

EM Algorithm (Cont'd)

RECAP

- Generic EM Algorithm
- Apply to GMMs (MESSAGE: Gmm algo is EM)
- START FACTOR ANALYSIS

WE EXAMINED GMMs

GIVEN: $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$ AND $k > 0$

DO: find $P(z^{(i)} = j)$ for $i = 1 \dots n, j = 1 \dots k$

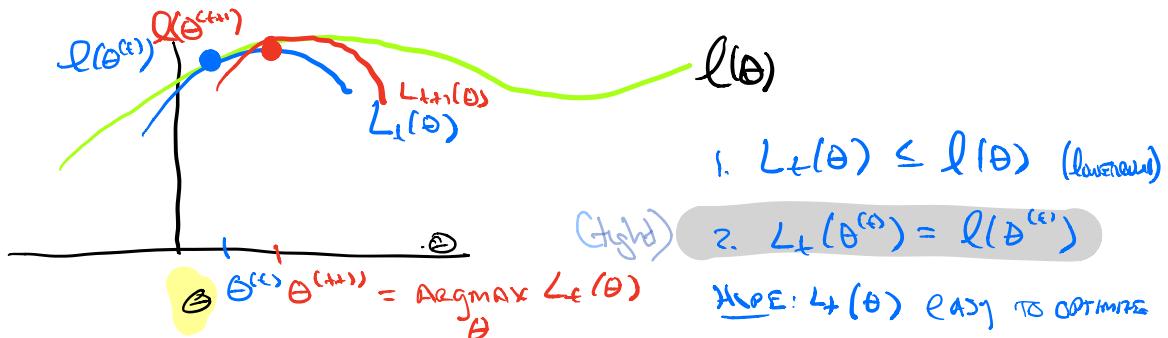
"A SOFT Assignment"

WE CALLED $z^{(i)}$ A LATENT VARIABLE NOT OBSERVED DIRECTLY

TODAY: Generalize model, relate it to MLE, EM Algorithm

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^n \log P(x^{(i)}; \theta) \quad \text{PARAMETERS} \\ &= \sum_{i=1}^n \log \sum_z P(x^{(i)}, z; \theta)\end{aligned}$$

Picture of Generic Algorithm



Rough Algorithm

(E-STEP) 1. Find $L_t(\theta)$ given $\theta^{(t)}$

$$(\text{m-step}) \quad 2. \quad \theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \quad L_t(\theta)$$

WE go term $\log \sum_z P(x^{(i)}, z; \theta)$ term fix a particular i

$$\log \sum_z P(x^{(i)}, z; \theta) = \log \sum_z \frac{Q^{(i)}(z)}{Q^{(i)}(z)} P(x^{(i)}, z; \theta)$$

Pick $Q^{(i)}(z)$ s.t. $\sum_z Q^{(i)}(z) = 1 \quad \frac{1}{z} \quad Q^{(i)}(z) \geq 0 \quad \forall z$. by

$$= \log \mathbb{E}_{z \sim Q^{(i)}} \left[\frac{P(x^{(i)}, z; \theta)}{Q^{(i)}(z)} \right]$$

Recall: $\log(\mathbb{E}(x)) \geq \mathbb{E}[\log(x)]$. Great.

ns: this holds
for any choice $Q^{(i)}$

that satisfies (x)

$$\geq \mathbb{E}_{z \sim Q^{(i)}} \left[\log \frac{P(x^{(i)}, z; \theta)}{Q^{(i)}(z)} \right]$$

$$= \sum_z Q^{(i)}(z) \log \frac{P(x^{(i)}, z; \theta)}{Q^{(i)}(z)}$$

Property 1

Pick $Q^{(i)}(z)$ satisfy Property 2?

$$L(\theta^{(i)}) = L_t(\theta^{(i)}) \quad \text{"tight"}$$

want $\log \sum_z P(x^{(i)}, z; \theta) = \sum_z Q^{(i)}(z) \log \frac{P(x^{(i)}, z; \theta)}{Q^{(i)}(z)}$

TAKE: $Q^{(i)}(z) = P(z | x^{(i)}; \theta)$

- See piazza post 822 for more discussion.

+ $P(x^{(i)}, z; \theta) = P(x^{(i)}; \theta)$ does not depend

$$\bullet = \log P(x^{(i)}; \theta) - \sum_z Q(z^{(i)}) \log \frac{P(z|x^{(i)}; \theta)}{Q(z^{(i)})}$$

NOTE $Q^{(i)}(z)$ varies for every point.

WE CALL ELBO($x, Q; \theta$) = $\sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)}$

WE'VE SHOWN $\ell(\theta) \geq \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta)$

$$\ell(\theta^{(t)}) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta^{(t)})$$

\uparrow with our iteration of θ^t

WARM UP: Mixture of Gaussians. EM RECOVERS AD HOC ALGORITHM

$$P(x^{(i)}, z^{(i)}) = P(x^{(i)}|z^{(i)}) P(z^{(i)})$$

"In clusters" $z^{(i)} \sim \text{Multinomial}(\Phi)$ $\Phi_i \geq 0, \sum \Phi_i = 1$

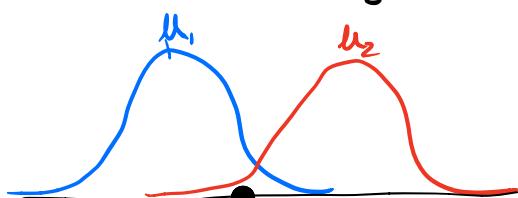
"cluster means" $x^{(i)}|z^{(i)}=j \sim N(\mu_j, \Sigma_j)$

$z^{(i)}$ LATENT VARIABLE.

WHAT IS EM HERE?

$$\Phi^{(i)}(j) = P(z^{(i)}=j | x^{(i)}; \theta)$$

"GMM" IN TERMS OF
 $\sim P(x^{(i)}|z^{(i)}=j; \theta)$



VIA BAYES RULE

M-STEP

MAX
 ϕ, μ, Σ

$$f_i(\theta) = \sum_{c=1}^n \sum_{z^{(i)}} Q^{(i)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_c(z^{(i)})}$$

$$\omega_j^{(i)} = Q^{(i)}(j)$$

$$f_i(\theta) = \sum_j \omega_j^{(i)} \log \frac{1}{2\pi |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma^{-1} (x^{(i)} - \mu_j) \right\} \phi_j$$

$$\nabla_{\theta_j} f_i(\theta) = \sum_j \omega_j^{(i)} \log \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma^{-1} (x^{(i)} - \mu_j) \right\}$$

derivative $\frac{\partial}{\partial \theta_j} = \sum_j -\frac{1}{2} \omega_j^{(i)} \Sigma^{-1} (x^{(i)} - \mu_j)$

Assume Σ^{-1} exists (full rank)

setting to 0 $\mu_j = \frac{\sum \omega_j^{(i)} x^{(i)}}{\sum \omega_j^{(i)}}$ if 6mm update

$$f_i(\theta) = \sum_j \omega_j^{(i)} \log \frac{1}{2\pi |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma^{-1} (x^{(i)} - \mu_j) \right\} \phi_j$$

$$\nabla_{\phi_j} f_i(\theta) = \sum_j \nabla_{\phi_j} \omega_j^{(i)} \log \phi_j \quad \phi_j \text{ is constrained}$$

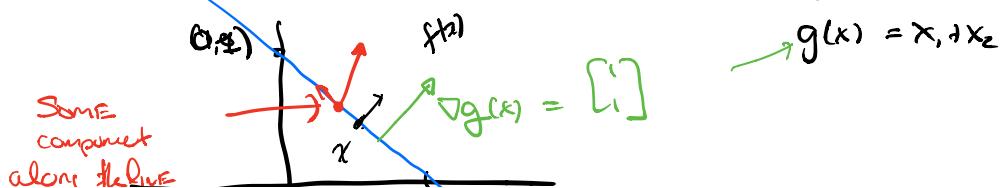
$$\sum \phi_j = 1. \quad \text{constraint}$$

Lagrangian

$$\Rightarrow \nabla_{\phi_j} f_i(\theta) = \sum_j \nabla_{\phi_j} \left(\omega_j^{(i)} \log \phi_j + \lambda \left(\sum_{m=1}^K \phi_m - 1 \right) \right)$$

Detail: want to find critical points of

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad \text{s.t. } x_1 + x_2 = 1 \quad (g(x) = 1)$$



\Rightarrow NOT A CRITICAL POINT

$$(1,0) \quad x_1 + x_2 = 1$$

If z is a CRITICAL POINT then $\nabla f(z)$ is PARALLEL to $\nabla g(z)$. $\nabla f(z) = -\lambda \nabla g(z)$

exists

$$L(x, \lambda) = f(x) + \lambda(g(x) - 1)$$

$$\nabla_x L = \nabla f(x) + \lambda \nabla g(x) = 0 \rightarrow \text{Parallel Condition}$$

$$\nabla_\lambda L = g(x) - 1 = 0 \rightarrow \text{Constraint is satisfied}$$

$$\nabla_{\phi_j} \left(\sum_{i=1}^n w_j^{(i)} \log \phi_j + \lambda \left(\sum_{m=1}^k \phi_m - 1 \right) \right)$$

$$= \sum_{i=1}^n \frac{w_j^{(i)}}{\phi_j} + \lambda = 0 \Rightarrow \phi_j = -\frac{1}{\lambda} \sum_{i=1}^n w_j^{(i)}$$

$$\text{Since } \sum_j \phi_j = 1 \Rightarrow \lambda = \sum_j \phi_j = -\frac{1}{\lambda} \sum_i \sum_j w_j^{(i)} = -\frac{n}{\lambda}$$

$$\Rightarrow \phi_j = \frac{1}{n} \sum_{i=1}^n w_j^{(i)}$$

MESSAGE: EM RECOVERS GMM ALGORITHM

NB: z is discrete here; can replace sums by integrals

FACTOR ANALYSIS

MANY FEWER points than dimensions "n << d"

$n \gg d$ GMMs

How does this happen?

PLACE SENSORS ALL OVER CANVAS RECORD TEMPERATURE

@ 1000s of locations \rightarrow 1000s
 But only record for 30 days ($N < d$)
 Want to fit a density but seems hopeless...

Key Idea Assume there is some latent r.v. that
 is not complex AND explains behavior

1st let's see the problems w/ fitting a single Gaussian..

Given: $x^{(1)} - x^{(n)} \in \mathbb{R}^d$

$\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}$ → seems ok..

↑ smaller \rightarrow more

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

$\text{rank}(\hat{\Sigma}) \leq n < d$ not full rank

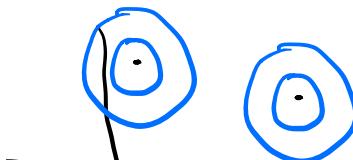
$$P(x; \mu, \Sigma) = \frac{1}{2\pi |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right\}$$

↑ $|\Sigma| = 0$ not defined
↑ zero

We will fix issues by examining three simpler models
 ⇒ SPOILER: Combine all the ~~one~~ final model

Building Block 1

Suppose independent & identical r.v.



COVARIANCE ARE CIRCLES

T

$$\hat{\Sigma} = \sigma^2 I$$

scalar.

$$\begin{aligned}\hat{\Sigma} &\in \mathbb{R}^{d \times d} \\ &= \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix}\end{aligned}$$

$$\min_{\mu, \Sigma} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T \hat{\Sigma}^{-1} (\mathbf{x}^{(i)} - \mu) + \log |\hat{\Sigma}|$$

(Assumption
 $\hat{\Sigma} = \sigma^2 I$)

$$= \sigma^2 \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T (\mathbf{x}^{(i)} - \mu) + \log \sigma^{2d}$$

$$Z = \sigma^2 = 2^{-1} \left(\sum_{i=1}^n \underbrace{\|\mathbf{x}^{(i)} - \mu\|^2}_{C} + d \log Z \right)$$

$$\nabla_Z = -Z^{-2} C + n \frac{d}{Z} = 0 \quad -C = dZ$$

$$Z = \frac{C}{nd}$$

$$\sigma^2 = \frac{C}{nd}$$

$$|Z| = \left| \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix} \right| = \sigma^{2d}$$