**Case Study 3: fMRI Prediction**

# Coping with Large Covariances: Graphical Models, Graphical LASSO

today

Machine Learning for Big Data

CSE547/STAT548, University of Washington

Emily Fox

May 5th, 2015

1

---

# Multivariate Normal Models

- So far, we looked at univariate multiple regression

$$y^i = \beta_0 + \beta_1 y^i + \dots \beta_p x_p^i + \epsilon^i \qquad \epsilon^i \sim \mathcal{N}(0, \sigma^2) \qquad y^i \in \mathbb{R}$$

$$= \beta^T x^i + \epsilon^i$$

$$\Longrightarrow \quad y^i \sim N(\beta^T x^i, \sigma^2)$$

- If one has a multivariate response $y^i \in \mathbb{R}^d$ ← # of semantic features
  - Assuming independence between dimensions

So far

$$\in \mathbb{R}^d \left\{ y^i \sim N\left( \begin{bmatrix} - \beta^{(1)T} - \\ - \beta^{(2)T} - \\ \vdots \\ - \beta^{(d)T} - \end{bmatrix} x^i, \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix} \right) \right.$$

B^T

leads to d ind. problems

$$\beta^{(\ell)} \text{ are coeff. for the } \ell^{th} \text{ semantic feature}$$

2

# Multivariate Normal Models

- If one has a multivariate response $y^i \in \mathbb{R}^d$
  - Assuming correlation between the output dimensions

  "dog" and "furry"

  $$y^i \sim N(B^\top x^i, \Sigma)$$

  non-diagonal

  $$\text{recall}: \ \text{Cov}(y_s, y_t) = \Sigma_{st}$$

- Assume linear (or other mean regression) is removed and focus on the correlation structure

  $$y^i \sim N(0, \Sigma)$$

  $\Sigma$ sym. pos. def.

- Matrix valued parameter!

  see more on matrix valued params in Case Study 4

©Emily Fox 2015                              3

# Low-Rank Approximations

- In pictures…

  $$\Sigma_0 = \text{diag}(\sigma_1^2, \ldots, \sigma_d^2)$$

  $d \times d$      $d \times k$

  $$\Sigma = \Lambda\Lambda' + \Sigma_0$$

  very big matrix

  $d \{$ ▮ $= d\{$ ▮ $+$ ◫

  $k$   low rank

  $d$ diag terms

  maintains pos. def.

- Number of parameters:

  $$dk + d = d(k+1) \ll \frac{d(d+1)}{2}$$

  sig. reduction in param. for $k \ll d$

©Emily Fox 2015                              4

# Latent Factor Models

- Original multivariate regression

*here: assume linear term is removed*

$$\mathbf{y}^i = B^T x^i + \epsilon^i, \qquad \epsilon^i \sim N(0, \Sigma)$$

- Latent factor model assumption: $\Sigma = \Lambda\Lambda' + \Sigma_0$
- Low-rank approximation arises from a latent factor model

$$y^i = \Lambda \eta^i + \tilde{\epsilon}^i \qquad \eta^i \sim N_k(0, I)$$

*dxk kx1*

*obs.* *"Factor loadings"* *"latent factors"* *ind. across dims*

$$\tilde{\epsilon}^i \sim N_d(0, \Sigma_0)$$

*← diag.*

- Proof:
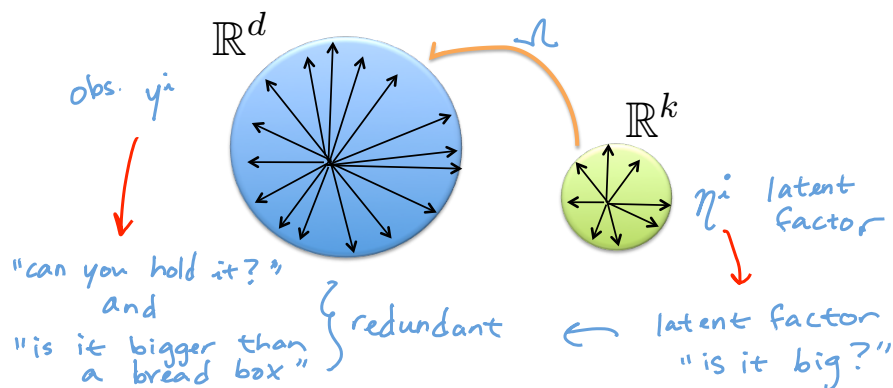
$$Cov(y, \Lambda, \Sigma_0) = E[(y - E[y])(y - E[y])^T] = E[yy^T]$$

$$= E[(\Lambda\eta + \tilde{\epsilon})(\Lambda\eta + \tilde{\epsilon})^T] = \Lambda E[\eta\eta^T]\Lambda$$

$$+ 2E[\eta^T]\Lambda E[\tilde{\epsilon}]$$

$$= \Lambda I \Lambda^T + \Sigma_0 \qquad + E[\tilde{\epsilon}\tilde{\epsilon}^T] \to \Sigma_0$$

©Emily Fox 2015    5

# Lower-dim Embeddings

## Sharing information in
## *low-dim subspace*

$$\mathbb{R}^d$$

$$\mathbb{R}^k$$

*obs. $y^i$*

$$\eta^i \text{ latent factor}$$

*"can you hold it?"*
*and*
*"is it bigger than a bread box"*
*} redundant*

*← latent factor*
*"is it big?"*

©Emily Fox 2015    6

3

# Sparsity Assumptions

- What if we assume $\Sigma$ is sparse?

$$(i \neq j) \quad \Sigma_{ij} = 0 \xrightarrow{\text{Gaussian}} \quad y_i \perp\!\!\!\perp y_j \quad \leftarrow \text{ind.}$$

$$\text{cov}(y_i, y_j) = 0$$

Could assume $\Sigma$ sparse to reduce # params, but each 0 encodes an <u>indep.</u> statement
→ often too strong of an assumption

- More often, we can reasonably make statements about *conditional independence*

$$\text{"cat"} \perp\!\!\!\perp \text{"dog"} \mid \text{"animal", "furry", "pet"} \quad \leftarrow \text{cond. on}$$

©Emily Fox 2015    7

# Information Form

- Motivations for considering "information form" of multivariate normal
  - ☐ Easier to read off conditional densities
  - ☐ Has log-linear form in terms of "information parameters"

$$\frac{1}{\sqrt{2\pi |\Sigma|}} \, e^{-\frac{1}{2}(y-\mu)^T \Sigma^{-1} (y-\mu)} \qquad \leftarrow y \sim N(\mu, \Sigma)$$

$$\Updownarrow \quad \begin{aligned} \Omega &= \Sigma^{-1} \\ \eta &= \Sigma^{-1}\mu \end{aligned}$$

$$y^T \Sigma^{-1} y \checkmark$$
$$-2 y^T \Sigma^{-1} \mu \checkmark$$
$$+ \mu^T \Sigma^{-1} \mu$$
$$\underbrace{\phantom{xxx}}_{\substack{\text{const.} \\ \text{wrt } y}}$$

$$\propto e^{\eta^T y - \frac{1}{2} y^T \Omega y} \qquad \leftarrow y \sim N^{-1}(\eta, \Omega)$$
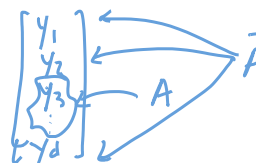
©Emily Fox 2015    8

4

# Conditional Densities

- Assume a model with

$$y \sim \mathcal{N}^{-1}(\eta, \Omega)$$

and divide the dimensions into two sets $A, \bar{A}$

Submatrix of $\Omega$ with rows indexed by indices in $A$ and columns by $\bar{A}$

- Then,

$$\begin{bmatrix} y_A \\ y_{\bar{A}} \end{bmatrix} \sim \mathcal{N}^{-1}\left( \begin{bmatrix} \eta_A \\ \eta_{\bar{A}} \end{bmatrix}, \begin{bmatrix} \Omega_{AA} & \Omega_{A\bar{A}} \\ \Omega_{\bar{A}A} & \Omega_{\bar{A}\bar{A}} \end{bmatrix} \right)$$

$$p(y_A | y_{\bar{A}}) = \mathcal{N}^{-1}\left( \eta_A - \Omega_{A\bar{A}} y_{\bar{A}}, \Omega_{AA} \right)$$

©Emily Fox 2015

9

---

# Conditional Densities

- Let $A = \{s, t\}$   "cat"  "dog"   $\bar{A} = $ everything else   $A = \{s, t\}$

$y_s, y_t$     $y_{\backslash st}$

$$p(y_A \mid y_{\bar{A}}) = \mathcal{N}^{-1}(\eta_A - \Omega_{A\bar{A}} y_{\bar{A}}, \Omega_{AA})$$

$$\begin{bmatrix} \Omega_{ss} & \Omega_{st} \\ \Omega_{ts} & \Omega_{tt} \end{bmatrix}$$

what if $\Omega_{st} = 0$? $\Rightarrow \begin{bmatrix} \Omega_{ss} & 0 \\ 0 & \Omega_{tt} \end{bmatrix}$

$$\text{cov}(y_s, y_t | y_{\backslash st}) = \Omega_{AA}^{-1} = \begin{bmatrix} \Omega_{ss}^{-1} & 0 \\ 0 & \Omega_{tt}^{-1} \end{bmatrix}$$
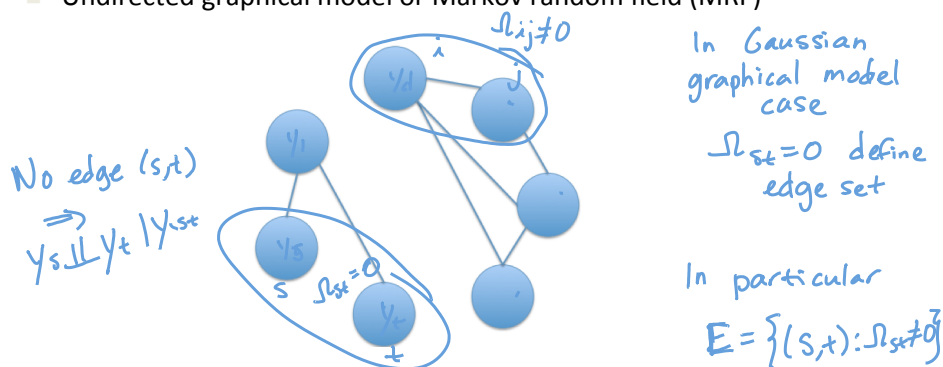
- Therefore,   $\Leftarrow$   $\boxed{y_s \perp\!\!\!\perp y_t \mid y_{\backslash st} \Leftrightarrow \Omega_{st} = 0}$

©Emily Fox 2015

10

5

# Connection with Graphical Models

- Undirected graphical model or Markov random field (MRF)



$\Omega_{ij} \neq 0$

In Gaussian graphical model case

$\Omega_{st} = 0$ define edge set

No edge $(s,t)$
$\Rightarrow y_s \perp\!\!\!\perp y_t \mid y_{\setminus st}$

$s \quad \Omega_{st} = 0$

In particular

$E = \{(s,t) : \Omega_{st} \neq 0\}$

$$p(y \mid \eta, \Omega) \propto \prod_t \psi_t(y_t) \prod_{(s,t) \in E} \psi_{st}(y_s, y_t)$$

node potentials      edge potentials

$$\psi_t(y_t) \propto e^{\eta_t y_t}$$
$$\psi_{st}(y_s, y_t) \propto e^{-\frac{1}{2} y_s \Omega_{st} y_t}$$

©Emily Fox 2015

11

# Sparse Precision vs. Covariance

- For a sparse precision matrix, the covariance need not be

zeros encode statements of cond. ind.

read graph structure from here



$\Rightarrow y$ is still fully correlated

does not imply sparsity in cov. (ind. statements)

©Emily Fox 2015

12

# ML Estimation for Given Graph

- Assume a known graph $G = \{V, E\}$    (nodes, edges)

$$\frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(y-\mu)^T \Sigma^{-1}(y-\mu)}$$

- Rewrite log likelihood:

$$\log p(y|\theta) = \frac{N}{2}\log|\Omega| - \frac{1}{2}\sum_i \underbrace{(y_i-\mu)^T}_{x^T}\underbrace{\Omega}_{A}\underbrace{(y_i-\mu)}_{x}$$

$$= \frac{N}{2}\log|\Omega| - \frac{1}{2}\sum_i tr\left[(y_i-\mu)(y_i-\mu)^T \Omega\right]$$

$$\triangleq \frac{N}{2}\log|\Omega| - \frac{1}{2}tr(S_\mu \Omega)$$

$$\sum_i (y_i-\mu)(y_i-\mu)^T$$

trace trick:
$$x^T A x = tr(x^T A x)$$
$$= tr(x x^T A)$$

$$L(\Omega) = \log|\Omega| - tr(S\Omega)$$

$$\frac{1}{N}S_\mu = \text{sample cov.}$$

In our case, $\mu = 0$

# ML Estimation for Given Graph

$$L(\Omega) = \log|\Omega| - \text{tr}(S\Omega)$$

- Take gradient:

$$\nabla L(\Omega) = \Omega^{-1} - S$$

constraint on our ML esti. from graph structure

$$\text{s.t. } \Omega_{st} = 0 \quad \text{if } (s,t) \notin E \quad \leftarrow \text{ linear constraint}$$

$$\Omega \quad \text{pos. def., sym. matrix} \quad \leftarrow \text{ hard}$$

- Many approaches to solving:
  - ☐ Barrier method – add penalty discouraging $\Omega$ from leaving the positive definite cone (Dahl et al. 2008)
  - ☐ Coordinate descent method (cf., Hastie et al. 2009)
  - ☐ …

# ML Estimation for Given Graph

- Can show that the optimal solution satisfies

$$\hat{\Sigma}_{st}^{ML,G} = S_{st} \qquad \begin{array}{l} \text{if } (s,t) \in E \\ \text{if } s = t \end{array} \qquad \text{match sample cov.}$$

$$\Omega_{st} = 0 \qquad \text{if } (s,t) \notin E$$

- Example:

adj. matrix
1 = edge

$$G = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} \qquad S = \begin{pmatrix} 10 & 1 & 5 & 4 \\ 1 & 10 & 2 & 6 \\ 5 & 2 & 10 & 3 \\ 4 & 6 & 3 & 10 \end{pmatrix}$$

$$\Omega = \begin{pmatrix} \cdot & \cdot & 0 & \cdot \\ \cdot & \cdot & \cdot & 0 \\ 0 & \cdot & \cdot & \cdot \\ \cdot & 0 & \cdot & \cdot \end{pmatrix} \qquad \hat{\Sigma}^{ML,G} = \begin{pmatrix} 10 & 1 & 1.31 & 4 \\ 1 & 10 & 2 & 0.87 \\ 1.31 & 2 & 10 & 3 \\ 4 & 0.87 & 3 & 10 \end{pmatrix}$$

"=$\Sigma^{-1}$"

15

---

# Estimating Graph Structure

- To learn the structure of the Gaussian graphical model, we want to trade off fit and sparsity
  - Measure of fit:  log-likelihood

    $$\log|\Omega| - tr(S\Omega) + const.$$

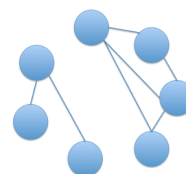  - Encouraging sparsity:  $\Omega_{st} = 0 \implies$ no edge "sparsity"

    $$\|\Omega\|_1 = \sum_{s,t} |\Omega_{s,t}| \quad \longleftarrow \quad \text{want to min}$$

- Overall objective = "graphical LASSO" or "Glasso"

  $$F(\Omega) = -\log|\Omega| + tr(S\Omega) + \lambda \|\Omega\|_1$$

  Just as in LASSO, but w/ a matrix parameter <u>and</u>  s.t. $\Omega \succ 0$

  ↑ pos. def.

16

# Solving the Graphical LASSO

- Objective is convex, but non-smooth as in LASSO   *... subgrad.*
- Also, positive definite constraint!

- There are many approaches to optimizing the objective
  - ☐ Most common = coordinate descent akin to shooting algorithm (Friedman et al. 2008)   *See HW 3*

- Some issues…
  - ☐ Ballpark: several minutes for a 1000-variable problem
  - ☐ Algorithms scale as $O(d^3)$

- Other approach = ADMM   *also HW 3*

©Emily Fox 2015                                                    17

# Faster Computations

From Daniela Witten's talk at JSM 2012:

1. The $j$th variable is unconnected from all others in the graphical lasso solution if and only if $|S_{ij}| \leq \lambda$ for all $i = 1, \dots, j-1, j+1, \dots, p$.   *← Sample cov is small relative to chosen penalty*
2. Let **A** denote the $p \times p$ matrix whose elements take the form $A_{ii} = 1$, $A_{ij} = 1_{|S_{ij}| > \lambda}$. Then the connected components of **A** are the same as the connected components of the graphical lasso solution.   *ind. on the threshholded values*

   We can obtain the *exact* right answer by solving the graphical lasso on each connected component separately!
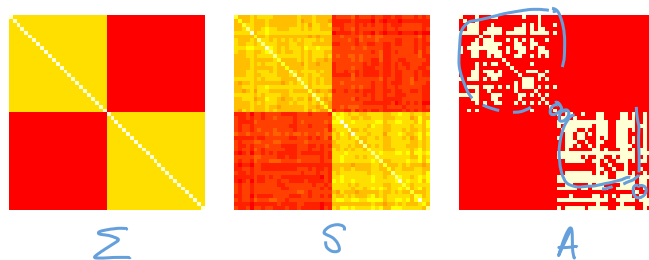
   Citations: Witten et al. JCGS 2011, Mazumder and Hastie JMLR 2012

©Emily Fox 2015                                                    18

9

# Covariance Screening for Glasso

From Daniela Witten's talk at JSM 2012:



- ▸ The solution to the graphical lasso problem with $\lambda = 0.7$ has five connected components (why 5?!)
- ▸ Perform graphical lasso on each component separately!
- ▸ Reduction in computational time: From $O(50^3)$ to $O(24^3)$.

©Emily Fox 2015                                                          19