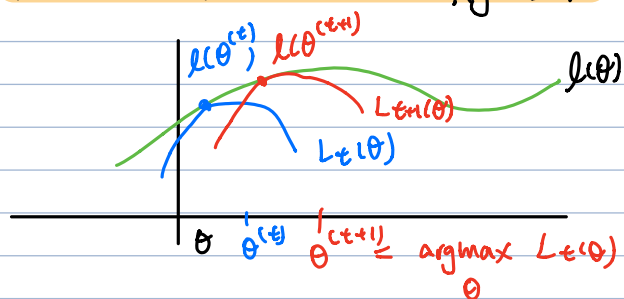


Overview

- General Algo
- EM Algorithm, relate to MLE
- Factor Analysis

Picture of General Algorithm



Property

$$1. L_t(\theta) \leq l(\theta) \quad \text{c (lower bound)}$$

$$2. L_t(\theta^{(t)}) = l(\theta^{(t)}) \quad \text{c (tight)}$$

Hope: $L_t(\theta)$ easy to optimize

Rough Algo

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

(M-STEP) 2. $\theta^{(t+1)} = \arg\max_{\theta} L_t(\theta)$

We go term $\log \sum_z P(x^{(i)}, z; \theta)$ by term

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

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

Any prob distribution for Q

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

Recall: $\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$ (Jensen's)

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

(found lower bound)

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

Property 1

Note this holds for any Q that satisfy $(*)$

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

$$l(\theta^{(t)}) = L_t(\theta^{(t)}) \quad \text{"tight"}$$

(7)

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

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

7 $\frac{P(x^{(i)}, z; \theta)}{P(z | x^{(i)}; \theta)} = P(x^{(i)}; \theta)$ does not depend on z

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

#Note $Q^{(i)}$ 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 $\ell(\theta) \geq \sum_{i=1}^n ELBO(x^{(i)}, Q^{(i)}, \theta)$

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

Warm Up: mixture of Gaussian

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

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

"cluster means" $x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \sigma_j^2)$

$z^{(i)}$: latent variable

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

What is EM here?

$Q^{(i)}_{ij} = P(z^{(i)} = j | x; \theta) \sim P(x^{(i)} | z^{(i)} = j; \theta)$

via Bayes rule

M-step.

$f_i(\theta) = \sum_j Q^{(i)}_{ij} \log \frac{P(x^{(i)}, z^{(i)} = j; \theta)}{Q^{(i)}_{ij}}$
 $w^{(i)}_j = Q^{(i)}_{ij}$

$f_i(\theta) = \sum_j w^{(i)}_j \log \frac{1}{2\pi|\Sigma_j|^{1/2}} \exp\left\{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right\} \Phi_j$

$\nabla_{\mu_j} \sum_i f_i(\theta) = \sum_i \nabla_{\mu_j} w^{(i)}_j \log \exp\left\{-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right\}$

$$= -\frac{1}{2} \sum_j w_j^{(i)} \sum_j^{-1} (x_j^{(i)} - \mu_j) = -\frac{1}{2} \sum_j^{-1} \left(\sum_i w_j^{(i)} (x_j^{(i)} - \mu_j) \right)$$

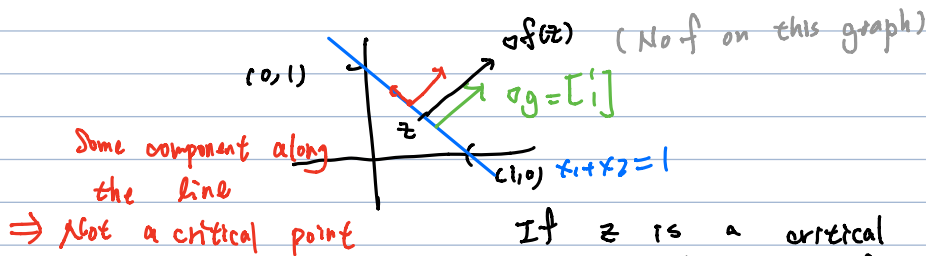
Assume Σ_j^{-1} exists (full rank)

Setting to 0. $\mu_j = \frac{\sum_i w_j^{(i)} x_j^{(i)}}{\sum_i w_j^{(i)}}$

$$\nabla_{\phi_j} f(\theta) = \sum_j \nabla_{\phi_j} w_j^{(i)} \log \phi_j \quad \phi_j \text{ is constrained s.t. } \sum \phi_j = 1$$

$$\Rightarrow \nabla_{\phi_j} f(\theta) = \sum_j \nabla_{\phi_j} (w_j^{(i)} \log \phi_j + \lambda \left(\sum_{m=1}^k \phi_m - 1 \right)) \quad \text{Lagrangian}$$

Detour: Want to find critical points of
 $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ s.t. $x_1 + x_2 = 1$ ($g(x) = 1$, $g(x) = x_1 + x_2$)



If z is a critical point then $\nabla f(z)$ is parallel to $\nabla g(z)$. $\nabla f(z) = -\lambda \nabla g(z)$

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

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

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

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

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

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

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

Message: EM recovers GMM Algorithm

NB: z is discrete, can replace with sum of integrals

Factor Analysis

when " $n \ll d$ " ($n \gg d$, GMMs)

How do this happen?

Place sensors all over the campus, record temp.
But only records for 30 days ($n \ll d$)

Key Idea: Assume there is some latent variable that is not complex and explains behavior

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

Given: $x^{(1)} \dots x^{(n)} \in \mathbb{R}^d$ (small n , large d)

$$\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

n rank 2 vecs

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T \in \mathbb{R}^{d \times d}$$

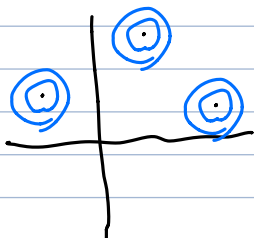
$\text{Rank}(\Sigma) \leq n < d \Rightarrow$ Not full rank

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

\uparrow zero \uparrow Not defined

Building Block 1.

Suppose independent & identical r.v.



Covariance are circles

$$\Sigma = \sigma^2 I$$

\uparrow scalar

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

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

let $z = \sigma^2$

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

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

$$= z^{-1} \sum_{i=1}^n \|x^{(i)} - \mu\|^2 + d \log z$$

$$\nabla_z = -z^{-2} c + n \frac{d}{z} \Rightarrow -c = d z$$

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

$$z = \frac{c}{nd}$$