

# EM and MM Algorithms

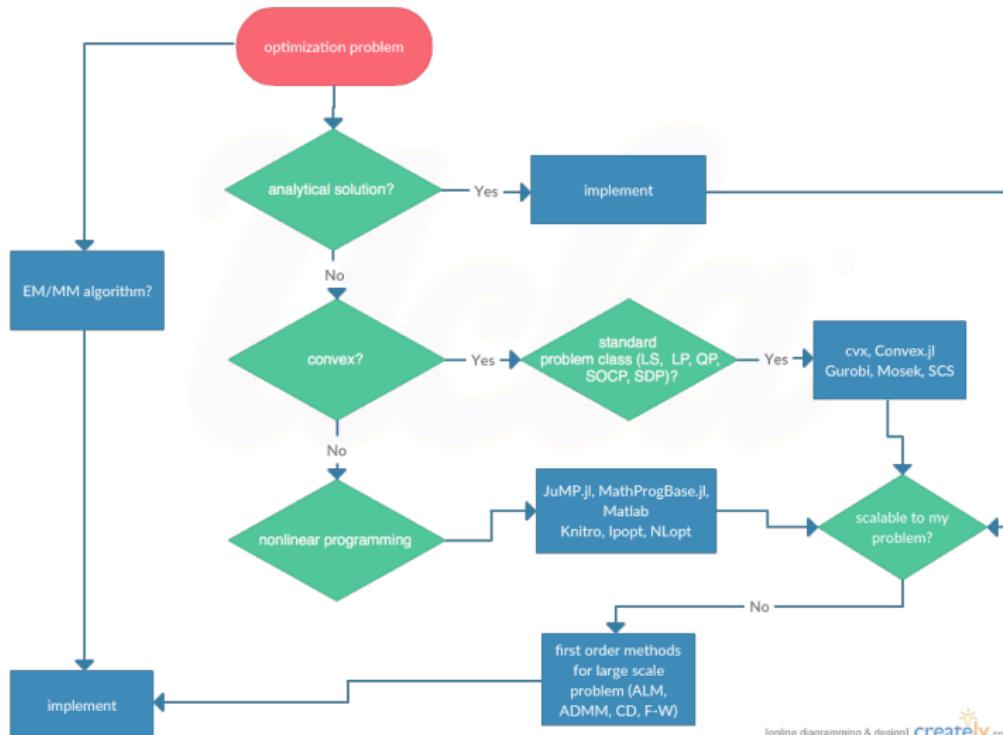
Hua Zhou

Department of Biostatistics  
University of California, Los Angeles

Aug 10, 2016  
SAMSI OPT Summer School

Slides and Jupyter Notebook for the lab session are available at:  
<https://github.com/Hua-Zhou/Public-Talks>

# My flowchart for solving optimization problems



# Outline

## Overview

Problem Setup

Review of Newton's method

## EM algorithm

Introduction

Canonical Example - Finite Mixture Model

In the Zoo of EM

## MM algorithm

Introduction

MM examples

## Acceleration of EM/MM

SQUAREM

Quasi-Newton acceleration

## Problem setup

We consider continuous optimization in this lecture

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \text{constraints on } \mathbf{x}, \end{aligned}$$

where both objective function  $f$  and constraint set are continuous.

I'll switch back and forth between maximization and minimization during lecture (sorry!)

- ▶ Applied mathematicians talk about minimization exclusively
- ▶ Statisticians talk about maximization due to maximum likelihood estimation

They are fundamentally same problem: minimizing  $f(\mathbf{x})$  is equivalent to maximizing  $-f(\mathbf{x})$

## Newton's method

Consider maximizing log-likelihood  $L(\theta)$ ,  $\theta \in \Theta \subset \mathbb{R}^p$ .

- ▶ Newton's method was originally developed for finding roots of nonlinear equations  $f(\mathbf{x}) = \mathbf{0}$
- ▶ Newton's method (aka *Newton-Raphson method*) is considered the **gold standard** for its fast (quadratic) convergence

$$\frac{\|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^*\|}{\|\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^*\|^2} \rightarrow \text{constant}$$

- ▶ Idea: iterative quadratic approximation

- ▶ Notations.  $\nabla L(\theta)$ : gradient vector,  $d^2 L(\theta)$ : Hessian matrix
- ▶ Statistical jargon.  $\nabla L(\theta)$ : score vector,  $-d^2 L(\theta)$ : observed information matrix,  $E[-d^2 L(\theta)]$ : expected (Fisher) information matrix
- ▶ Taylor expansion around the current iterate  $\theta^{(t)}$

$$L(\theta) \approx L(\theta^{(t)}) + \nabla L(\theta^{(t)})^T (\theta - \theta^{(t)}) + \frac{1}{2} (\theta - \theta^{(t)})^T d^2 L(\theta^{(t)}) (\theta - \theta^{(t)})$$

and then maximize the quadratic approximation

- ▶ To maximize the quadratic function, we equate its gradient to zero

$$\nabla L(\theta^{(t)}) + [d^2 L(\theta^{(t)})] (\theta - \theta^{(t)}) = \mathbf{0}_p,$$

which suggests the next iterate

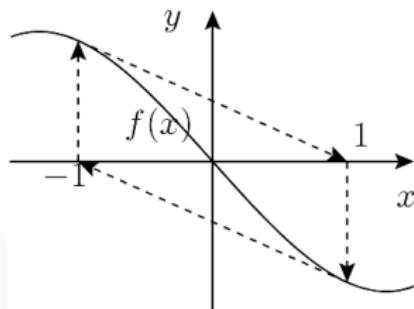
$$\begin{aligned}\theta^{(t+1)} &= \theta^{(t)} - [d^2 L(\theta^{(t)})]^{-1} \nabla L(\theta^{(t)}) \\ &= \theta^{(t)} + [-d^2 L(\theta^{(t)})]^{-1} \nabla L(\theta^{(t)})\end{aligned}$$

## Issues with Newton's method and remedies

- ▶ Need to derive, evaluate, and “invert” the observed information matrix. In statistical problems, often evaluating Hessian costs  $O(np^2)$  flops and inverting it costs  $O(p^3)$  flops

### Remedies:

1. automatic (analytical) differentiation
2. numerical differentiation (works for small problems)
3. quasi-Newton methods (BFGS,  $\ell$ -BFGS probably the most popular black-box nonlinear optimizer)
4. exploit structures whenever possible to reduce the cost of evaluating and inverting the Hessian



- ▶ Stability: Naïve Newton's iterate is not guaranteed to be an ascent algorithm. It's equally happy to head uphill or downhill

### Remedies:

1. approximate  $-d^2 L(\theta^{(t)})$  by a positive definite  $\mathbf{A}$  (if it's not), **and**
2. line search (backtracking)

By first-order Taylor expansion,

$$\begin{aligned}
 & L(\theta^{(t)} + s\Delta\theta^{(t)}) - L(\theta^{(t)}) \\
 = & s\nabla L(\theta^{(t)})^T \Delta\theta^{(t)} + o(s) \\
 = & s\nabla L(\theta^{(t)})^T \mathbf{A}^{-1} \nabla L(\theta^{(t)}) + o(s).
 \end{aligned}$$

For  $s$  sufficiently small, right hand side is strictly positive

In summary, **Newton scheme** iterates according to

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + s[\mathbf{A}^{(t)}]^{-1} \nabla L(\boldsymbol{\theta}^{(t)}) = \boldsymbol{\theta}^{(t)} + s\Delta\boldsymbol{\theta}^{(t)}$$

where  $\mathbf{A}^{(t)}$  is a pd approximation of  $-d^2L(\boldsymbol{\theta}^{(t)})$  and  $s$  is a step length

- ▶ Line search strategy: step-halving ( $s = 1, 1/2, \dots$ ), golden section search, cubic interpolation, Amijo rule, ...  
Newton direction  $\Delta\theta^{(t)}$  only need to be calculated once. Cost of line search mainly lies in objective function evaluation
- ▶ How to approximate  $-d^2L(\theta)$ ? More of an art than science.  
Often requires problem specific analysis.
- ▶ Taking  $A^{(t)} = I$  leads to the method of **steepest ascent**, aka **gradient ascent**

## Fisher's scoring

Fisher's scoring method replaces  $-d^2 L(\theta)$  by the expected (Fisher) information matrix

$$\mathbf{I}(\theta) = \text{E}[-d^2 L(\theta)] = \text{E}[\nabla L(\theta) \nabla L(\theta)^T] \succeq \mathbf{0}_{p \times p},$$

which is psd under exchangeability of expectation and differentiation.  
Therefore the Fisher's scoring algorithm iterates according to

$$\boxed{\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + s[\mathbf{I}(\boldsymbol{\theta}^{(t)})]^{-1} \nabla L(\boldsymbol{\theta}^{(t)})}$$

## Constrained optimization

Main strategies for dealing with constrained optimization, e.g.,

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && c_i(\mathbf{x}) \geq 0, i = 1, \dots, m \end{aligned}$$

- ▶ Interior point method (barrier method)

$$\text{minimize } f(\mathbf{x}) - \mu \sum_{i=1}^m \log(c_i(\mathbf{x}))$$

and push  $\mu$  to 0

- ▶ Sequential quadratic programming (SQP)

$$\begin{aligned} & \text{minimize} && \text{quadratic approximation of } f \\ & \text{subject to} && c_i(\mathbf{x}) \geq 0 \end{aligned}$$

- ▶ Active set method

Interior point method is the most successful for large problems

# Outline

## Overview

Problem Setup

Review of Newton's method

## EM algorithm

Introduction

Canonical Example - Finite Mixture Model

In the Zoo of EM

## MM algorithm

Introduction

MM examples

## Acceleration of EM/MM

SQUAREM

Quasi-Newton acceleration

# EM algorithm

- ▶ Which are the most cited statistical papers?

Paper	Citations	Per Year
Kaplan-Meier (Kaplan and Meier, 1958)	46886	808
EM (Dempster et al., 1977)	44050	1129
Cox model (Cox, 1972)	40920	930
Metropolis (Metropolis et al., 1953)	31284	497
FDR (Benjamini and Hochberg, 1995)	30975	1450
Unit root test (Dickey and Fuller, 1979)	18259	493
Lasso (Tibshirani, 1996)	15306	765
bootstrap (Efron, 1979)	12992	351
FFT (Cooley and Tukey, 1965)	11319	222
Gibbs sampler (Gelfand and Smith, 1990)	6531	251

Citation counts from Google Scholar on Feb 17, 2016.

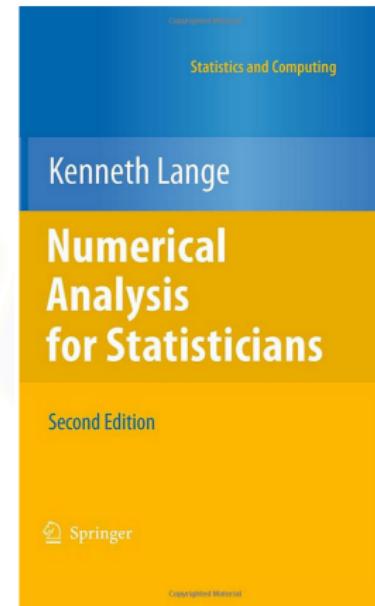
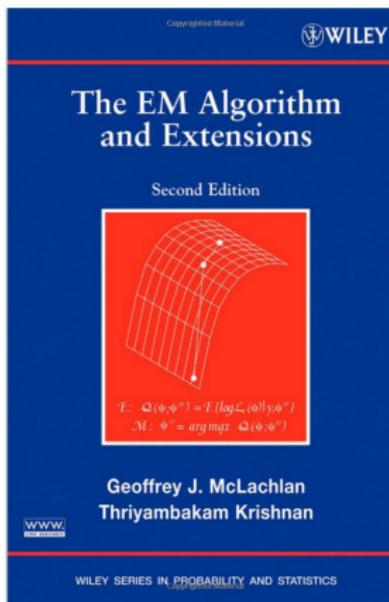
- ▶ EM is one of the most influential statistical ideas, finding applications in various branches of science

EM and MM Algorithms

└ EM algorithm

└ Introduction

## General reference books on EM



▶ Notations

- ▶  $\mathbf{Y}$ : observed data
- ▶  $\mathbf{Z}$ : missing data
- ▶  $\mathbf{X} = (\mathbf{Y}, \mathbf{Z})$ : complete data

- ▶ Goal: maximize the log-likelihood of the observed data  $\ln g(\mathbf{y}|\theta)$  (optimization!)
- ▶ Idea: choose  $\mathbf{Z}$  such that MLE for the complete data is trivial
- ▶ Let  $f(\mathbf{x}|\theta) = f(\mathbf{y}, \mathbf{z}|\theta)$  be the density of complete data
- ▶ Each iteration of EM contains two steps
  - ▶ E step: calculate the conditional expectation

$$Q(\theta|\theta^{(t)}) = \mathbf{E}_{\mathbf{z}|\mathbf{y}=\mathbf{y}, \theta^{(t)}} [\ln f(\mathbf{Y}, \mathbf{Z}|\theta) \mid \mathbf{Y} = \mathbf{y}, \theta^{(t)}]$$

- ▶ M step: maximize  $Q(\theta|\theta^{(t)})$  to generate the next iterate

$$\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta|\theta^{(t)})$$

► Fundamental inequality of EM

$$\begin{aligned}
 & Q(\theta \mid \theta^{(t)}) - \ln g(\mathbf{y} \mid \theta) \\
 = & \mathbf{E}[\ln f(\mathbf{Y}, \mathbf{Z} \mid \theta) \mid \mathbf{Y} = \mathbf{y}, \theta^{(t)}] - \ln g(\mathbf{y} \mid \theta) \\
 = & \mathbf{E}\left\{\ln\left[\frac{f(\mathbf{Y}, \mathbf{Z} \mid \theta)}{g(\mathbf{Y} \mid \theta)}\right] \mid \mathbf{Y} = \mathbf{y}, \theta^{(t)}\right\} \\
 (\text{inform. ineq.}) \leq & \mathbf{E}\left\{\ln\left[\frac{f(\mathbf{Y}, \mathbf{Z} \mid \theta^{(t)})}{g(\mathbf{Y} \mid \theta^{(t)})}\right] \mid \mathbf{Y} = \mathbf{y}, \theta^{(t)}\right\} \\
 = & Q(\theta^{(t)} \mid \theta^{(t)}) - \ln g(\mathbf{y} \mid \theta^{(t)})
 \end{aligned}$$

Rearranging shows

$$\ln g(\mathbf{y} \mid \theta) \geq Q(\theta \mid \theta^{(t)}) - Q(\theta^{(t)} \mid \theta^{(t)}) + \ln g(\mathbf{y} \mid \theta^{(t)})$$

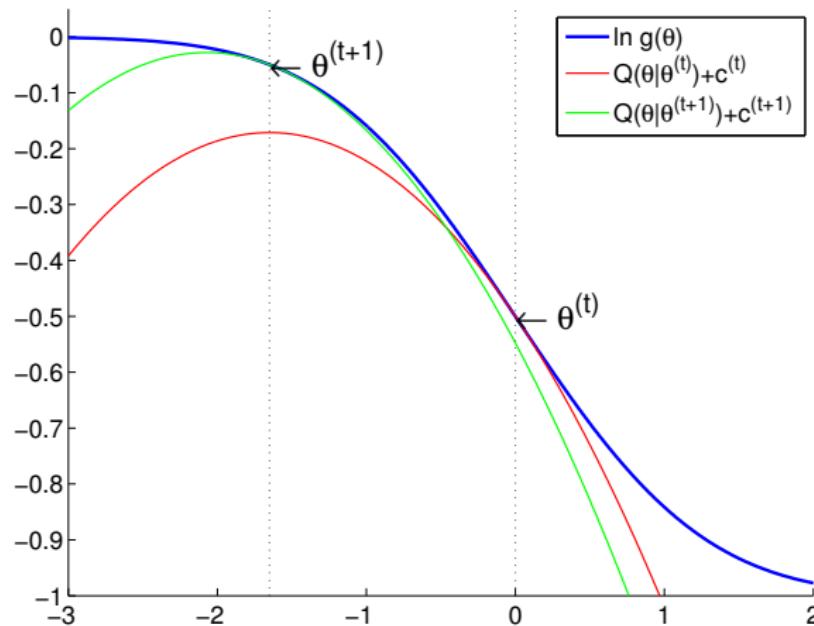
for all  $\theta$ . Equality is achieved at  $\theta = \theta^{(t)}$

► Ascent property of EM

$$\begin{aligned}\ln g(\mathbf{y} \mid \boldsymbol{\theta}^{(t+1)}) &\geq Q(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - Q(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) + \ln g(\mathbf{y} \mid \boldsymbol{\theta}^{(t)}) \\ &\geq \ln g(\mathbf{y} \mid \boldsymbol{\theta}^{(t)})\end{aligned}$$

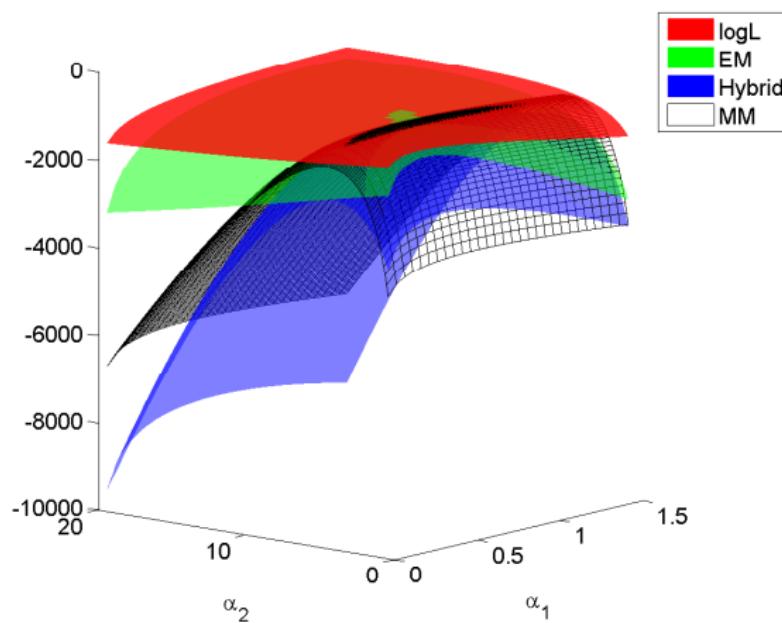
- Obviously we only need  $Q(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - Q(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) \geq 0$  for this ascent property to hold (Generalized EM)

## A picture to keep in mind



# One more ...

2D function



## Remarks on EM

- ▶ Work magically (at least to applied mathematicians). No need for gradient and Hessian
- ▶ Concavity of the observed log-likelihood function (objective function) is **not** required
- ▶ Under regularity conditions, EM converges to a stationary point of  $\ln g(\mathbf{y}|\boldsymbol{\theta})$ . For a non-concave objective function, it can be a local mode or even a saddle point
- ▶ Original idea also appeared in the HMM paper (Baum-Welch algorithm) (Baum et al., 1970)
- ▶ Tend to separate variables – parallel computing

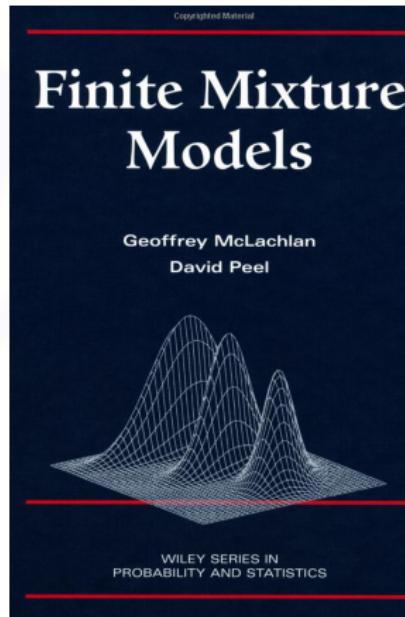
EM and MM Algorithms

└ EM algorithm

└ Canonical Example - Finite Mixture Model

# A canonical EM example – finite mixture model

A canonical example for EM



- ▶ Each data point  $\mathbf{y}_i \in \mathbb{R}^d$  comes from candidate distribution  $h_j$  with probability  $\pi_j$ ,  $\sum_{j=1}^k \pi_j = 1$
- ▶ Mixture density  $h(\mathbf{y}) = \sum_{j=1}^k \pi_j h_j(\mathbf{y} \mid \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j)$ ,  $\mathbf{y} \in \mathbb{R}^d$ , where

$$h_j(\mathbf{y} \mid \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j) = \left( \frac{1}{2\pi} \right)^{d/2} |\det(\boldsymbol{\Omega}_j)|^{-1/2} e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\mu}_j)^T \boldsymbol{\Omega}_j^{-1} (\mathbf{y}-\boldsymbol{\mu}_j)}$$

are  $N(\boldsymbol{\mu}_j, \boldsymbol{\Omega}_j)$

- ▶ Given data  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , want to estimate

$$\boldsymbol{\theta} = (\pi_1, \dots, \pi_k, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Omega}_1, \dots, \boldsymbol{\Omega}_k).$$

Observed (incomplete) data log-likelihood is

$$\ln g(\mathbf{y}_1, \dots, \mathbf{y}_n | \boldsymbol{\theta}) = \sum_{i=1}^n \ln h(\mathbf{y}_i) = \sum_{i=1}^n \ln \sum_{j=1}^k \pi_j h_j(\mathbf{y}_i \mid \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j)$$

## EM for fitting mixture model

- ▶ **Missing data:** let  $z_{ij} = I\{\mathbf{y}_i \text{ comes from group } j\}$
- ▶ Complete data likelihood is

$$f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{j=1}^k [\pi_j h_j(\mathbf{y}_i | \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j)]^{z_{ij}}$$

and thus complete log-likelihood is

$$\ln f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} [\ln \pi_j + \ln h_j(\mathbf{y}_i | \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j)]$$

# E step

- ▶ Conditional expectation

$$Q(\theta|\theta^{(t)}) = \mathbf{E} \left\{ \sum_{i=1}^n \sum_{j=1}^k z_{ij} [\ln \pi_j + \ln h_j(\mathbf{y}_i | \boldsymbol{\mu}_j, \boldsymbol{\Omega}_j) \mid \mathbf{Y} = \mathbf{y}, \boldsymbol{\pi}^{(t)}, \right. \\ \left. \boldsymbol{\mu}_1^{(t)}, \dots, \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Omega}_1^{(t)}, \dots, \boldsymbol{\Omega}_k^{(t)}] \right\}$$

- ▶ By Bayes rule, we have

$$w_{ij}^{(t)} := \mathbf{E}[z_{ij} \mid \mathbf{y}, \boldsymbol{\pi}^{(t)}, \boldsymbol{\mu}_1^{(t)}, \dots, \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Omega}_1^{(t)}, \dots, \boldsymbol{\Omega}_k^{(t)}] \\ = \frac{\pi_j^{(t)} h_j(\mathbf{y}_i | \boldsymbol{\mu}_j^{(t)}, \boldsymbol{\Omega}_j^{(t)})}{\sum_{j'=1}^k \pi_{j'}^{(t)} h_{j'}(\mathbf{y}_i | \boldsymbol{\mu}_{j'}^{(t)}, \boldsymbol{\Omega}_{j'}^{(t)})}$$

- ▶ Q function becomes

$$\sum_{i=1}^n \sum_{j=1}^k w_{ij}^{(t)} \ln \pi_j + \sum_{i=1}^n \sum_{j=1}^k w_{ij}^{(t)} \left[ -\frac{1}{2} \ln \det \boldsymbol{\Omega}_j - \frac{1}{2} (\mathbf{y}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Omega}_j^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) \right]$$

## M step

- ▶ Maximizer of the Q function gives the next iterate

$$\pi_j^{(t+1)} = \frac{\sum_{i=1}^n w_{ij}^{(t)}}{n}$$

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^n w_{ij}^{(t)} \mathbf{y}_i}{\sum_{i=1}^n w_{ij}^{(t)}}$$

$$\Omega_j^{(t+1)} = \frac{\sum_{i=1}^n w_{ij}^{(t)} (\mathbf{y}_i - \mu_j^{(t+1)}) (\mathbf{y}_i - \mu_j^{(t+1)})^\top}{\sum_i w_{ij}^{(t)}}$$

- ▶ Why? Multinomial MLE and weighted multivariate normal MLE

## Remarks

- ▶ No gradient and Hessian calculation at all
- ▶ Simplex constraint on  $\pi$  and psd constraint on  $\Omega_j$  are satisfied by EM iterates
- ▶ The “membership variables”  $w_{ij}^{(t)}$  can be computed in parallel.  
Suchard et al. (2010) observe > 100 speed up when fitting massive mixture model (huge  $n$ ) on GPU
- ▶ It is a common theme EM/MM type algorithms are amenable to massively parallel computing
- ▶ Similar EM applies to other component distributions, as far as its MLE is easy to compute

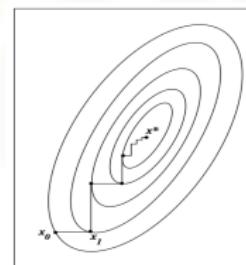
## Other classical EM algorithms

- ▶ Hidden Markov model (Baum-Welch, or forward-backward algorithm)
- ▶ Factor analysis
- ▶ Variance component model
- ▶ Hyper-parameter estimation in empirical Bayes procedure. MAP (maximum *a posteriori* estimation)
- ▶ Missing data
- ▶ Grouped/censorized/truncated model
- ▶ Robust regression ( $t$ -distribution and  $t$ -regression)
- ▶ ...

# ECM

Expectation Conditional Maximization (Meng and Rubin, 1993)

- ▶ In some problems the M step is difficult (no analytic solution)
- ▶ Conditional maximization is easy (block ascent)
  - ▶ partition parameter vector into blocks  $\theta = (\theta_1, \dots, \theta_B)$
  - ▶ alternate update  $\theta_b, b = 1, \dots, B$



- ▶ Ascent property still holds. Why?
- ▶ ECM may converge slower than EM (more iterations) but the total computer time may be shorter due to ease of the CM step

## ECME

### ECM Either (Liu and Rubin, 1994)

- ▶ Each CM step maximizes either the  $Q$  function or the original incomplete observed log-likelihood
- ▶ Ascent property still holds. Why?
- ▶ Faster convergence than ECM

## AECM

Alternating ECM (Meng and van Dyk, 1997)

- ▶ The specification of the complete-data is allowed to be different on each CM-step
- ▶ Ascent property still holds. Why?

## PX-EM and efficient data augmentation

- ▶ Parameter-EXpanded EM (Liu et al., 1998; Liu and Wu, 1999)
- ▶ Efficient data augmentation (Meng and van Dyk, 1997)
- ▶ Idea: Speed up the convergence of EM algorithm by efficiently augmenting the observed data (introducing a working parameter in the specification of complete data)

## Example: $t$ distribution

- ▶  $\mathbf{W} \in \mathbb{R}^p$  is a multivariate  $t$ -distribution  $t_p(\mu, \Sigma, \nu)$  if  
 $\mathbf{W} \sim \text{Normal}(\mu, \Sigma/u)$  and  $u \sim \text{Gamma}(\frac{\nu}{2}, \frac{\nu}{2})$
- ▶ Widely used for robust modeling
- ▶ Gamma( $\alpha, \beta$ ) has density

$$f(u | \alpha, \beta) = \frac{\beta^\alpha u^{\alpha-1}}{\Gamma(\alpha)} e^{-\beta u}, \quad u \geq 0,$$

with mean  $\mathbf{E}(U) = \alpha/\beta$  and  $\mathbf{E}(\ln U) = \psi(\alpha) - \ln \beta$

- ▶ Density of  $\mathbf{W}$  is

$$f_p(\mathbf{w} | \mu, \Sigma, \nu) = \frac{\Gamma(\frac{\nu+p}{2}) |\Sigma|^{-1/2}}{(\pi\nu)^{p/2} \Gamma(\nu/2) [1 + \delta(\mathbf{w}, \mu; \Sigma)/\nu]^{(\nu+p)/2}},$$

where  $\delta(\mathbf{w}, \mu; \Sigma) = (\mathbf{w} - \mu)^T \Sigma^{-1} (\mathbf{w} - \mu)$  is the Mahalanobis squared distance between  $\mathbf{w}$  and  $\mu$

- ▶ Given iid data  $\mathbf{w}_1, \dots, \mathbf{w}_n$ , the log-likelihood is

$$\begin{aligned} L(\mu, \Sigma, \nu) &= -\frac{np}{2} \ln(\pi\nu) + n \left[ \ln \Gamma \left( \frac{\nu+p}{2} \right) - \ln \Gamma \left( \frac{\nu}{2} \right) \right] \\ &\quad - \frac{n}{2} \ln |\Sigma| + \frac{n}{2} (\nu + p) \ln \nu \\ &\quad - \frac{1}{2} (\nu + p) \sum_{j=1}^n \ln[\nu + \delta(\mathbf{w}_j, \mu; \Sigma)] \end{aligned}$$

- ▶ How to compute MLE  $(\hat{\mu}, \hat{\Sigma}, \hat{\nu})$ ?

## Multivariate $t$ MLE – EM

- ▶  $\mathbf{W}_j|u_j$  independent  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma}/u_j)$  and  $U_j$  iid gamma( $\frac{\nu}{2}, \frac{\nu}{2}$ )
- ▶ **Missing data:**  $\mathbf{z} = (u_1, \dots, u_n)^\top$
- ▶ Log-likelihood of complete data is

$$\begin{aligned} L_c(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) &= -\frac{1}{2}np \ln(2\pi) - \frac{n}{2} \ln |\boldsymbol{\Sigma}| \\ &\quad - \frac{1}{2} \sum_{j=1}^n u_j (\mathbf{w}_j - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w}_j - \boldsymbol{\mu}) \\ &\quad - n \ln \Gamma\left(\frac{\nu}{2}\right) + \frac{n\nu}{2} \ln\left(\frac{\nu}{2}\right) \\ &\quad + \frac{\nu}{2} \sum_{j=1}^n (\ln u_j - u_j) - \sum_{j=1}^n \ln u_j \end{aligned}$$

## E step

- ▶ Since gamma is conjugate prior for  $U$ , conditional distribution of  $U$  given  $\mathbf{W} = \mathbf{w}$  is gamma( $(\nu + p)/2, (\nu + \delta(\mathbf{w}, \mu; \Sigma))/2$ ). Thus

$$\mathbf{E}(U_j | \mathbf{w}_j, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}, \nu^{(t)}) = \frac{\nu^{(t)} + p}{\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})} =: u_j^{(t)}$$

$$\mathbf{E}(\ln U_j | \mathbf{w}_j, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}, \nu^{(t)}) = \ln u_j^{(t)} + \left[ \psi\left(\frac{\nu^{(t)} + p}{2}\right) - \ln\left(\frac{\nu^{(t)} + p}{2}\right) \right]$$

- ▶ Overall  $Q$  function (up to an additive constant) takes the form

$$-\frac{n}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{j=1}^n u_j^{(t)} (\mathbf{w}_j - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w}_j - \boldsymbol{\mu})$$

$$-n \ln \Gamma\left(\frac{\nu}{2}\right) + \frac{n\nu}{2} \ln\left(\frac{\nu}{2}\right)$$

$$+ \frac{n\nu}{2} \left[ \frac{1}{n} \sum_{j=1}^n (\ln u_j^{(t)} - u_j^{(t)}) + \psi\left(\frac{\nu^{(t)} + p}{2}\right) - \ln\left(\frac{\nu^{(t)} + p}{2}\right) \right]$$

## M step

- ▶ Maximization over  $(\mu, \Sigma)$  is simply a weighted multivariate normal problem

$$\begin{aligned}\boldsymbol{\mu}^{(t+1)} &= \frac{\sum_{j=1}^n u_j^{(t)} \mathbf{w}_j}{\sum_{j=1}^n u_j^{(t)}} \\ \boldsymbol{\Sigma}^{(t+1)} &= \frac{1}{n} \sum_{j=1}^n u_j^{(t)} (\mathbf{w}_j - \boldsymbol{\mu}^{(t+1)}) (\mathbf{w}_j - \boldsymbol{\mu}^{(t+1)})^\top\end{aligned}$$

Notice down-weighting of outliers is obvious in the update

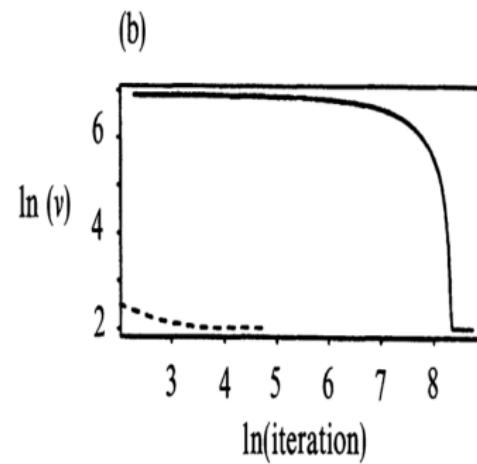
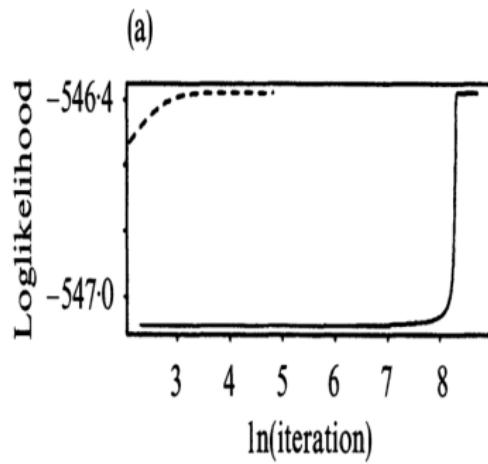
- ▶ Maximization over  $\nu$  is a univariate problem – equating derivative to 0 and find the root

## Multivariate $t$ MLE – ECM and ECME

Partition parameter  $(\mu, \Sigma, \nu)$  into two blocks  $(\mu, \Sigma)$  and  $\nu$

- ▶ ECM = EM for this example. Why?
- ▶ ECME: In the second CM step, maximize over  $\nu$  in terms of the original log-likelihood function instead of the  $Q$  function. They have similar difficulty since both are univariate optimization problems!

An example from (Liu and Rubin, 1994).  $n = 79$ ,  $p = 2$ , with missing entries. EM=ECM: solid line. ECME: dashed line.



## Multivariate $t$ MLE – efficient data augmentation

Assume  $\nu$  known

- ▶ Write  $\mathbf{W}_j = \boldsymbol{\mu} + \mathbf{C}_j / U_j^{1/2}$ , where  $\mathbf{C}_j$  is  $N(\mathbf{0}, \boldsymbol{\Sigma})$  independent of  $U_j$  which is  $\text{Gamma}(\frac{\nu}{2}, \frac{\nu}{2})$
- ▶  $\mathbf{W}_j = \boldsymbol{\mu} + |\boldsymbol{\Sigma}|^{-a/2} \mathbf{C}_j / U_j^{1/2}(a)$ , where  $U_j(a) = |\boldsymbol{\Sigma}|^{-a} U_j$
- ▶ Then the complete data is  $(\mathbf{w}, \mathbf{z}(a))$ ,  $a$  the working parameter
- ▶  $a = 0$  corresponds to the vanilla EM
- ▶ Meng and van Dyk (1997) recommend using  $a_{\text{opt}} = 1/(\nu + p)$  to maximize the convergence rate
- ▶ Exercise: work out the EM update for this special case
- ▶ The only change to the vanilla EM is to replace the denominator  $n$  in the update of  $\boldsymbol{\Sigma}$  by  $\sum_{j=1}^n u_j^{(t)}$
- ▶ PX-EM (Liu et al., 1998) leads to the same update

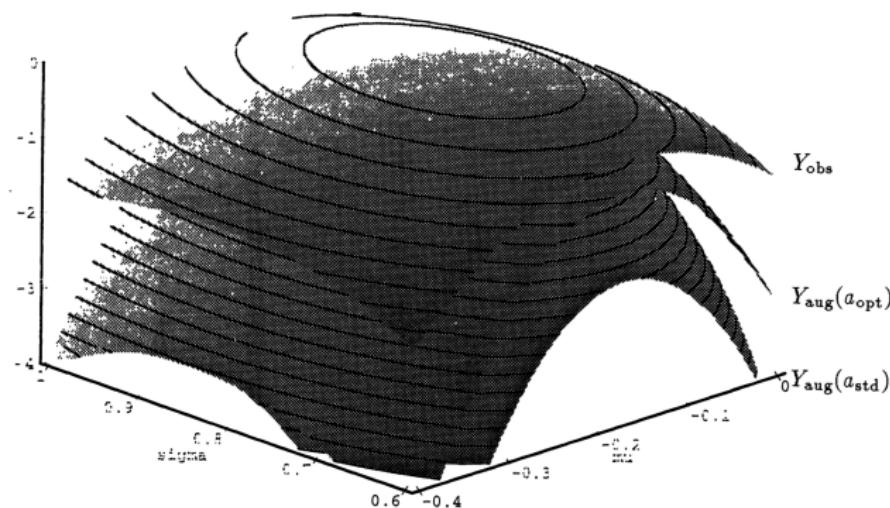


Fig. 5. Comparing the log-likelihoods: the plot shows  $L(\theta|Y_{\text{obs}})$ , as well as  $E[L(\theta|Y_{\text{aug}})|Y_{\text{obs}}, \theta^*]$  for both the standard and the optimal augmentations (each adjusted by their maximum value for comparison); notice that the optimal augmentation results in a flatter log-likelihood that better approximates  $L(\theta|Y_{\text{obs}})$

Assume  $\nu$  unknown

- ▶ Version 1:  $a = a_{\text{opt}}$  in both updating of  $(\mu, \Sigma)$  and  $\nu$
- ▶ Version 2:  $a = a_{\text{opt}}$  for updating  $(\mu, \Sigma)$  and taking the observed data as complete data for updating  $\nu$

Conclusion in (Meng and van Dyk, 1997): Version 1 is  $8 \sim 12$  faster than EM=ECM or ECME. Version 2 is only slightly more efficient than Version 1

## MCEM

### Monte Carlo EM (Wei and Tanner, 1990)

- ▶ Hard to calculate  $Q$  function? Simulate!

$$Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}) \approx \frac{1}{m} \sum_{j=1}^m \ln f(\mathbf{y}, \mathbf{z}_j \mid \boldsymbol{\theta}), \quad \circledcirc$$

where  $\mathbf{z}_j$  are iid from conditional distribution of missing data given  $\mathbf{y}$  and previous iterate  $\boldsymbol{\theta}^{(t)}$

- ▶ Ascent property may be lost due to Monte Carlo errors
- ▶ Example: capture-recapture model (Robert and Casella, 2004, p184), generalized linear mixed model (Booth and Hobert, 1999)

## SEM and SAEM

### Stochastic EM (SEM) (Celeux and Diebolt, 1985)

- ▶ Same as MCEM with  $m = 1$ . A single draw of missing data  $\mathbf{z}$  from the conditional distribution
- ▶  $\theta^{(t)}$  forms a Markov chain which converges to a stationary distribution. No definite relation between this stationary distribution and the MLE
- ▶ In some specific cases, it can be shown that the stationary distribution concentrates around the MLE with a variance inversely proportional to the sample size
- ▶ Advantage: can escape the attraction to inferior mode/saddle point in some mixture model problems

### Simulated Annealing EM (SAEM) (Celeux and Diebolt, 1989)

- ▶ Increase  $m$  with the iterations, ending up with an EM algorithm

# DA

Data Augmentation (DA) algorithm (Tanner and Wong, 1987)

- ▶ Aim for sampling from  $p(\theta|\mathbf{y})$  instead of maximization
- ▶ Idea: incomplete data posterior density is complicated, but the complete-data posterior density is relatively easy to sample
- ▶ Data augmentation algorithm
  - ▶ draw  $\mathbf{z}^{(t+1)}$  conditional on  $(\theta^{(t)}, \mathbf{y})$
  - ▶ draw  $\theta^{(t+1)}$  conditional on  $(\mathbf{z}^{(t+1)}, \mathbf{y})$
- ▶ A special case of Gibbs sampler
- ▶  $\theta^{(t)}$  converges to the distribution  $p(\theta|\mathbf{y})$  under general conditions
- ▶ Ergodic mean converges to the posterior mean  $\mathbf{E}(\theta|\mathbf{y})$ , which may perform better than MLE in finite sample

## EM as a maximization-maximization procedure

Neal and Hinton (1999)

- ▶ Consider the objective function

$$F(\theta, q) = \mathbf{E}_q[\ln f(\mathbf{Z}, \mathbf{y} | \theta)] + H(q)$$

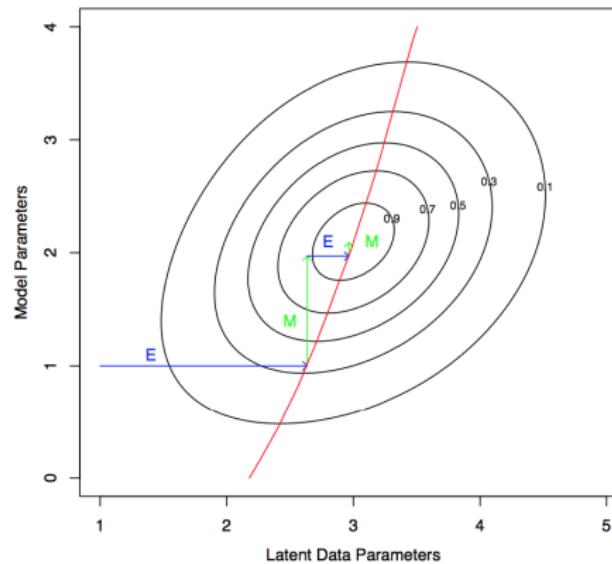
over  $\Theta \times \mathcal{Q}$ , where  $\mathcal{Q}$  is the set of *all* conditional pdfs of the missing data  $\{q(\mathbf{z}) = p(\mathbf{z} | \mathbf{y}, \theta), \theta \in \Theta\}$  and  $H(q) = -\mathbf{E}_q \ln q$  is the entropy of  $q$

- ▶ EM is essentially performing coordinate ascent for maximizing  $F$ 
  - ▶ E step: At current iterate  $\theta^{(t)}$ ,

$$F(\theta^{(t)}, q) = \mathbf{E}_q[\ln p(\mathbf{Z} | \mathbf{y}, \theta^{(t)})] - \mathbf{E}_q \ln q + \ln g(\mathbf{y} | \theta^{(t)}).$$

The maximizing conditional pdf is given by  $q = \ln f(\mathbf{Z} | \mathbf{y}, \theta^{(t)})$

- ▶ M step: Substitute  $q = \ln f(\mathbf{Z} | \mathbf{y}, \theta^{(t)})$  into  $F$  and maximize over  $\theta$



**FIGURE 8.7.** Maximization–maximization view of the EM algorithm. Shown are the contours of the (augmented) observed data log-likelihood  $F(\theta', \tilde{P})$ . The E step is equivalent to maximizing the log-likelihood over the parameters of the latent data distribution. The M step maximizes it over the parameters of the log-likelihood. The red curve corresponds to the observed data log-likelihood, a profile obtained by maximizing  $F(\theta', \tilde{P})$  for each value of  $\theta'$ .

Hastie et al. (2009, Section 8.5.3)

## Incremental EM

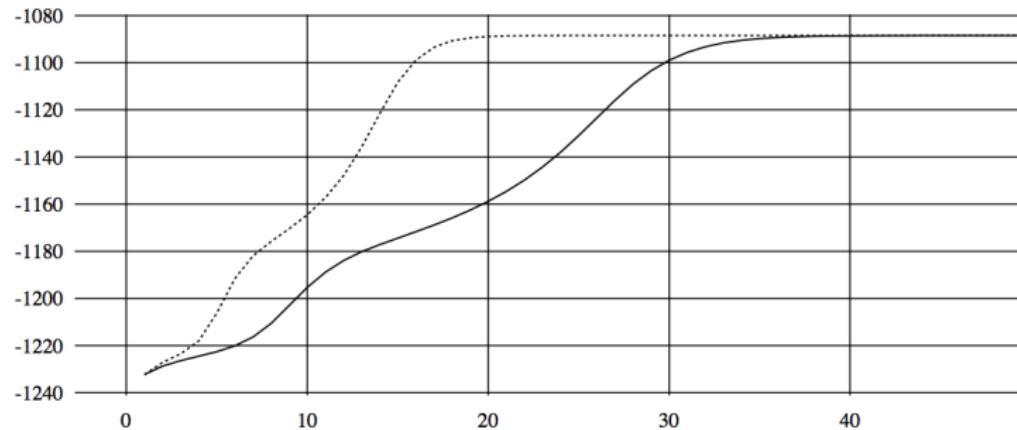
- ▶ Assume iid observations. Then the objective function is

$$F(\theta, q) = \sum_i [\mathbf{E}_{q_i} \ln p(\mathbf{Z}_i, \mathbf{y}_i | \theta) + H(q_i)],$$

where we search for  $q$  under factored form  $q(\mathbf{z}) = \prod_i q_i(\mathbf{z})$

- ▶ Maximizing  $F$  over  $q$  is equivalent to maximizing the contribution of each data with respect to  $q_i$
- ▶ Update  $\theta$  by visiting data items **sequentially** rather than from a global E step
- ▶ Finite mixture example: for each data point  $i$ 
  - ▶ evaluate membership variables  $w_{ij}, j = 1, \dots, k$
  - ▶ update parameter values  $(\pi^{(t+1)}, \mu_j^{(t+1)}, \Sigma_j^{(t+1)})$  (only need to keep sufficient statistics!)

Faster convergence than vanilla EM in some examples



*Figure 1.* Comparison of convergence rates for the standard EM algorithm (solid line) and the incremental algorithm (dotted line). The log likelihood is shown on the vertical axis, the number of passes of the algorithm on the horizontal axis.

# Outline

## Overview

Problem Setup

Review of Newton's method

## EM algorithm

Introduction

Canonical Example - Finite Mixture Model

In the Zoo of EM

## MM algorithm

Introduction

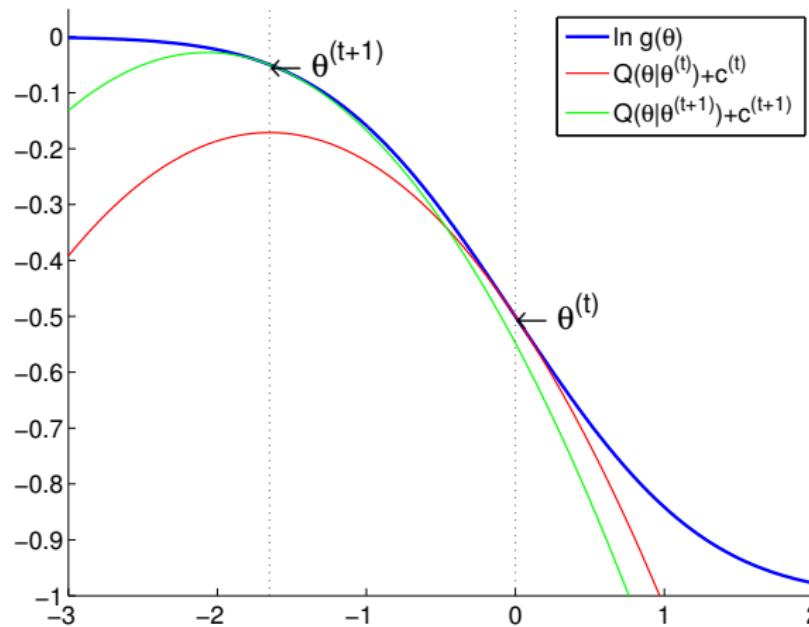
MM examples

## Acceleration of EM/MM

SQUAREM

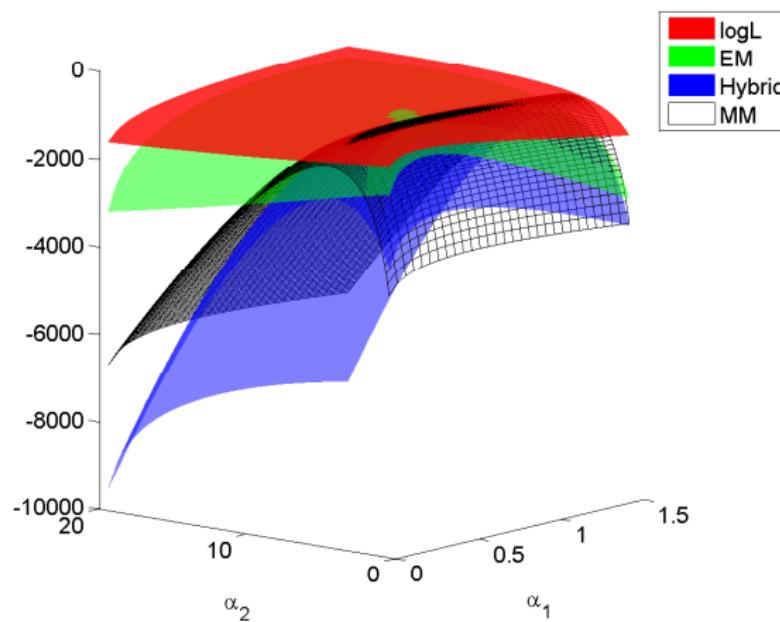
Quasi-Newton acceleration

## Our friend ...



# Our other friend ...

2D function



## EM as a minorization-maximization (MM) algorithm

Refer to the picture in previous slide

- ▶ The  $Q$  function constitutes a **minorizing** function of the objective function up to an additive constant

$$\begin{aligned} L(\theta) &\geq Q(\theta|\theta^{(t)}) + c^{(t)} \quad \text{for all } \theta \\ L(\theta^{(t)}) &= Q(\theta^{(t)}|\theta^{(t)}) + c^{(t)} \end{aligned}$$

- ▶ **Maximizing** the  $Q$  function generates an ascent iterate  $\theta^{(t+1)}$



## Questions:

- ▶ Is there any other way to produce such surrogate function?
- ▶ Can we flip the picture and apply same principle to **minimization** problem?

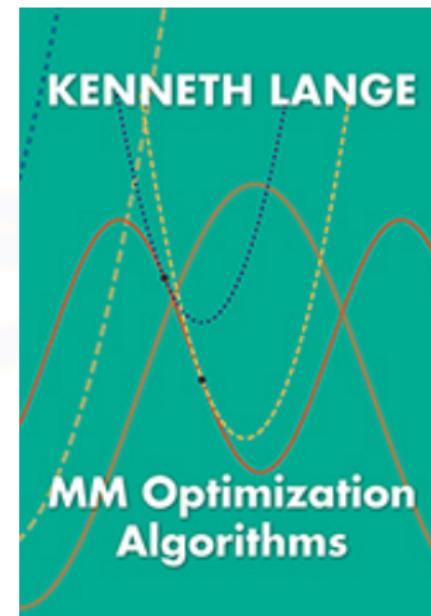
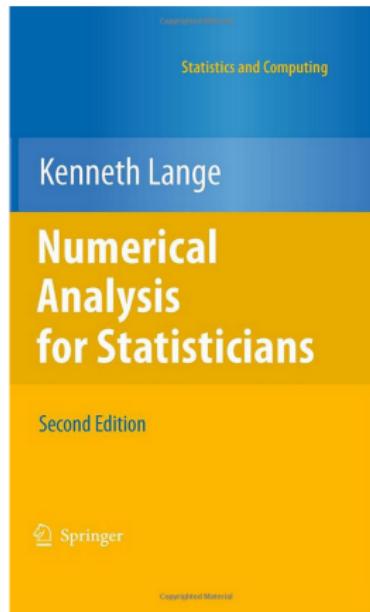
This generalization leads to a powerful tool – MM principle (Lange et al., 2000)

EM and MM Algorithms

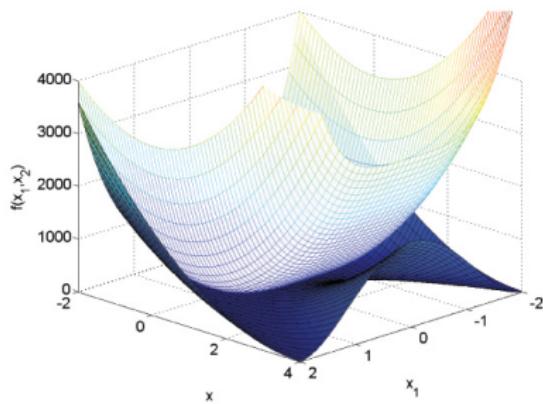
└ MM algorithm

└ Introduction

General reference book on MM



## A powerful tool - MM algorithm



- ▶ A prescription for constructing optimization algorithms
- ▶ **Majorize/Minimize or Minorize/Maximize**
- ▶ Creating a surrogate function that minorizes/majorizes the objective function
- ▶ Optimizing the surrogate function drives the objective function uphill or downhill as needed

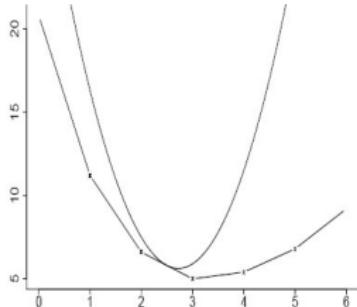
Famous special cases in statistics: EM algorithm (Dempster et al., 1977), multidimensional scaling (de Leeuw and Heiser, 1977)

## Majorization and definition of the algorithm

- ▶ A function  $g(\theta|\theta^n)$  is said to **majorize** the function  $f(\theta)$  at  $\theta^n$  provided

$$f(\theta^n) = g(\theta^n|\theta^n)$$

$$f(\theta) \leq g(\theta|\theta^n) \quad \text{for all } \theta.$$



- ▶ So we choose a “nice/separable” majorizing function  $g(\theta|\theta^n)$  and minimize it. This produces the next point  $\theta^{n+1}$  in the algorithm.
- ▶ The **descent property** follows from

$$f(\theta^{n+1}) \leq g(\theta^{n+1}|\theta^n) \leq g(\theta^n|\theta^n) = f(\theta^n).$$

This makes MM algorithms very stable.

## Rationale for the MM principle

- ▶ Free the derivation from missing data structure
- ▶ Avoid matrix inversion
- ▶ Linearize an optimization problem
- ▶ Deal gracefully with certain equality and inequality constraints
- ▶ Turn a non-differentiable problem into a smooth problem
- ▶ Separate the parameters of a problem (perfect for massive, **fine-scale** parallelization)

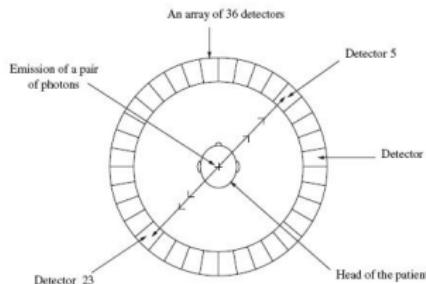
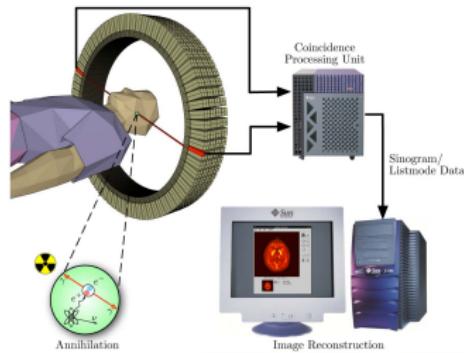


## Generic methods of majorization and minorization

- ▶ Jensen's inequality – EM algorithms
- ▶ The Cauchy-Schwartz inequality - multidimensional scaling
- ▶ Supporting hyperplane property of a convex function
- ▶ Arithmetic-geometric mean inequality
- ▶ Quadratic upper bound principle - Böhning and Lindsay
- ▶ ...

A great source of inequalities: *The Cauchy-Schwarz Master Class* by Michael Steele.

# Positron Emission Tomography (PET)



- ▶ Data: tube readings  $\mathbf{y} = (y_1, \dots, y_d)$
- ▶ Estimate: photon emission intensities (pixels)  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$
- ▶ Poisson Model:  $Y_i \sim \text{Poisson}(\sum_{j=1}^p c_{ij} \lambda_j)$  where  $c_{ij} = \text{cond. prob.}$  that a photon emitted by  $j$ -th pixel is detected by  $i$ -th tube

- ▶ Log-likelihood

$$L(\lambda|\mathbf{y}) = \sum_i \left[ y_i \ln \left( \sum_j c_{ij} \lambda_j \right) - \sum_j c_{ij} \lambda_j \right] + \text{const.}$$

Essentially a Poisson regression with constraint  $\lambda_j \geq 0$

- ▶ Regularized log-likelihood for smoother image:

$$L(\lambda|\mathbf{y}) - \frac{\mu}{2} \sum_{\{j,k\} \in \mathcal{N}} (\lambda_j - \lambda_k)^2$$

- ▶ EM algorithm does not apply now: no likelihood model for the regularization term

## Minorization step

- ▶ By concavity of the  $\ln s$  function

$$\ln \left( \sum_j c_{ij} \lambda_j \right) \geq \sum_j \frac{c_{ij} \lambda_j^{(t)}}{\sum_j c_{ij'} \lambda_{j'}^{(t)}} \ln \left( \frac{\sum_{j'} c_{ij'} \lambda_{j'}^{(t)}}{c_{ij} \lambda_j^{(t)}} c_{ij} \lambda_j \right) = \sum_j \frac{c_{ij} \lambda_j^{(t)}}{\sum_j c_{ij'} \lambda_{j'}^{(t)}} \ln \lambda_j + c^{(t)}$$

- ▶ By concavity of the  $-s^2$  function

$$-(\lambda_j - \lambda_k)^2 \geq -\frac{1}{2}(2\lambda_j - \lambda_j^{(t)} - \lambda_k^{(t)})^2 - \frac{1}{2}(2\lambda_k - \lambda_j^{(t)} - \lambda_k^{(t)})^2$$

- ▶ Combining minorizing terms gives an overall surrogate function

$$\begin{aligned} g(\boldsymbol{\lambda} | \boldsymbol{\lambda}^{(t)}) &= \sum_i y_i \sum_j \frac{c_{ij} \lambda_j^{(t)}}{\sum_{j'} c_{ij'} \lambda_{j'}^{(t)}} \ln \lambda_j - \sum_i \sum_j c_{ij} \lambda_j \\ &\quad - \frac{\mu}{4} \sum_{\{j,k\} \in \mathcal{N}} [(2\lambda_j - \lambda_j^{(t)} - \lambda_k^{(t)})^2 + (2\lambda_k - \lambda_j^{(t)} - \lambda_k^{(t)})^2] \end{aligned}$$

## Maximization step

- ▶  $g(\lambda|\lambda^{(t)})$  is trivial to maximize because  $\lambda_j$  are separated!
- ▶ Solving for the root of

$$\begin{aligned} & \frac{\partial}{\partial \lambda_j} g(\lambda|\lambda^{(t)}) \\ = & \sum_i y_i \frac{c_{ij} \lambda_j^{(t)}}{\sum_{j'} c_{ij} \lambda_{j'}^{(t)}} \lambda_j^{-1} - \sum_i c_{ij} - \mu \sum_{k \in \mathcal{N}_j} (2\lambda_j - \lambda_j^{(t)} - \lambda_k^{(t)}) \\ = & 0 \end{aligned}$$

gives  $\lambda_j^{(t+1)}$

► MM algorithm for PET:

Initialize:  $\lambda_j^{(0)} = 1$

**repeat**

$$z_{ij}^{(t)} = (y_i c_{ij} \lambda_j^{(t)}) / (\sum_k c_{ik} \lambda_k^{(t)})$$

**for**  $j = 1$  to  $p$  **do**

$$a = -2\mu|\mathcal{N}_j|, b = \mu(|\mathcal{N}_j| \lambda_j^{(t)} + \sum_{k \in \mathcal{N}_j} \lambda_k^{(t)}) - 1, c = \sum_i z_{ij}^{(t)}$$

$$\lambda_j^{(t+1)} = (-b - \sqrt{b^2 - 4ac}) / (2a)$$

**end for**

**until** convergence occurs

- Parameter constraints  $\lambda_j \geq 0$  are satisfied when start from positive initial values
- The loop for updating pixels can be carried out independently – **massive parallelism**

## Multivariate $t$ MLE (revisited)

### Problem

- ▶ Given iid data  $\mathbf{w}_1, \dots, \mathbf{w}_n$ , the log-likelihood is

$$\begin{aligned} f(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) &= -\frac{np}{2} \ln(\pi\nu) + n[\ln \Gamma(\frac{\nu+p}{2}) - \ln \Gamma(\frac{\nu}{2})] \\ &\quad - \frac{n}{2} \ln |\boldsymbol{\Sigma}| + \frac{n}{2}(\nu+p) \ln \nu \\ &\quad - \frac{1}{2}(\nu+p) \sum_{j=1}^n \ln[\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})], \end{aligned}$$

where  $\delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma}) = (\mathbf{w}_j - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w}_j - \boldsymbol{\mu})$

- ▶ Want to find MLE  $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}, \hat{\nu})$

## Minorization step

- ▶ Apply the supporting hyperplane inequality

$$\begin{aligned} & -\ln[\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})] \\ \geq & -\frac{\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})}{\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})} + 1 - \ln[\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})] \end{aligned}$$

- ▶ Summing over observations gives an overall minorizing function

$$\begin{aligned} & g(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu | \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}, \nu^{(t)}) \\ = & -\frac{np}{2} \ln(\pi\nu) + n[\ln \Gamma(\frac{\nu+p}{2}) - \ln \Gamma(\frac{\nu}{2})] \\ & -\frac{n}{2} \ln |\boldsymbol{\Sigma}| + \frac{n}{2}(\nu+p) \ln \nu \\ & -\frac{\nu+p}{2} \sum_{j=1}^n \frac{\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})}{\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})} \\ & + \frac{\nu+p}{2} \left\{ n - \sum_{j=1}^n \ln[\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})] \right\} \end{aligned}$$

## Maximization step

Block ascent for maximizing the minorizing function

- ▶ Given  $\nu^{(t)}$ , update  $(\mu, \Sigma)$  by

$$\max_{\mu, \Sigma} -\frac{n}{2} \ln |\Sigma| - \frac{1}{2} \sum_{j=1}^n \frac{\nu^{(t)} + p}{\nu^{(t)} + \delta(\mathbf{w}_j, \mu^{(t)}; \Sigma^{(t)})} \delta(\mathbf{w}_j, \mu; \Sigma)$$

A weighted multivariate normal problem – same update as in the vanilla EM(=ECM)

- ▶ Given  $(\mu^{(t+1)}, \Sigma^{(t+1)})$ , update of  $\nu$  is easy since it is a univariate optimization problem. This gives a different update from the vanilla EM(=ECM)

- ▶ Above derivation gives an **MCM** (minorization conditional maximization) algorithm, which turns out to be different from ECM
- ▶ Wait a minute, why don't we maximize the original objective function directly when updating  $\nu$ ? **MCM Either!**
- ▶ MCME = ECME. Why?

## Yet another MM for $t \dots$

Assume  $\nu$  known

- ▶ Notice

$$\begin{aligned} & -\frac{1}{2}|\boldsymbol{\Sigma}| - \frac{\nu + p}{2} \ln[\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})] \\ &= -\frac{\nu + p}{2} \ln\{|\boldsymbol{\Sigma}|^a [\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})]\} \end{aligned}$$

where  $a = 1/(\nu + p)$  is the working parameter

- ▶ Minorize via

$$\begin{aligned} & -\ln\{|\boldsymbol{\Sigma}|^a [\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})]\} \\ &\geq -\frac{|\boldsymbol{\Sigma}|^a [\nu + \delta(\mathbf{w}_j, \boldsymbol{\mu}; \boldsymbol{\Sigma})]}{|\boldsymbol{\Sigma}^{(t)}|^a [\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})]} + 1 \\ & \quad -\ln\{|\boldsymbol{\Sigma}^{(t)}|^a [\nu^{(t)} + \delta(\mathbf{w}_j, \boldsymbol{\mu}^{(t)}; \boldsymbol{\Sigma}^{(t)})]\} \end{aligned}$$

- ▶ Now working out the maximization step (do it!) yields the efficient data augmentation algorithm in (Meng and van Dyk, 1997)

## Non-negative matrix factorization (NNMF)

Lee and Seung (1999, 2001)

- ▶ **Goal:** Find a low rank  $r$  factorization  $\mathbf{V}_{m \times r} \mathbf{W}_{r \times n}$  of non-negative data matrix  $\mathbf{X}_{m \times n}$ , such that

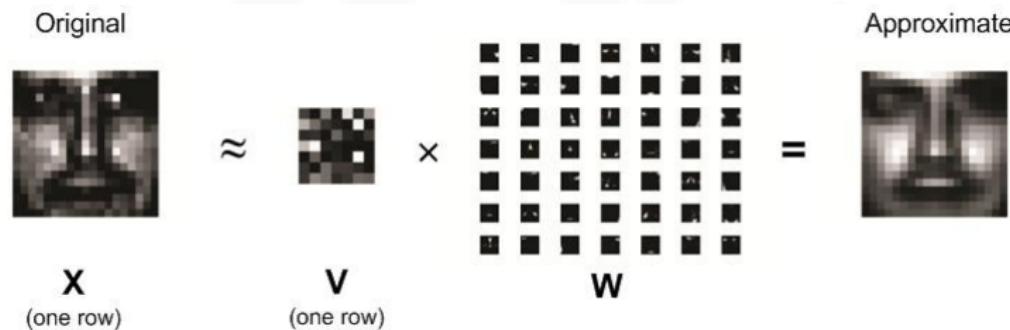
$$\|\mathbf{X} - \mathbf{V}\mathbf{W}\|_F$$

is minimized subject to constraints  $v_{ij}, w_{jk} \geq 0$

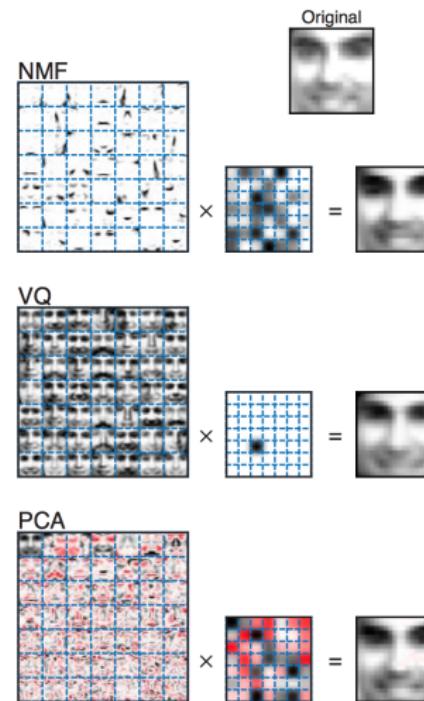
- ▶ Resembles SVD except the non-negativity constraint
- ▶ A popular alternative to PCA and vector quantization
- ▶ Used for data compression, clustering, ...

# Image Factorization

- ▶ MIT CBCL Face Images: 2429 faces, each  $19 \times 19 = 361$  pixels
  - ▶ A striking feature of NNMF basis images is they correspond to different parts of face (eyes, nose, ...)



approximation based on a rank-49 NNMF



**FIGURE 14.33.** Non-negative matrix factorization (NMF), vector quantization (VQ, equivalent to k-means clustering) and principal components analysis (PCA) applied to a database of facial images. Details are given in the text. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

## MM for NMF

- ▶ Objective

$$\|\mathbf{X} - \mathbf{V}\mathbf{W}\|_F = \sum_i \sum_j \left( x_{ij} - \sum_k v_{ik} w_{kj} \right)^2$$

- ▶ Majorization: by convexity of  $s^2$  function

$$\left( x_{ij} - \sum_k v_{ik} w_{kj} \right)^2 \leq \sum_k \frac{v_{ik}^{(t)} w_{kj}^{(t)}}{\sum_l v_{il}^{(t)} w_{lj}^{(t)}} \left( x_{ij} - \frac{\sum_l v_{il}^{(t)} w_{lj}^{(t)}}{\sum_k v_{ik}^{(t)} w_{kj}^{(t)}} v_{ik} w_{kj} \right)^2$$

- ▶ Minimization: alternate update of  $\mathbf{V}$  and  $\mathbf{W}$

## MM Algorithm for NNMF

Initialize:  $\mathbf{V} = \text{rand}(m, r)$  and  $\mathbf{W} = \text{rand}(r, n)$ .

**repeat**

$\mathbf{V}_{..} \leftarrow \mathbf{V}_{..} \times (\mathbf{X}\mathbf{W}^T)_{..} / (\mathbf{V}\mathbf{W}\mathbf{W}^T)_{..}$  for all  $m \times r$  elements

$\mathbf{W}_{..} \leftarrow \mathbf{W}_{..} \times (\mathbf{V}^T\mathbf{X})_{..} / (\mathbf{V}^T\mathbf{V}\mathbf{W})_{..}$  for all  $r \times n$  elements

**until** convergence

- ▶ Updates are extremely simple
- ▶ Non-negativity constraints are satisfied
- ▶ Utilize high throughput of BLAS 3 routines, either on CPU or on GPU
- ▶ Lasso or other sparsity penalty can be incorporated

## Exercise

- ▶ In the original paper, Lee and Seung (1999) pose a Poisson model for the image pixels  $x_{ij} \sim \text{Poisson}(\sum_k v_{ik} w_{kj})$
- ▶ The log-likelihood is

$$L(\mathbf{V}, \mathbf{W}) = \sum_i \sum_j \left[ x_{ij} \ln \left( \sum_k v_{ik} w_{kj} \right) - \sum_k v_{ik} w_{kj} \right]$$

- ▶ MLE under the Poisson model is amenable to another MM algorithm
- ▶ Derive it

# Netflix movie rating data and matrix completion

ID	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6	...	movie 17,770
user 1	5	3	4	3	3	NA	...	1
user 2	4	NA	NA	NA	NA	NA	...	NA
user 3	NA	NA	NA	NA	NA	NA	...	NA
user 4	4	NA	NA	NA	NA	2	...	4
user 5	NA	NA	NA	5	NA	NA	...	NA
user 6	3	NA	NA	5	1	NA	...	3
user 7	NA	NA	NA	NA	NA	NA	...	NA
user 8	5	NA	5	NA	NA	NA	...	NA
user 9	NA	NA	NA	NA	3	NA	...	NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
user 480,189	NA	5	NA	NA	NA	NA	...	NA

- ▶ Snapshot of the kind of data collected by Netflix. Only 100,480,507 ratings (1.2% entries of the matrix) are observed
- ▶ **Goal:** impute the unobserved ratings for personalized recommendation



### Leaderboard

10.05% Display top 20 leaders.

Rank	Team Name	Best Score	% Improvement	Last Submit Time
1	Bellkor's Pragmatic Chaos	0.8558	10.05	2009-06-26 18:42:37
<b>Grand Prize - RMSE &lt;= 0.8563</b>				
2	PragmaticTheory	0.8582		2009-06-25 22:15:51
3	Bellkor in BigChaos	0.8590	9.71	2009-05-13 08:14:09
4	Grand Prize Team	0.8593	9.68	2009-06-12 08:20:24
5	Dace	0.8604	9.56	2009-04-22 05:57:03
6	BigChaos	0.8613	9.47	2009-06-23 23:06:52

- ▶ Observe a very sparse matrix  $\mathbf{Y} = (y_{ij})$ . Want to impute all the missing entries. It is possible when the matrix is structured and of **low rank**
- ▶ Let  $\Omega = \{(i, j) : \text{observed entries}\}$  collect the observed entries and  $P_\Omega(\mathbf{M})$  denote the projection of matrix  $\mathbf{M}$  to  $\Omega$ . The problem

$$\min_{\text{rank}(\mathbf{X}) \leq r} \frac{1}{2} \|P_\Omega(\mathbf{Y}) - P_\Omega(\mathbf{X})\|_F^2 = \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2$$

unfortunately is non-convex and difficult

- ▶ **Convex relaxation** (Mazumder et al., 2010)

$$\min_{\mathbf{X}} f(\mathbf{X}) = \frac{1}{2} \|P_\Omega(\mathbf{Y}) - P_\Omega(\mathbf{X})\|_F^2 + \lambda \|\mathbf{X}\|_*,$$

where  $\|\mathbf{X}\|_* = \|\sigma(\mathbf{X})\|_1 = \sum_i \sigma_i(\mathbf{X})$  is the nuclear norm

# Majorization step

A majorizing function

$$\begin{aligned}
 f(\mathbf{X}) &= \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2 + \frac{1}{2} \sum_{(i,j) \notin \Omega} 0 + \lambda \|\mathbf{X}\|_* \\
 &\leq \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2 + \frac{1}{2} \sum_{(i,j) \notin \Omega} (x_{ij}^{(t)} - x_{ij})^2 + \lambda \|\mathbf{X}\|_* \\
 &= \frac{1}{2} \|\mathbf{X} - \mathbf{Z}^{(t)}\|_{\text{F}}^2 + \lambda \|\mathbf{X}\|_* \\
 &= g(\mathbf{X} | \mathbf{X}^{(t)}),
 \end{aligned}$$

where  $\mathbf{Z}^{(t)} = P_{\Omega}(\mathbf{Y}) + P_{\Omega^\perp}(\mathbf{X}^{(t)})$

## Minimization step

- ▶ Majorizing function

$$g(\mathbf{X}|\mathbf{X}^{(t)}) = \frac{1}{2}\|\mathbf{X}\|_F^2 - \text{tr}(\mathbf{X}^\top \mathbf{Z}^{(t)}) + \frac{1}{2}\|\mathbf{Z}^{(t)}\|_F^2 + \lambda\|\mathbf{X}\|_*$$

- ▶ Observe

$$\|\mathbf{X}\|_F^2 = \|\sigma(\mathbf{X})\|_2^2 = \sum_i \sigma_i^2 \quad \circledcirc$$

$$\|\mathbf{Z}^{(t)}\|_F^2 = \|\sigma(\mathbf{Z}^{(t)})\|_2^2 = \sum_i \omega_i^2$$

and by Fan-von Neuman's inequality

$$\text{tr}(\mathbf{X}^\top \mathbf{Z}^{(t)}) \leq \sum_i \sigma_i \omega_i$$

with equality achieved if and only if the left and right singular vectors of the two matrices coincide

- ▶ Thus we can assume  $\mathbf{X}$  has same singular vectors as  $\mathbf{Z}^{(t)}$  and

$$\begin{aligned} g(\mathbf{X}|\mathbf{X}^{(t)}) &= \frac{1}{2} \sum_i \sigma_i^2 - \sum_i \sigma_i \omega_i + \frac{1}{2} \omega_i^2 + \lambda \sum_i \sigma_i \\ &= \frac{1}{2} \sum_i (\sigma_i - \omega_i)^2 + \lambda \sum_i \sigma_i, \end{aligned}$$

with minimizer given by  $\sigma_i^{(t+1)} = (\omega_i - \lambda)_+$

## MM algorithm for matrix completion

Initialize  $\mathbf{X}^{(0)} \in \mathbb{R}^{m \times n}$

**repeat**

$$\mathbf{Z}^{(t+1)} \leftarrow P_{\Omega}(\mathbf{Y}) + P_{\Omega^\perp}(\mathbf{X}^{(t)})$$

$$\text{Compute SVD } \mathbf{Z}^{(t+1)} = \mathbf{U}\text{diag}(\mathbf{w})\mathbf{V}^T$$

$$\mathbf{X}^{(t+1)} \leftarrow \mathbf{U}\text{diag}[(\mathbf{w} - \lambda)_+] \mathbf{V}^T$$

**until** objective value converges

- ▶ “Golub-Kahan-Reinsch” algorithm takes  $4m^2n + 8mn^2 + 9n^3$  flops for a  $m \geq n$  matrix and is not going to work for 480K-by-18K Netflix matrix. Notice only top singular values are needed.  
Lanczos/Arnoldi algorithm is the way to go

# Outline

## Overview

Problem Setup

Review of Newton's method

## EM algorithm

Introduction

Canonical Example - Finite Mixture Model

In the Zoo of EM

## MM algorithm

Introduction

MM examples

## Acceleration of EM/MM

SQUAREM

Quasi-Newton acceleration

# Acceleration

- ▶ Slow convergence of EM/MM (**linear rate** at best) is a serious concern in many problems
- ▶ Some proposed acceleration methods
  - ▶ Aitken's and Louis' methods (Louis, 1982)
  - ▶ Conjugate gradient method (Jamshidian and Jennrich, 1993)
  - ▶ Quasi-Newton acceleration (Lange, 1995; Jamshidian and Jennrich, 1997)
  - ▶ Ikeda acceleration (Ikeda, 2000)
- ▶ We review two recent ones which are particularly attractive in **high-dimensional** setting

# SQUAREM

Varadhan and Roland (2008)

- ▶  $F : \mathbb{R}^p \mapsto \mathbb{R}^p$  algorithm mapping

$$\mathbf{x} \rightarrow F(\mathbf{x}) \rightarrow F \circ F(\mathbf{x}) \rightarrow F \circ F \circ F(\mathbf{x}) \rightarrow \dots$$

- ▶ SQUAREM extrapolation

$$\begin{aligned}\mathbf{x}^{(t+1)} &= \mathbf{x}^{(t)} - 2s[F(\mathbf{x}^{(t)}) - \mathbf{x}^{(t)}] + s^2[F \circ F(\mathbf{x}^{(t)}) - 2F(\mathbf{x}^{(t)}) + \mathbf{x}^{(t)}] \\ &= \mathbf{x}^{(t)} - 2s\mathbf{u} + s^2\mathbf{v}\end{aligned}$$

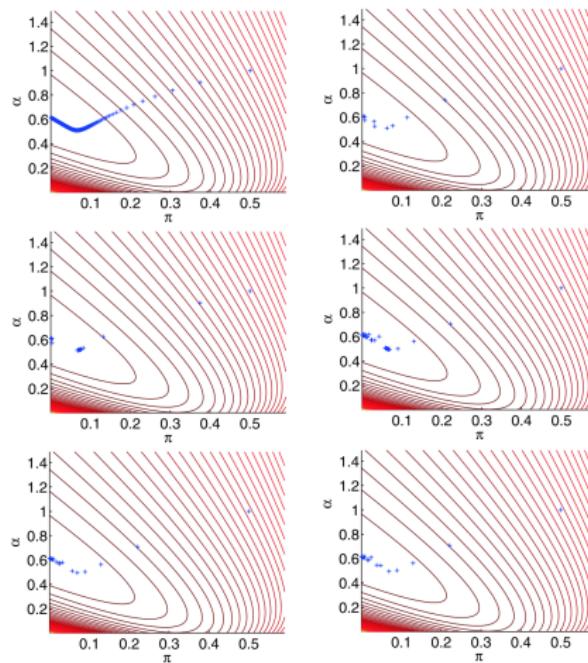
where  $s$  is a scalar

- ▶ SqS1:  $s = \frac{\mathbf{u}^\top \mathbf{u}}{\mathbf{u}^\top (\mathbf{v} - \mathbf{u})}$
- ▶ SqS2:  $s = \frac{\mathbf{u}^\top (\mathbf{v} - \mathbf{u})}{(\mathbf{v} - \mathbf{u})^\top (\mathbf{v} - \mathbf{u})}$
- ▶ SqS3:  $s = -\sqrt{\frac{\mathbf{u}^\top \mathbf{u}}{(\mathbf{v} - \mathbf{u})^\top (\mathbf{v} - \mathbf{u})}}$

- ▶ Extremely simple and effective
- ▶ Ascent property may be lost by the extrapolation; but we can always revert back to the EM/MM update as safeguard

# An example with 2 parameters

**Fig. 1** Ascent of the different algorithms for the Lidwell and Somerville household type (a) data starting from  $(\pi^0, \alpha^0) = (0.5, 1)$  with stopping criterion  $\varepsilon = 10^{-9}$ . Top left: naive MM; Top right:  $q = 1$ ; Middle left:  $q = 2$ ; Middle right: SqS1; Bottom left: SqS2; Bottom right: SqS3



## Quasi-Newton acceleration

Zhou et al. (2011)

- ▶  $F : \mathbb{R}^p \mapsto \mathbb{R}^p$  algorithm mapping

$$\mathbf{x} \rightarrow F(\mathbf{x}) \rightarrow F \circ F(\mathbf{x}) \rightarrow F \circ F \circ F(\mathbf{x}) \rightarrow \dots$$

- ▶ Fixed points of  $F \Leftrightarrow$  roots of  $G(\mathbf{x}) = \mathbf{x} - F(\mathbf{x})$
- ▶ Newton's method for finding the roots of the mapping  $G(\mathbf{x}) = \mathbf{x} - F(\mathbf{x})$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - [I - dF(\mathbf{x}^{(t)})]^{-1}[\mathbf{x}^{(t)} - F(\mathbf{x}^{(t)})]$$

- ▶ Idea: Approximate  $dF(\mathbf{x}^{(t)})$  by a low-rank matrix  $M$  and explicitly form the inverse  $(I - M)^{-1}$

- ▶ **Secant condition:**  $u = F(x^{(t)}) - x^{(t)}$ ,  $v = F \circ F(x^{(t)}) - F(x^{(t)})$ .

$$Mu = v,$$

where  $M = dF(x^{(t)})$

- ▶ For best results, we collect  $q$  secant pairs  $u_i, v_i, i = 1, \dots, q$  and require that

$$Mu_i = v_i, \quad i = 1, \dots, q, \quad \text{or}$$

$$MU = V.$$

## Lemma

The minimum of the strictly convex function  $\|M\|_F^2$  subject to the constraint  $MU = V$  is attained by the choice  $M = V(U^T U)^{-1}U^T$ .

- ▶ Fortunately, we have the explicit inverse by Woodbury formula

$$(I - M)^{-1} = [I - V(U^T U)^{-1}U^T]^{-1} = I + V[U^T(U - V)]^{-1}U^T$$

- ▶ Quasi-Newton update

$$\begin{aligned}x^{(t+1)} &= x^{(t)} - (I - M)^{-1}[x^{(t)} - F(x^{(t)})] \\&= x^{(t)} - [I - V(U^T U)^{-1}U^T]^{-1}[x^{(t)} - F(x^{(t)})] \\&= F(x^{(t)}) - V[U^T(U - V)]^{-1}U^T[x^{(t)} - F(x^{(t)})]\end{aligned}$$

## Advantages of the new quasi-Newton acceleration scheme

$$x^{(t+1)} = F(x^{(t)}) - V[U^T(U - V)]^{-1} U^T[x^{(t)} - F(x^{(t)})]$$

- ▶ The effort per iteration is light:  $O(pq^2) + O(q^3)$  ⚡
- ▶ Linear constraints are preserved
- ▶ Do not need to store the whole approximate (inverse) Hessian matrix:

$$O(p^2) \text{ vs } O(pq)$$

$p = \# \text{ parameters}$ ,  $q = \# \text{ secant conditions}$

## An example with 2771 parameters

**Table 8** Comparison of accelerations for the movie rating problem. Here the starting point is  $\pi_i = \alpha_i = \beta_j = 0.5$ , the stopping criterion is  $\varepsilon = 10^{-9}$ , and the number of parameters equals 2,771

Algorithm	ln $L$	Evals	Time
EM	-119085.2039	671	189.3020
$q = 1$	-119085.2020	215	64.1149
$q = 2$	-119085.1983	116	36.6745
$q = 3$	-119085.1978	153	46.0387
$q = 4$	-119085.1961	156	46.9827
$q = 5$	-119085.1974	161	48.6629
SqS1	-119085.2029	341	127.9918
SqS2	-119085.2019	301	110.9871
SqS3	-119085.2001	157	56.7568

- Baum, L. E., Petrie, T., Soules, G., and Weiss, N. (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Ann. Math. Statist.*, 41:164–171.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Ser. B*, 57:289–300.
- Booth, J. G. and Hobert, J. P. (1999). Maximizing generalized linear mixed model likelihoods with an automated Monte Carlo EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodology)*, 61(1):265–285.
- Celeux, G. and Diebolt, J. (1985). The SEM algorithm: A probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Computational Statistics Quarterly*, 2:73–82.
- Celeux, G. and Diebolt, J. (1989). Une version de type recuit simulé de l'algorithme EM. Rapport de Recherche INRIA.
- Cooley, J. W. and Tukey, J. W. (1965). An algorithm for the machine calculation of complex Fourier series. *Math. Comp.*, 19:297–301.
- Cox, D. R. (1972). Regression models and life-tables. *J. Roy. Statist. Soc. Ser. B*, 34:187–220. With discussion by F. Downton, Richard Peto, D. J. Bartholomew, D. V. Lindley, P. W. Glassborow, D. E. Barton, Susannah Howard, B. Benjamin, John J. Gart, L. D. Meshalkin, A. R. Kagan, M. Zelen, R. E. Barlow, Jack Kalbfleisch, R. L. Prentice and Norman Breslow, and a reply by D. R. Cox.
- de Leeuw, J. and Heiser, W. (1977). Convergence of correction matrix algorithms for multidimensional scaling. *Geometric Representations of Relational Data*.

- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B.*, 39(1-38).
- Dickey, D. A. and Fuller, W. A. (1979). Distribution of the estimators for autoregressive time series with a unit root. *J. Amer. Statist. Assoc.*, 74(366, part 1):427–431.
- Efron, B. (1979). Bootstrap methods: another look at the jackknife. *Ann. Statist.*, 7(1):1–26.
- Gelfand, A. E. and Smith, A. F. M. (1990). Sampling-based approaches to calculating marginal densities. *J. Amer. Statist. Assoc.*, 85(410):398–409.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Series in Statistics. Springer, New York, second edition.
- Ikeda, S. (2000). Acceleration of the EM algorithms. *Systems and Computers in Japan*, 31(2):10–18.
- Jamshidian, M. and Jennrich, R. I. (1993). Conjugate gradient acceleration of the EM algorithm. *J. Amer. Statist. Assoc.*, 88(421):221–228.
- Jamshidian, M. and Jennrich, R. I. (1997). Acceleration of the EM algorithm by using quasi-Newton methods. *J. Roy. Statist. Soc. Ser. B*, 59(3):569–587.
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *J. Amer. Statist. Assoc.*, 53:457–481.
- Lange, K. (1995). A quasi-Newton acceleration of the EM algorithm. *Statist. Sinica*, 5(1):1–18.

- Lange, K., Hunter, D. R., and Yang, I. (2000). Optimization transfer using surrogate objective functions. *J. Comput. Graph. Statist.*, 9(1):1–59. With discussion, and a rejoinder by Hunter and Lange.
- Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791.
- Lee, D. D. and Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In *NIPS*, pages 556–562. MIT Press.
- Liu, C. and Rubin, D. B. (1994). The ECME algorithm: a simple extension of EM and ECM with faster monotone convergence. *Biometrika*, 81(4):633–648.
- Liu, C., Rubin, D. B., and Wu, Y. N. (1998). Parameter expansion to accelerate EM: the PX-EM algorithm. *Biometrika*, 85(4):755–770.
- Liu, J. S. and Wu, Y. N. (1999). Parameter expansion for data augmentation. *J. Amer. Statist. Assoc.*, 94(448):1264–1274.
- Louis, T. A. (1982). Finding the observed information matrix when using the EM algorithm. *J. Roy. Statist. Soc. Ser. B*, 44(2):226–233.
- Mazumder, R., Hastie, T., and Tibshirani, R. (2010). Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research*, 11:2287–2322.
- Meng, X.-L. and Rubin, D. B. (1993). Maximum likelihood estimation via the ECM algorithm: a general framework. *Biometrika*, 80(2):267–278.
- Meng, X.-L. and van Dyk, D. (1997). The EM algorithm—an old folk-song sung to a fast new tune. *J. Roy. Statist. Soc. Ser. B*, 59(3):511–567. With discussion and a reply by the authors.

- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–1092.
- Neal, R. M. and Hinton, G. E. (1999). A view of the EM algorithm that justifies incremental, sparse, and other variants. In Jordan, M. I., editor, *Learning in graphical models*, pages 355–368. MIT Press, Cambridge, MA, USA.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo Statistical Methods*. Springer Texts in Statistics. Springer-Verlag, New York, second edition.
- Suchard, M. A., Wang, Q., Chan, C., Frelinger, J., Cron, A., and West, M. (2010). Understanding GPU programming for statistical computation: studies in massively parallel massive mixtures. *Journal of Computational and Graphical Statistics*, 19(2):419–438.
- Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *J. Amer. Statist. Assoc.*, 82(398):528–550. With discussion and with a reply by the authors.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288.
- Tropp, J. A. (2006). Just relax: convex programming methods for identifying sparse signals in noise. *IEEE Transactions on Information Theory*, 52(3):1030–1051.
- Varadhan, R. and Roland, C. (2008). Simple and globally convergent methods for accelerating the convergence of any EM algorithm. *Scand. J. Statist.*, 35(2):335–353.

- Wei, G. C. G. and Tanner, M. A. (1990). A Monte Carlo implementation of the EM algorithm and the poor man's data augmentation algorithms. *Journal of the American Statistical Association*, 85(411):699–704.
- Zhou, H., Alexander, D., and Lange, K. (2011). A quasi-Newton acceleration for high-dimensional optimization algorithms. *Statistics and Computing*, 21:261–273.