i - \hat{\mu}_j^{(t+1)})(x_i - \hat{\mu}_j^{(t+1)})^\top}{\sum_{i=1}^n \gamma_{ij}^{(t)}}$$

- initialize: e.g. $\mu_j^{(0)}$ from k-means++

$\Sigma_j^{(0)}$ big spherical covariance $\Sigma_j^{(0)} = \sigma^2 I$

$\pi_j^{(0)}$: proportions from k-means++

- EM step in GMM with fixed $\Sigma_j = \sigma^2 I$ with $\sigma^2 \rightarrow 0$

→ get k-means alg
"hand EM"