

Lecture 10 - scribbles

Tuesday, October 8, 2019 14:35

today: finish MLE for Gaussian
 * K-means, GMM, EM

continuation: MLE for multivariate Gaussian

example 2: derivative of $f(A) \triangleq \log \det(A)$ where $\begin{pmatrix} A \text{ symmetric} \\ A \succ 0 \end{pmatrix}$

can represent derivative of a fct. from matrix to scalar, as a matrix

$$f(A+\Delta) - f(A) = \text{tr}(f'(A)^T \Delta) + o(\|\Delta\|)$$

$$= \langle f'(A), \Delta \rangle + \dots$$

$$\begin{aligned} & \log \det(A+\Delta) - \log \det(A) \quad A \succ 0 \Rightarrow \text{invertible \& has unique positive square root} \\ &= \log \det(A^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) A^{1/2}) - \log \det(A) \\ &= \log |A|^{1/2} |I + A^{-1/2} \Delta A^{-1/2}| |A|^{1/2} - \log |A| \\ &= \log \det |I + A^{-1/2} \Delta A^{-1/2}| \quad \text{use } \det(A) = \prod_i \lambda_i(A) \quad \leftarrow \text{e-values} \\ &= \sum_i \log \lambda_i(I + A^{-1/2} \Delta A^{-1/2}) \\ &= \sum_i \log(1 + \lambda_i(A^{-1/2} \Delta A^{-1/2})) \quad \log(1+x) = x + O(x^2) \text{ for } |x| < 1 \\ &= \sum_i \lambda_i(A^{-1/2} \Delta A^{-1/2}) + O(\lambda_i(\dots)^2) \\ & \quad \quad \quad o(\|\Delta\|) \quad \text{tr}(A) = \sum_i \lambda_i(A) \\ &= \text{tr}(A^{-1/2} \Delta A^{-1/2}) + o(\|\Delta\|) \\ &= \text{tr}(A^{-1} \Delta) + o(\|\Delta\|) \quad \text{(recall } A \text{ is symmetric)} \\ & \quad \quad \quad \langle A^{-1}, \Delta \rangle \quad \Rightarrow \quad \boxed{\frac{d}{dA} \log \det(A) = A^{-1}} \end{aligned}$$

back to log-likelihood for Gaussian:

$$+\frac{n}{2} \log |\Sigma^{-1}| - \frac{n}{2} \langle \Sigma^{-1}, \tilde{\Sigma}(\mu) \rangle$$

concave as function of $\Lambda = \Sigma^{-1}$

derivative $\rightarrow n(\Sigma^{-1})^{-1} \quad n \tilde{\Sigma}(\dots)$ want \rightarrow

$$\Lambda \succeq \tilde{\Sigma}(\mu)$$

derivative w.r.t to Σ^{-1}

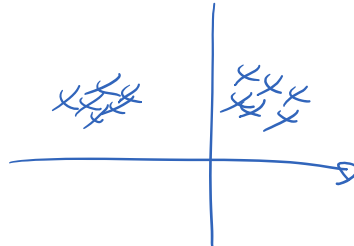
$$\frac{n}{2} (\Sigma^{-1})^{-1} - \frac{n}{2} \tilde{\Sigma}(\mu) = 0 \quad \text{want}$$

$$\Rightarrow \hat{\Sigma}_{MLE} = \tilde{\Sigma}(\mu_{MLE}) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^T$$

15h

Unsupervised Learning

here X without any labels Y



consider the Gaussian mixture model (GMM)
(can be obtained from FLD:)

$$Y \sim \text{Mult}(\pi) \quad \pi \in \Delta_K$$

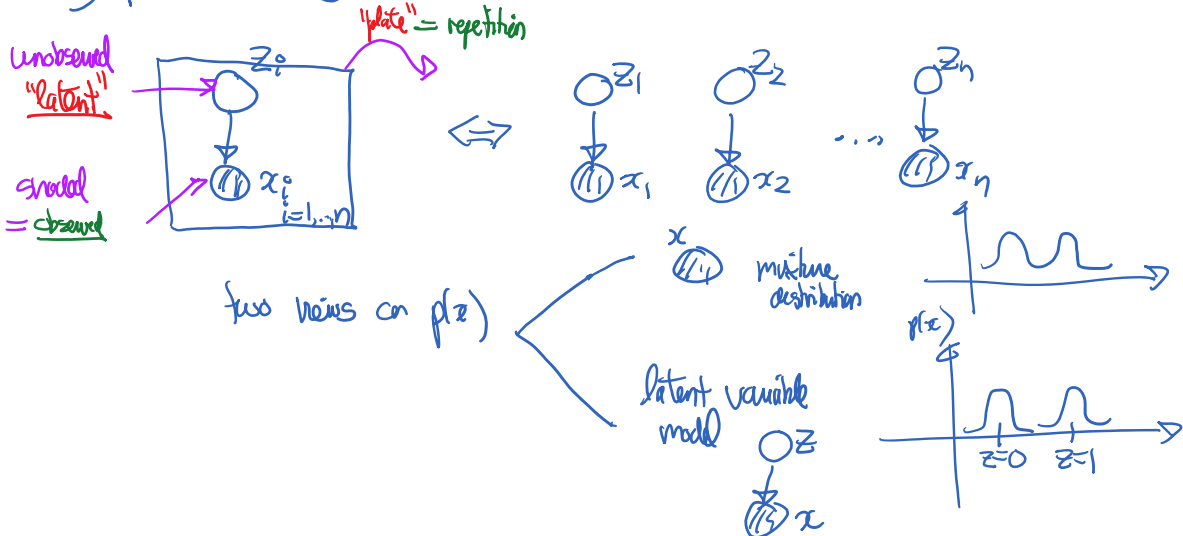
$$X | Y=j \sim N(\mu_j, \Sigma)$$

$$p(x) = \sum_y p(x|y) = \sum_y p(x|y)p(y) = \sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma)$$

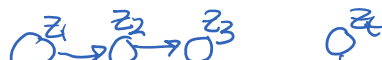
"GMM model"

more generally, can have different Σ_j per class

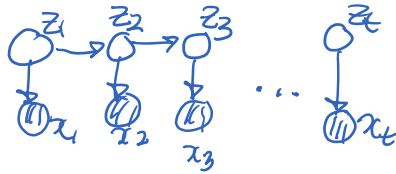
graphical model for this "latent variable model"



(later in class, we will add time structure: HMM)



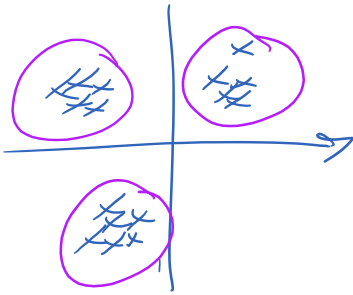
(x-axis in cross, we would add time structure: HMM)



K-means → to do clustering i.e. group data

(can be seen as a limit of GMM $\sigma^2 \rightarrow 0$)

we want to get a cluster assignment for every data pt. x_i



represent $z_{ij} = 1$ to mean that x_i in cluster j

$$j = 1, \dots, K$$

of clusters (specified in advance for K-means)

applications: • vector quantization

[5h3]

• in computer vision: use K-means to get "bag of visual words" representation of image patches

K-mean alg.

→ can be derived as a block-coordinate minimization alg. of objective f.t.

$$J(z, \mu) \triangleq \sum_{i=1}^n \left(\sum_{j=1}^K z_{ij} \|x_i - \mu_j\|^2 \right) = \sum_i \|x_i - \mu_{z_i}\|^2$$

cluster assign. $z_1, \dots, z_n \in \text{corners of } \Delta_K$

$\mu_1, \dots, \mu_K \in \mathbb{R}^d$ cluster means

"distance measure"

alg: 1) initialize $\mu^{(1)}$

2) iterate until convergence

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

$$\Rightarrow z_{ij}^{(t+1)} = 1 \text{ for } j^* = \underset{j}{\text{argmin}} \|x_i - \mu_j^{(t)}\|$$

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

$$\Rightarrow \mu_i^{(t+1)} = \left| \sum_{z_{ij} > 0} x_i \right| \text{ empirical}$$

demo: <http://web.stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>

$$\mu_j^{(t+1)} = \frac{\sum_i z_{ij} x_i}{\sum_i z_{ij}}$$

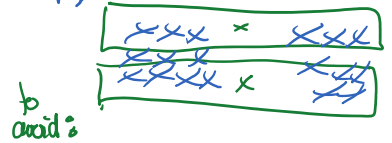
empirical mean of cluster

properties of k-means:

- 1) converges in finite # of iterations to a local min
- 2) NP hard in general to find best z

kmeans++: clever initialization scheme which guarantees that obj. is within logK of global opt. (w.h.p.)

→ idea: spread out as much as possible the initial means



3) choice of k?

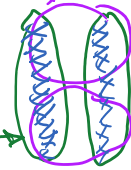
• one heuristic is: $J(\mu, z, k) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|x_i - \mu_j\|^2 + \lambda k$

→ we'll see often in class "non-parametric" models where "k" is usually infinite and can get f(k|data)

hyperparameters

e.g. Dirichlet process mixture model

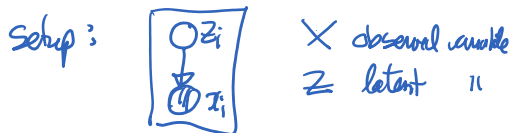
4) k-mean is very sensitive on distance measure: it assumes spherical cluster



↳ GMM fixes that

16/12

EM - maximum likelihood in latent variable model



$$\begin{aligned} \log \text{likelihood } \log p(x_{1:n}, \theta) &= \log \prod_{i=1}^n p(x_i; \theta) \\ &= \sum_{i=1}^n \log p(x_i; \theta) \\ &= \sum_{i=1}^n \log \left[\sum_z p(x_i, z_i; \theta) \right] \end{aligned}$$

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

problem?
→ gives multi-modal difficult opt. problem (non-convex)

options for ML in latent variable model

1) do gradient ascent on non-convex obj.

2) EM alg. → block-coordinate ascent on auxiliary ~~set~~ which lower bounds $\log p(x_{1:n}; \theta)$

nice interpretation in terms of filling "missing data"

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

M step → max w.r. to θ for fully observed model

Trick overview:

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

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

Jensen's inequality
 $\mathbb{E}_q [f(z)] \leq f(\mathbb{E}_q [z])$
 when f is concave



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

$$\triangleq \mathcal{L}(q, \theta) \triangleq \mathbb{E}_{q(z)} [\log p(x, z; \theta)] + H(q)$$

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

entropy of q

in Jensen's inequality, we get equality only if $f(z) = \text{constant}$.

$$\text{i.e. } \frac{p(x, z)}{q(z)} = \text{constant} \quad \forall z \Rightarrow q^*(z) \propto p(x, z)$$

$$\text{i.e. } q^*(z) = p(z|x; \theta)$$

this means argmax
 $q \in \text{all dist. over } z$
 $\mathcal{L}(q, \theta_t) = p(z|x; \theta_t)$

ignore

EM algorithm: E step: $Q_{t+1} \stackrel{\text{def}}{=} \arg \max_q \mathcal{J}(q, \Theta_t) \Rightarrow Q_{t+1}(z) = p(z|x, \Theta_t)$

M step: $\Theta_{t+1} \stackrel{\text{def}}{=} \arg \max_{\Theta} \mathcal{J}(Q_{t+1}, \Theta)$

$$= \arg \max_{\Theta} \mathbb{E}_{Q_{t+1}(z)} \log(p(x, z; \Theta_{t+1}))$$

"expected complete log-likelihood"

this is another ML problem, but for complete information

[usually, replace z with $\mathbb{E}_q[z]$]