

- Math stuff

- Gaussian Mixture Model (GMM)

- Expectation-Maximization Algo (EM)

- EM applied to GMM

- demo

Section notes on EM  
Algo + Gaussian mixtures

6/13/14

Douglas Rubin

Math stuff

① Indicator Functions:  $\mathbb{1}\{\text{Condition}\} = \begin{cases} 1 & \text{if condition is true} \\ 0 & \text{if condition is false} \end{cases}$

examples:  $\mathbb{1}\{2 \neq 5\} = 1$

$\mathbb{1}\{-3 > 0\} = 0$

- Convenient in summations and integrals bc sometimes they "kill" the sum/integral

- Also convenient for counting things, for example the # of misclassifications:

$$\text{misclassification rate} = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\hat{y}(x^{(i)}) \neq y^{(i)}\}$$

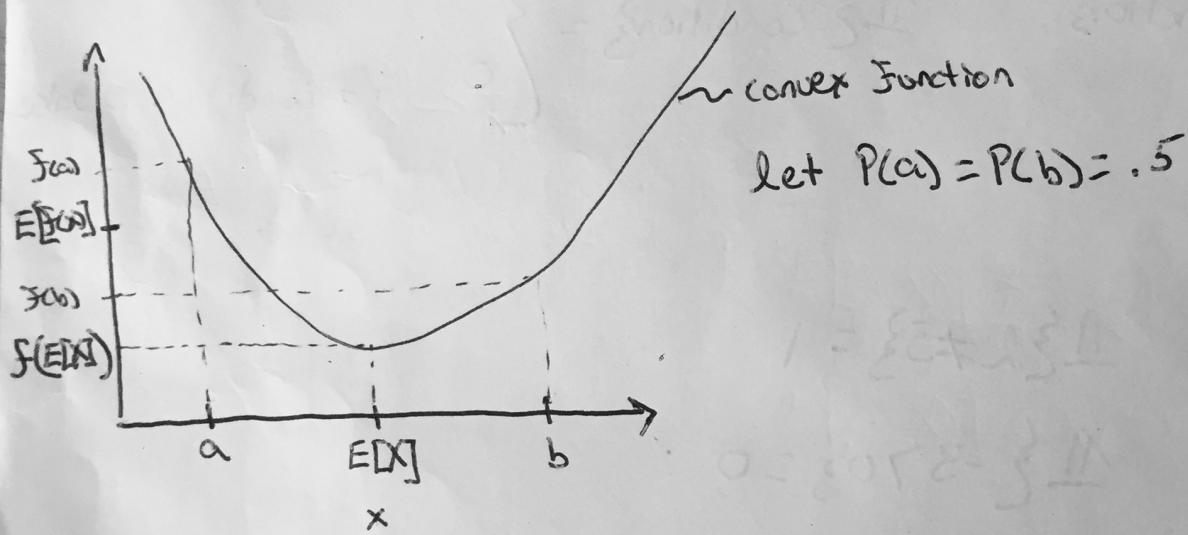
## ② Jensen's inequality

let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $X \in \mathbb{R}^n$  be a random variable. Then

$$E[f(X)] \geq f(E[X]) \quad (\text{if } f \text{ is convex})$$

$$E[f(X)] \leq f(E[X]) \quad (\text{if } f \text{ is concave})$$

Moreover, if  $f$  is strictly convex/concave  $E[f(X)] = f(E[X])$   
holds only if  $X$  is a constant.



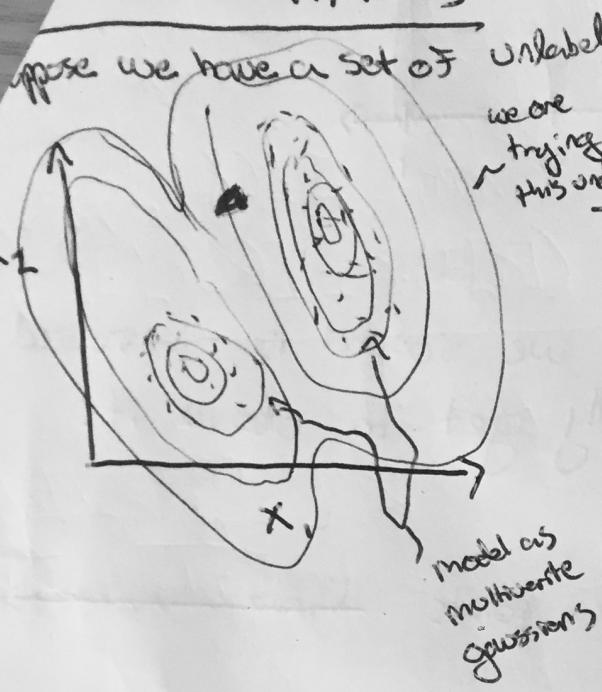
$$\text{clearly } E[f(a)] \geq f(E[X])$$

e.g. let  $P$  be a prob. dist., then

$$\sum_{j=1}^N p_j \log(x_j) \leq \log \sum_{j=1}^N p_j x_j \quad (\log \text{ is strictly concave})$$

## Gaussian Mixtures

(3)

- Suppose we have a set of unlabeled cluster  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ;  $x^{(i)} \in \mathbb{R}^n$
- 
- we are trying to learn this underlying density field
- GMM achieves 2 things
- ① estimates underlying density Field  
(possibly good for outlier detection)
  - ② clusters the data
    - does so in a "soft" way,  
i.e., assigns probs. of being associated to each cluster. may be useful in some scenarios.

\* note: We will work through the case of  $K=2$  clusters, but the generalization to arbitrary  $K$  is very straight forward.

let's model this

- ① Flip an <sup>on fair</sup> coin
- ② if heads (class 1) draw a data point from Gaussian 1  
" tails (class 0) " " " " "

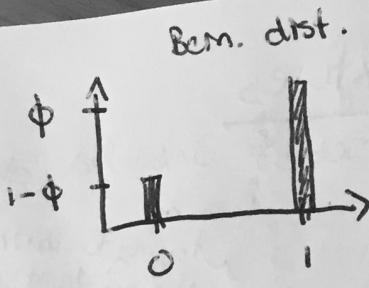


(in math:

$$Z^{(i)} \sim \text{Bern}(\phi)$$

$$X^{(i)} | Z^{(i)} = j \sim N(\mu_j, \Sigma_j)$$

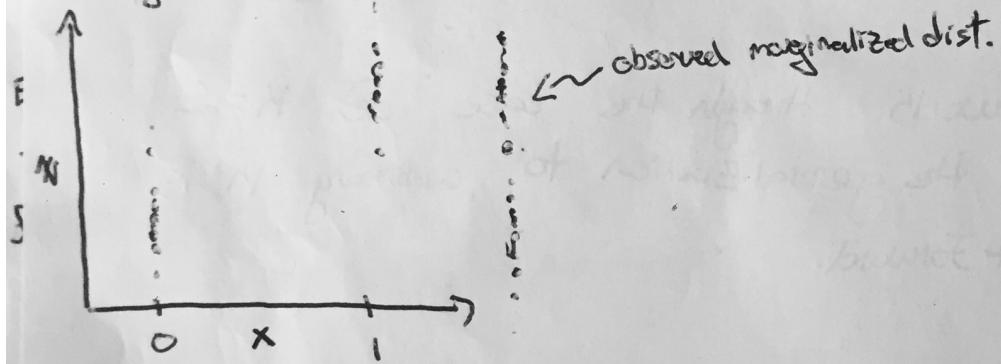
$$\mu_j \in \mathbb{R}^n, \Sigma_j \in \mathbb{R}^{n \times n}$$



called a "latent" RV. that is a RV. we suspect is involved in the process, but we don't actually get to observe it.

Consider a 1D observed set of data points

joint density dist



$$\text{let } \Theta = \{\phi, \mu_0, \mu_1, \Sigma_0, \Sigma_1\} \text{ (for convenience)}$$

lets write down all the prob. dists:

$$P(Z^{(i)}; \Theta) = \phi^{Z^{(i)}} (1-\phi)^{1-Z^{(i)}}$$

$$P(X^{(i)} | Z^{(i)} = j) = \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left\{ -\frac{1}{2} (X^{(i)} - \mu_j)^T \Sigma_j^{-1} (X^{(i)} - \mu_j) \right\}$$

$$P(X^{(i)}, Z^{(i)}; \Theta) = P(X^{(i)} | Z^{(i)}; \Theta) P(Z^{(i)}; \Theta) \quad (\text{By defn. of cond. prob})$$

$$P(Z^{(i)} | X^{(i)}; \Theta) = \frac{P(X^{(i)} | Z^{(i)}; \Theta) P(Z^{(i)}; \Theta)}{\sum_{j=1}^J P(X^{(i)} | Z^{(i)} = j) P(Z^{(i)} = j; \Theta)} \quad (\text{Bayes rule})$$

(5)

solve for all the parameters  $\Theta$  in the usual way.

MML

$$l(\Theta) = \log \prod_{i=1}^m p(X^{(i)}; \Theta)$$

$$= \sum_{i=1}^m \log p(X^{(i)}; \Theta) \leftarrow \text{we need to introduce our model here, so we realize that this is the joint dist. marginalized over the } Z_i \text{.}$$

$$= \sum_{i=1}^m \log \sum_{Z^{(i)}=0}^{\infty} p(X^{(i)}, Z^{(i)}; \Theta)$$

Problem - plugging in the dists. and trying to maximize  $l(\Theta)$  by taking derivatives wpt  $\Theta$ , we find that this is analytically intractable, mostly due to summation in the log. We need to try to push the summation out of the log (let's try Jensen's ineq.)



## EM algo

- in a somewhat unmotivated fashion I will introduce  $m$  different prob. dists.  $Q_1(Z^{(i)}), Q_2(Z^{(i)}), \dots, Q_m(Z^{(i)})$

( $\sum_{Z^{(i)}=0}^{Q_i} Q_i(Z^{(i)}) = 1, Q_i(Z^{(i)}) \geq 0$ ), where I have yet to

specify the functional form of the  $Q_i$ 's? note it is important that each term have its own  $Q_i$ , as you will see later

~~$$= \log \left\{ \sum_{Z^{(i)}=0}^{Q_i} Q_i(Z^{(i)}) \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})} \right\}$$~~

$$= \log \left\{ \sum_{Z^{(i)}=0}^{Q_i} Q_i(Z^{(i)}) \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})} \right\} + \dots + \log \left\{ \sum_{Z^{(i)}=0}^{Q_i} Q_i(Z^{(i)}) \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})} \right\}$$

$$\geq \sum_{Z^{(i)}=1}^{\infty} Q_i(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})}$$

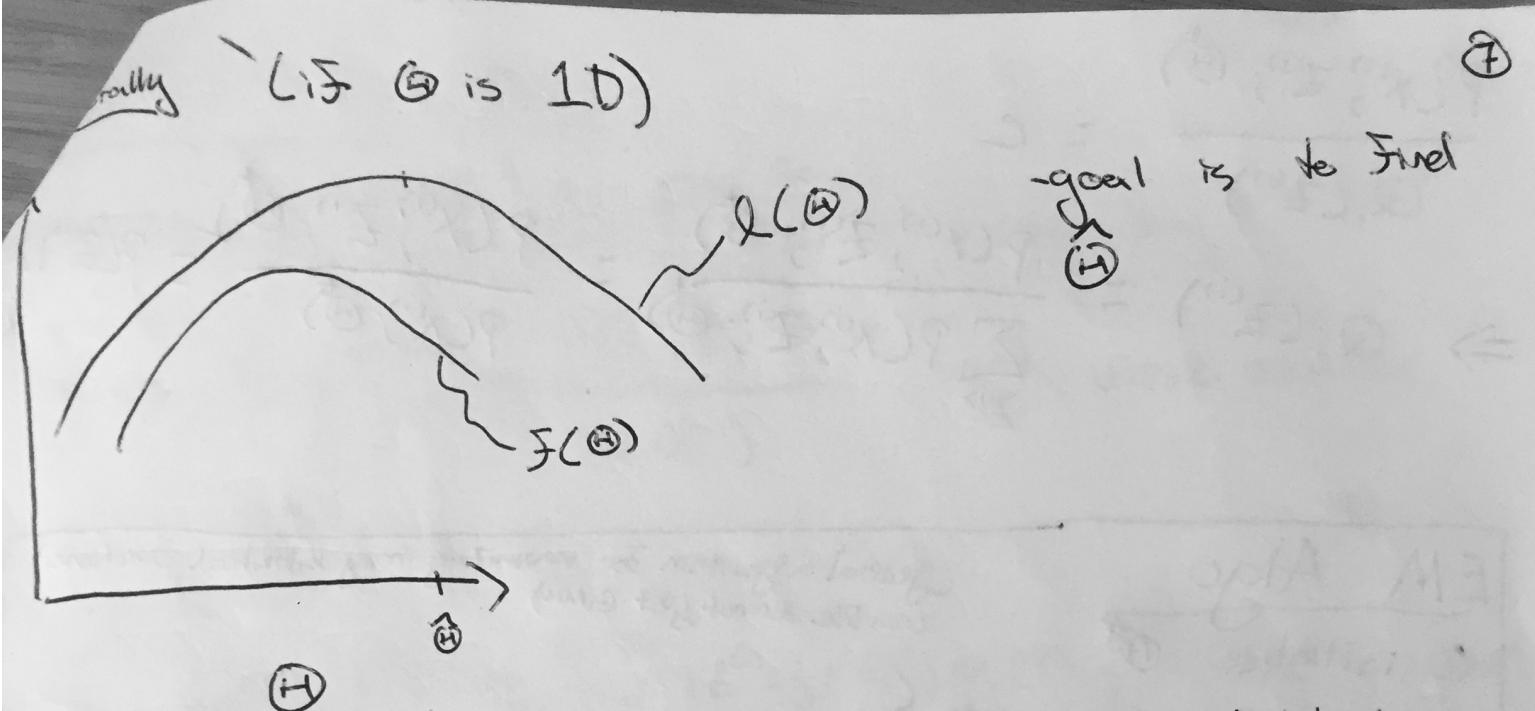
$$\geq \sum_{Z^{(i)}=1}^{\infty} Q_m(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_m(Z^{(i)})}$$

$$\Rightarrow \geq \sum_{Z^{(i)}=0}^{\infty} Q_i(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})} + \dots + \sum_{Z^{(i)}=0}^{\infty} Q_m(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_m(Z^{(i)})}$$

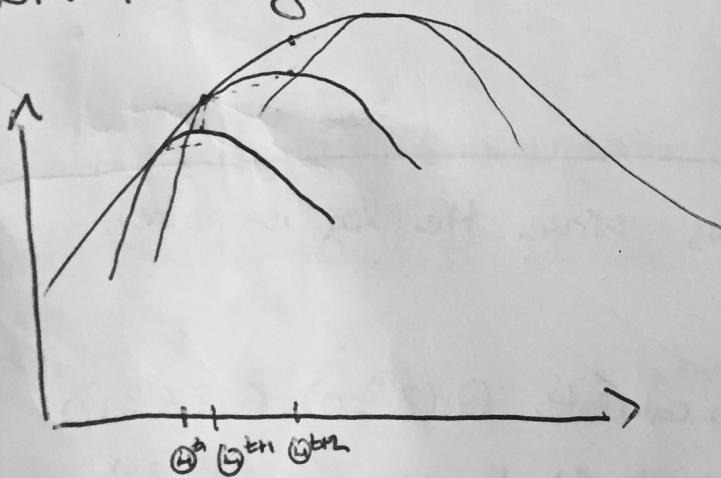
From Jensen's inequality

$$= \sum_{i=1}^m \sum_{Z^{(i)}=0}^{Q_i} Q_i(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})} \quad \} = f(\Theta)$$

(notice that we have been able to push the log inside the sum, but at the expense of an inequality sign)



we still need to specify all the  $Q_i$ 's. To get tightest lower bound on  $l(H)$  lets construct the  $Q_i$ 's to hold w/ equality at a current "best" guess of  $\hat{H}, \hat{H}^t$ . This will naturally lead to a method to find  $\hat{H}$ :



- construct  $Q_i(Z^0; \hat{H}^t)$
- maximize  $\mathcal{J}_0$  over  $\hat{H}$  ~~and~~
- construct  $Q_i(Z^0; \hat{H}^{t+1})$

### Constructing the $Q_i$ 's:

- Remember ~~Jens~~ Jensen's inequality only holds with equality if the RV is a constant



$$\frac{P(X^{(i)}, Z^{(i)}; \Theta)}{Q_i(Z^{(i)})} = C$$

$$\Rightarrow Q_i(Z^{(i)}) = \frac{P(X^{(i)}, Z^{(i)}; \Theta)}{\sum_{Z^{(i)}} P(X^{(i)}, Z^{(i)}; \Theta)} = \frac{P(X^{(i)}, Z^{(i)}; \Theta)}{P(X^{(i)}; \Theta)} = P(Z^{(i)})$$

## EM Algo

(General algorithm for maximizing many likelihood functions)  
 (can use for not just GMM)

① initialize  $\Theta^*$

② repeat until convergence {

$$(\text{E-step}) \quad 2a \quad Q_i(Z^{(i)}) := P(Z^{(i)} | X^{(i)}; \Theta^*)$$

$$(\text{M-step}) \quad 2b \quad \Theta^* := \underset{\Theta}{\operatorname{argmax}} \sum_i \sum_{Z^{(i)}} Q_i(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}; \Theta)}{Q_i(Z^{(i)})}$$

- Note 2b is now tractable w/ calculus since the log is pushed inside the sum

- Note that for GMM, we need to calculate  $Q_i(Z^{(i)}=0), Q_i(Z^{(i)}=1)$  for each data point, then plug this into 2b to sum over the  $Z^{(i)}$ 's.

## EM applied to GMM

- need to find  $\hat{\phi}^*, \hat{\mu}_1^*, \hat{\mu}_2^*, \hat{\Sigma}_1^*, \hat{\Sigma}_2^*$  by maximizing

$$\sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}; z^{(i)}; \Theta)}{Q_i(z^{(i)})} \quad \text{w.r.t. these variables.}$$

$$= \sum_{i=1}^m \sum_{j=1}^o Q_i(z^{(i)}) \log \cancel{p(x^{(i)}|z^{(i)}=j; \Theta)} p(z^{(i)}=j) - Q_i(z^{(i)}=j) \log Q_i(z^{(i)}=j)$$

$$= \left\{ \sum_{i=1}^m \sum_{j=1}^o Q_i(z^{(i)}=j) \log \left\{ \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp \left\{ \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\} \right\} \times \phi_j^{\prod_{k \neq j} w_k^{(i)}} (1-\phi_j)^{\prod_{k \neq j} (1-w_k^{(i)})} \right\} - Q_i \log Q_i(z^{(i)}=j)$$

take derivatives of this expression w/  $\frac{\partial}{\partial \phi}$ ,  $\nabla \mu_e$ ,  $\nabla \Sigma_e$ , and

set the expressions to 0, and solve to find:

$$\hat{\phi}^* = \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \quad \text{where } w_j^{(i)} = Q_i(z^{(i)}=j)$$

$$\hat{\mu}_j^* = \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}} \quad \begin{array}{l} \text{evaluated at old values} \\ \text{of } \Theta \end{array}$$

$$\hat{\Sigma}_j^* = \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

- notice that if we replace  $w_j^{(0)}$  w/  $\mathbb{1}\{Z^{(0)}=j\}$  we <sup>(10)</sup> get the sample  $\phi$ , sample  $\mu_{jS}$  and sample  $\Sigma_{jS}$ .

- thus the MML parameters closely resemble the sample estimates but w/ the hard  $(0,1)$   $\mathbb{1}$  indicator function replaced w/ a soft  $w_j^{(0)} \in [0,1]$  probability w/  $w_j^{(0)}$  being close to 1 if ~~it is more likely to be in~~ cluster 1, and closer to 0 if it is less likely to be in cluster 1.