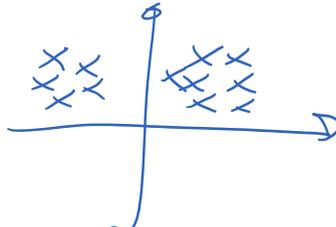


today: • Unsupervised learning
• K-means & EM

Unsupervised learning

here X without any labels Y



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

$$Y \sim \text{Mult}(z) \quad \pi \in \Delta^k$$

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

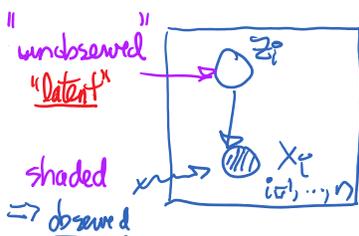
$$p(x) = \sum_y p(x, y) = \sum_y p(x|y)p(y) = \sum_{j=1}^k \pi_j N(x | \mu_j, \Sigma)$$

[extension of FLD to multiple classes]

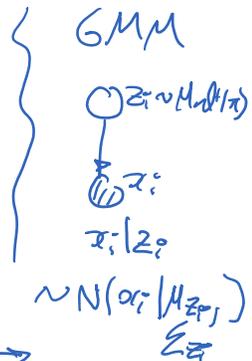
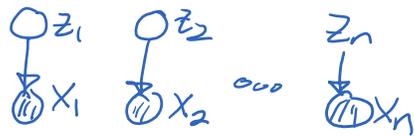
"GMM model"

more generally can have Σ_j per class

graphical models for this "latent variable model"



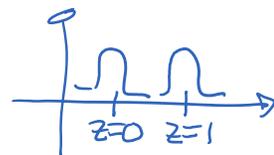
"plate" = repetition



two views on $p(x)$

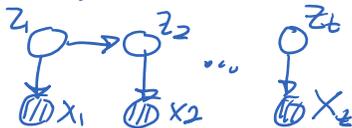


latent variable model



(structural representation)

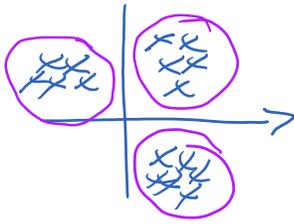
(later in class we will add time structure = HMM)



$$p(x) = \sum_z p(x|z)p(z) = \sum_z p(z) p(x|z)$$

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

$$p(\theta | \text{data})$$



we want to get a cluster assignment for every data point x_i

represent $z_{i,j} = 1$ to mean that x_i belongs to cluster j

$j=1, \dots, K$
 \downarrow
 # of clusters (specified in advance for K-means)

- Applications:
- vector quantization (compression)
 - in computer vision & use K-means to get "bag of visual words" representation of image patches
 - many others?

15h36

K-mean alg. → can be derived as a block-coordinate minimization alg. of objective function:

"distortion measure"

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

$z_1, \dots, z_n \in \text{corners of } \Delta^K$ ("one-hot encoding")
 $\mu_1, \dots, \mu_K \in \mathbb{R}^d$ cluster centers
 μ_{z_i} cluster index represented by z_i

- alg.:
- 1) initialize $\mu^{(1)}$
 - 2) iterate until convergence:

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

$$\Rightarrow z_{i,j^*} = 1 \text{ for } j^* = \text{argmin}_j \|x_i - \mu_j^{(t)}\|$$

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

$$\Rightarrow \mu_j^{(t+1)} = \frac{\sum_i z_{i,j} x_i}{\sum_i z_{i,j}}$$

empirical mean of cluster

Visualization:

<https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html>

properties of K-means:

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

- 1) converge in finite # of iterations to a local min
- 2) NP hard in general to compute the global minimum in Z

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

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



3) choice of k?

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

hyperparameter

→ we'll see later in class "non-parametric" models where "k" is basically infinite and can get p(k|data)

e.g. Dirichlet process mixture model

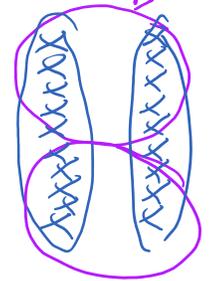
k-mean clusters

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, I was wrong on the author! ;)]

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

↳ GMM fixes that problem



Mahalanobis distance = $d(x, x') \cong \sqrt{(x-x')^T \Sigma^{-1} (x-x')}$