

"Partitioning the likelihood"[†]

SOMA DHAVALA
"Dept. of Statistics"
Texas A & M University

to be presented star 613 class

[†] based on "partitioned algorithms for maximum likelihood
other non-linear estimation" by Gordon Smyth,
Statistics & Computing, 6, 201-216, 1996

Scope & Context of the talk:

$\log(x)$ ← Model formulation (6/3)

present talk ← Compute ML estimates

(6/3) ← Analyze bias/variance/
robustness/
consistency/
model selection

Maximum likelihood estimation

$y_i \stackrel{iid}{\sim} p(\theta)$: θ : parameter vector of dimension " p "

y : observed data

p : a probability density function parametrized by θ

$$L(\theta) = \prod_{i=1}^n p_{\theta}(y_i) \quad (\text{likelihood})$$

$$l(\theta) = \log L(\theta) \quad \text{log likelihood}$$

$$\hat{\theta}_{MLE} = \max_{\theta} l(\theta)$$

So is MLE a turn-the-crank procedure??

$$\hat{\theta}_{mle} = \max_{\theta} l(\theta)$$

is just the beginning of the harsh reality

Explore partitioning as a means to
efficiently solve an otherwise difficult
optimization problem

Partitioning : data, parameter,
space objective function

ways: inferential
(eg: nuisance parameter)

: computational

(uncorrelated parameters /
orthogonal parameters)

:

essentially an optimization problem

- stability
- speed of convergence
- accuracy

Parameter space partition:

- Reduced objective function
- nested iterations
- Zig-zag iterations
- leap-frog iterations

Data-partition:

Distributed ML estimation

- EM
- Sub-object gradient methods

• reduced objective function

$$\gamma(\theta_2) = l(\hat{\theta}_1(\theta_2), \theta_2)$$

- dimensionality reduction
- possible when $\hat{\theta}_1(\theta_2)$ is in closed-form
- converges faster than -full objective function (curse-of-dimensionality)

* Profile likelihood when θ_1 is a nuisance parameter

$$\begin{aligned}\theta_1 &= f_1(\theta_1, \theta_2) \\ \theta_2 &= f_2(\theta_1, \theta_2)\end{aligned} \rightarrow (1)$$

Nested algorithms:

$$\theta_2^{k+1} = F_2(\hat{\theta}_1(\theta_2^k), \theta_2^k)$$

F_2 does involve θ_1 and θ_2
unlike reduced objective function

Zig-zag :

$$\theta_1^{k+1} = \hat{\theta}_1(\theta_2^k)$$

$$\theta_2^{k+1} = \hat{\theta}_2(\theta_1^{k+1})$$

eg: $\mu_i = \beta_1 + \beta_2 x_i$ and $\sigma_i^2 = \exp(\gamma_1 + z\gamma_2)$
 both β and γ are estimated regressing
 given others.

Leap-frog :

$$\theta_1^{k+1} = F_1(\theta_1^k, \theta_2^k)$$

$$\theta_2^{k+1} = F_2(\theta_1^{k+1}, \theta_2^k)$$

F_1 & F_2 are update equations
not estimates

Properties:

- 1) Reduced objective function: faster, stable
- 2) Nested; can not be worse than full-iterations
- 3) Zig-zag: stable but slow, faster if parameters are uncorrelated
- 4) leap-frog: less guaranteed, similar to zig-zag if parameters are orthogonal

Scoring Method :

$$\theta^{k+1} = \theta^k - \ddot{l}^T(\theta^k) \dot{l}(\theta^k)$$

\downarrow \downarrow
 observed score
 information matrix

$$\theta^{k+1} = \theta^k - \ddot{L}^T(\theta^k) \dot{l}(\theta^k)$$

\downarrow expected IM

Partitioned (θ_1, θ_2)

$$\theta_2^{k+1} = \theta_2^k - \ddot{l}_{2.1}^{-1} \dot{l}_2$$

Contd..

where
$$\ddot{l} = \begin{pmatrix} \ddot{l}_{11} & \ddot{l}_{12} \\ \ddot{l}_{21} & \ddot{l}_{22} \end{pmatrix} \text{ and}$$

$$\ddot{l}_{21} = \ddot{l}_{22} - \ddot{l}_{21} \ddot{l}_{11}^{-1} \ddot{l}_{12}$$

Newton Raphson of profile likelihood ←

nested Fisher-scoring method

$$A^{-1} = \begin{pmatrix} I & -A_{11}^{-1} A_{12} \\ 0 & I \end{pmatrix} \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & A_{22 \cdot 1}^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ -A_{21} A_{11}^{-1} & I \end{pmatrix}$$

$$A_{22 \cdot 1} = A_{22} - A_{21} A_{11}^{-1} A_{12}$$

then

$$\theta_2^{k+1} = \theta_2^k + A_{2,1}^{-1} \dot{l}_{2,1} \longrightarrow \textcircled{3}$$

however

$$\dot{l}_{2,1} = \dot{l}_2 \text{ at } \theta_1 = \hat{\theta}_1(\theta_2)$$

since

$$\dot{l}_{2,1} = \dot{l}_2 - A_{21} A_{11}^{-1} \underbrace{\dot{l}_1}_{\downarrow 0}$$

$$\text{or } \theta_1 = \hat{\theta}_1(\theta_2)$$

③

evaluated at

$$\theta_1^k = \hat{\theta}_1(\theta_2^k)$$

$$\theta_2 = \theta_2^k$$

Gauss-Newton :

$$\theta^{k+1} = \theta^k + (\tilde{\mu}^T \mu) \tilde{\mu}^T (y - \mu)$$

where $E(y) = \mu$

 Partitioned Gauss-Newton

$$\theta_2^{k+1} = \theta_2^k + \left[\tilde{\mu}_2^T (I - P) \tilde{\mu}_2 \right]^{-1} \tilde{\mu}_2^T (y - \mu)$$

$$\theta_1 = \hat{\theta}_1(\theta_2)$$

If θ_1 ~~an~~ is linear, then,
the above is called Seperable
least-squares ---

→ y $\mu = X(\theta_2) \theta_1$, then

$$\hat{\theta}_1 = (X^T X)^{-1} X^T y \quad \text{and}$$

$$\hat{\mu}_1 = X$$

$$P = \hat{\mu} (\hat{\mu}^T \hat{\mu})^{-1} \hat{\mu}^T$$

projection matrix of $\hat{\mu}$

→ Separable least-squares is
equivalently a Fisher-scaling if

nested

$$y \sim N(x(\theta_2) \theta_1, \sigma^2 I)$$

since $y \sim N(\mu(\theta), c(\theta))$

$$I(\theta) = \dot{\mu}^T c^{-1} \dot{\mu} + \frac{1}{2} \text{tr} \left[\left(c^{-1}(\theta) \frac{\partial c}{\partial \theta} \right)^2 \right]$$

stability! / ill-conditioned matrices

$$\theta^{k+1} = \theta^k + A^T S$$

① A is 'ill-conditioned'

then

$$A_n = (\alpha I + A)$$

α : small +ve constant

②

$$\theta^{k+1} = \theta^k + \underline{\alpha} \cdot A^T S$$

α : damping factor

Partitioning in EM algorithm:

Standard EM: in \mathbf{I}_c is over-informative

Rate-of-Convergence: \propto largest eigen vector of

observed
incomplete IM
↓

$$\mathbf{I}_c = \mathbf{I}_{IC} + \mathbf{I}_m$$

$$\mathbf{I}_c^{-1} \mathbf{I}_m$$

↓
missing
information
matrix

choose "missing" information s.t

Γ_c is as small as possible.

bride: partition:

$\theta \rightarrow (\theta_1, \theta_2)$

$\theta_1 \rightarrow$ separate hidden-space

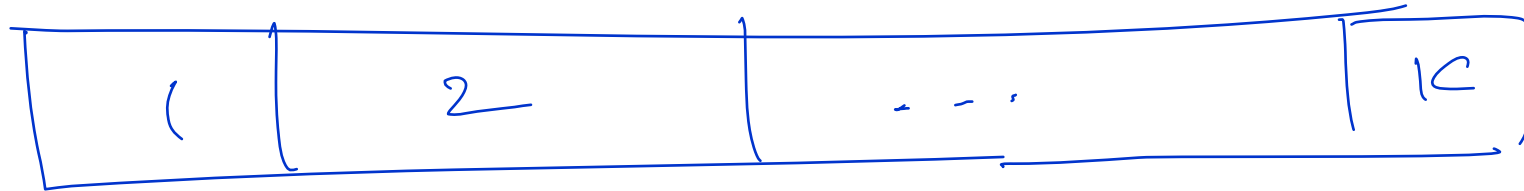
$\theta_2 \rightarrow$ separate hidden-space



SA GE

space-alternating generalized EM
algorithm

partition the data



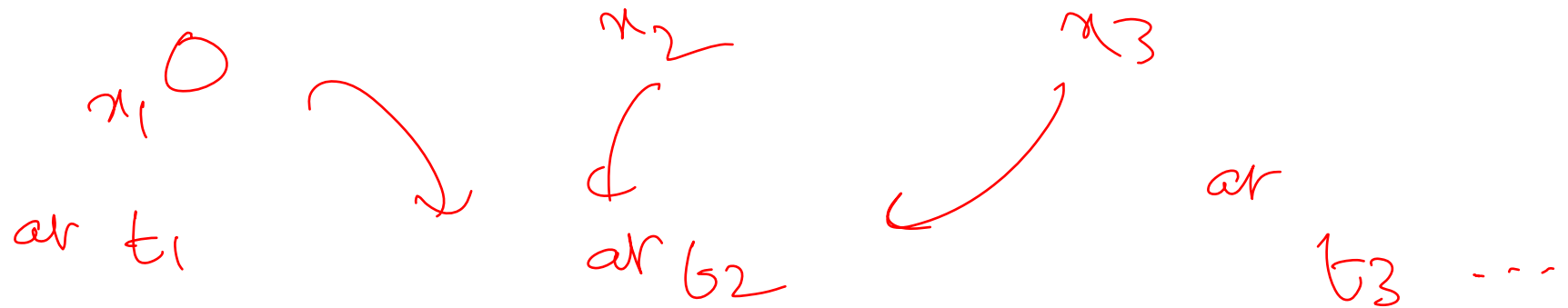
k - blocks of data:

Partial E - step

Partial M - step



useful in distributed MC estimation



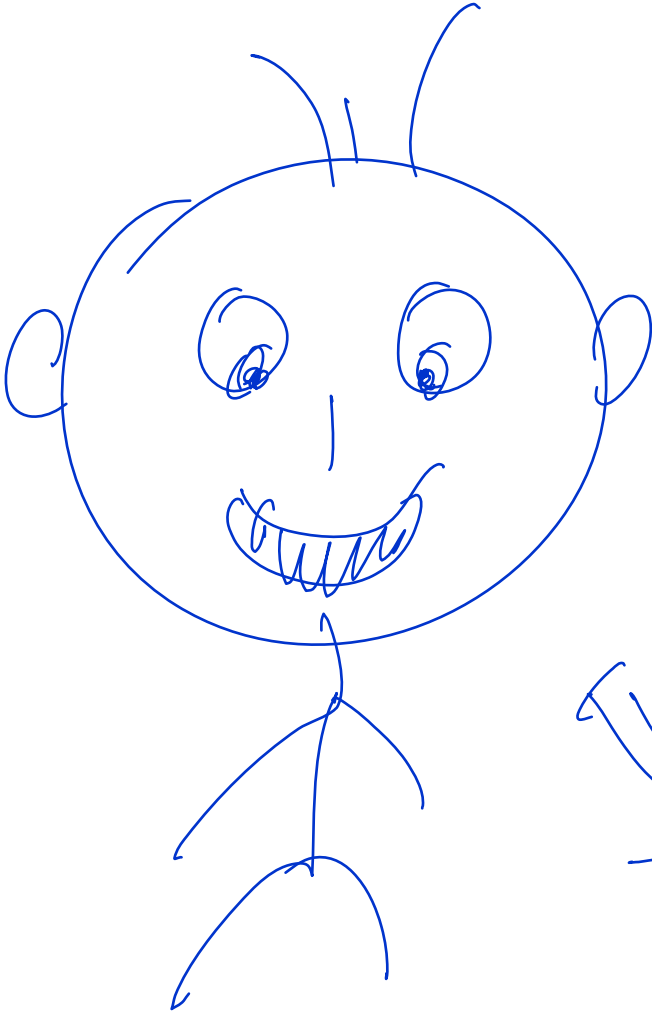
$$f(x) = \sum_i f_i(x)$$

Sub-objective functions

$$\theta^{k+1} = \theta^k - \underbrace{\alpha}_{\text{step-size}} \cdot \left\{ \sum_i \nabla f_i(\theta^k) \right\} \text{ gradient}$$

References: Google those keywords

- 1) Krishnan & McCulloch: "EM algorithm & its extensions"
wiley & sons
- 2) Fessler and A.O. Hero, "space alternating generalized
EM algorithm" SAGE
- 3) Neal and Hilton, "A view of EM algorithm that
justifies incremental view point" Incremental
EM algorithm
- 4) Dimitris (stanford university), "Incremental LMS/
steepest descent algorithms", distributed optimization
- 5) Robert Nowak, "Distributed EM algorithm
- 6) Golube, "review of variable projection separable
least-squares



Thank you
- Soma