

10-708

06/06/2020

19-66MS, Using models review

(\*) Matrix-inverse lemma

for a block-partitioned matrix  $M = \begin{pmatrix} E & F \\ G & H \end{pmatrix}$  with  $E$  and  $H$  invertible.

• Matrix-inverse lemma:-

$$(E - FH^{-1}G)^{-1} = E^{-1} + E^{-1}E(H - GE^{-1}E)^{-1}GE^{-1}$$

Method:

• diagonalise  $M$ :-

$$\begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} E & F \\ G & H \end{bmatrix} \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} = \begin{bmatrix} E - FH^{-1}G & 0 \\ 0 & H \end{bmatrix}$$

• Schur's complement: (notable in numerical analysis)

$$M/H = E - FH^{-1}G$$

• Note: for  $XYZ = W$ ,  $Y^{-1} = ZW^{-1}X$

$$M^{-1} = \begin{bmatrix} E & F \\ G & H \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} \begin{bmatrix} E - FH^{-1}G & 0 \\ 0 & H \end{bmatrix}^{-1} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix}$$

$$= \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} \begin{bmatrix} (M/H)^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix}$$

$$= \begin{bmatrix} (M/H)^{-1} & 0 \\ -H^{-1}G(M/H)^{-1} & H^{-1} \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix}$$

$$= \begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1}G(M/H)^{-1}FH^{-1} + H^{-1} \end{bmatrix}$$

- Applying A.2.6:

$$\underline{M}^{-1} = \begin{bmatrix} \underline{E}' + \underline{E}' \underline{F} (\underline{M}/\underline{E})^{-1} \underline{G} \underline{E}^{-1} & -\underline{E}' \underline{F} (\underline{M}/\underline{E})^{-1} \\ -(\underline{M}/\underline{E})^{-1} \underline{G} \underline{E}^{-1} & (\underline{M}/\underline{E})^{-1} \end{bmatrix} \quad (D)$$

- see Murphy (2012) 5.4.3-4.1

(E) is achieved by decomposing in terms of  $\underline{E}$  and

$$\underline{M}/\underline{E} = (\underline{H} - \underline{G} \underline{E}^{-1} \underline{F})$$

- we get MIL as a corollary

(\*) Covariance and precision matrices

$$\underline{\Sigma} = \begin{bmatrix} \sigma_{11} & \sigma_1^T \\ \sigma_1 & \underline{\Sigma}_{-1} \end{bmatrix}$$

$\sigma_{11}$  - scalar  
 $\sigma_1, \sigma_1^T \in \mathbb{R}^{N-1}$   
 $\underline{\Sigma}_{-1} \in \mathbb{R}^{(N-1) \times (N-1)}$

$$\underline{x} = \begin{bmatrix} x_1 \\ | \\ x_{-1} \\ | \end{bmatrix} \quad \underline{x} \in \mathbb{R}^N$$

$$\underline{M}^{-1} = \begin{bmatrix} (\underline{M}/\underline{H})^{-1} & -(\underline{M}/\underline{H})^{-1} \underline{F} \underline{H}^{-1} \\ -\underline{H}^{-1} \underline{G} (\underline{M}/\underline{H})^{-1} & \underline{H}^{-1} + \underline{H}^{-1} \underline{G} (\underline{M}/\underline{H})^{-1} \underline{F} \underline{H}^{-1} \end{bmatrix} \quad - \underline{Q} = \underline{\Sigma}^{-1}$$

$$(\underline{M}/\underline{H})^{-1} = \sigma_{11}^{-1} + \sigma_{11}^{-1} \sigma_1^T (\underline{\Sigma}_{-1} - \sigma_1 \sigma_{11}^{-1} \sigma_1^T)^{-1} \sigma_1 \sigma_{11}^{-1} = q_{11}$$

$$\underline{Q} = \begin{bmatrix} q_{11} & -q_{11} \sigma_1^T \underline{\Sigma}_{-1}^{-1} \\ -q_{11} \underline{\Sigma}_{-1}^{-1} \sigma_1 & \underline{\Sigma}_{-1}^{-1} (1 + q_{11} \sigma_1 \sigma_1^T \underline{\Sigma}_{-1}^{-1}) \end{bmatrix} = \begin{bmatrix} q_{11} & q_1^T \\ q_1 & \underline{Q}_{-1} \end{bmatrix}$$

(\*) single node conditional

$$\text{Partition } \underline{x} = \begin{bmatrix} x_1 \\ x_{-1} \end{bmatrix} \quad p\left(\begin{bmatrix} x_1 \\ x_{-1} \end{bmatrix} \mid \underline{M}, \underline{\Sigma}\right) = N\left(\begin{bmatrix} x_1 \\ x_{-1} \end{bmatrix} \mid \begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_{x_{-1}} \end{bmatrix}, \begin{bmatrix} \underline{\Sigma}_{x_1 x_1} & \underline{\Sigma}_{x_1 x_{-1}} \\ \underline{\Sigma}_{x_{-1} x_1} & \underline{\Sigma}_{x_{-1} x_{-1}} \end{bmatrix}\right)$$

Then we have:

$$p(x_i | x_{-i}) = N(x_i | \mu_{i,-i}, \Sigma_{i,-i})$$

- apply cond. Gaussian.

$$\mu_{i,-i} = \mu_i + \Sigma_{x_i x_{-i}} \Sigma_{x_{-i} x_{-i}}^{-1} (x_{-i} - \mu_{x_{-i}})$$

$$\Sigma_{i,-i} = \Sigma_{x_i x_i} - \Sigma_{x_i x_{-i}} \Sigma_{x_{-i} x_{-i}}^{-1} \Sigma_{x_{-i} x_i}$$

- set  $\mu_i = 0$

$$\Rightarrow p(x_i | x_{-i}) = N(\Sigma_{x_i x_{-i}} + \Sigma_{x_{-i} x_{-i}}^{-1} x_{-i} | \Sigma_{x_i x_i} - \Sigma_{x_i x_{-i}} \Sigma_{x_{-i} x_{-i}}^{-1} \Sigma_{x_{-i} x_i})$$

$$= N(\sigma_i^{-1} \Sigma_{-i}^{-1} x_{-i}, q_{i|-i})$$

✓

0/52

$$= N\left(\frac{q_i^T}{-q_{ii}} x_{-i}, q_{i|-i}\right)$$

✓

(\*)  $p(x_i | x_{-i})$  is a 1-D Gaussian dist. (univariate)!

$\Rightarrow$  mean, covariance i.e.  $\frac{q_i^T}{-q_{ii}} x_{-i}$ ,  $q_{i|-i}$  are both scalars

(\*) Amended lecture

notes; decipher now what is going on.

- conduct  $N$  conditional autoregressions (using  $L_1$ -regularisation LASSO) of  $x_i$  against  $x_{-i}$ .

- each autoregression will yield an estimate  $\frac{q_i^T}{-q_{ii}}$ , which

can be used to populate the precision matrix

$Q$  column by column (and as it is symmetric also)

As we are focusing on presence/absence of non-zero entries rather than value in itself

(\*) Note that we do not go for straightforward inversion of  $\Sigma$  to get precision  $Q$

- we attempt to estimate  $Q$  directly (or columns of it) by exploiting one-to-one correspondence of  $Q$  (in terms of sparseness) with param. estimates of  $N$  autoregressions



(\*) Achieved with a sparsity constraint

(\*) Structure learning  $\rightarrow$  focus on sparseness not value of components of precision.  $S_i = \{j: j \neq i, \theta_{ij} \neq 0\}$

(\*) known as graphical LASSO

(\*) This is the Meinhausen-Buhlman algorithm

(\*) Two alternatives: - (Friedman, Hastie, Tibshirani 2007)

1) - directly estimate  $\Omega$

- compute sample variance covariance matrix.
- use  $L_1$  reg. ML estimation.

evolving social networks

- Infer  $\mathcal{I}$

- time-evolving graphs are dependent
- small subset of nodes/edges modified between each time point
- find an algorithm to estimate time-specific graph using entire dataset of graphs evolving through time (one for each time point?)

KELLER - cyclically run graph algorithm that performs neighbourhood selection. (

- similar loss to earlier.

$$\hat{\theta}_i^t = \underset{\theta_i^t}{\operatorname{argmin}} L_w(\theta_i^t) + \lambda_1 \|\theta_i^t\|_1 \quad \forall t$$

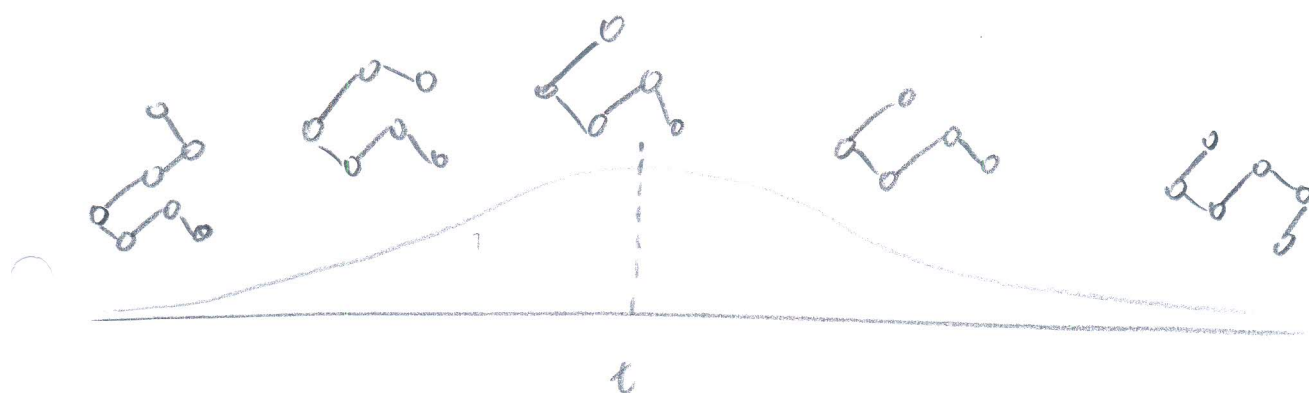
$$L_w(\theta_i^{(t)}) = \sum_{t'=1}^T w(x^{t'}; x^t) \log p(x_i^{t'} | x_{-i}^{t'}, \theta_i^t)$$

(\*) Summing over time-series of vertex examples

- weight function defined over two time points  $w(x^{t'}; x^t)$
- (kernel?) - considers 'distance' between network at two time points.
- use this to do a weighted sum of condit. likelihoods of a node of interest  $x_i$  given the rest of nodes. (over time points)

### (\*) Intuition:

- for graph of interest at time  $t$ , and you have one example at that time.
- Acknowledge that other examples of graphs not from time  $t$  may be relevant (due to dependence) to graph at time  $t$ .
- the closer the graph at another time  $t'$  is to graph of interest, the more relevant info this may contain.



- introducing this weight function intuitively allows for use of entire dataset at every time point (some kind of temporal-based proximity measure aka kernel)

- an example of a non-parametric neighborhood selection
- Recall  $\rightarrow$  kernel density estimation ~~by~~ weights samples according to a 'unit of gravity'
- see Kolar, Le Song, Ahmed, Xing (2010)
- nonparametric neighborhood selection

- cond. like

$$P_{\theta}^t(x_i^t | x_{\setminus i}^t) = \text{logistic}(2x_i^t \langle \theta_{\setminus i}^t, x_{\setminus i}^t \rangle)$$

- Neighborhood sel.  $S(x_i) = \{j | \theta_{i,j}^t \neq 0\}$

- Time specific graph reg.

- estimate at  $t^* \in [0, 1]$

$$\min_{\theta \in \mathbb{R}^{p-1}} \left\{ - \sum_{t \in T^n} w_t(t^*) \gamma(\theta_i; x^t) + \lambda \|\theta\|_1 \right\}$$

$$\gamma(\theta_i^t; x^t) = \log P_{\theta_i^t}^t(x_i^t | x_{\setminus i}^t)$$

$$w_t(t^*) = \frac{K_{hn}(t - t^*)}{\sum_{t' \in T^n} K_{hn}(t' - t^*)}$$



# (\*) Inference (II)

TESLA: temporally smoothed  $L_1$  reg. logistic regression

$$\hat{\theta}_i^1, \dots, \hat{\theta}_i^T = \underset{\theta_i^1, \dots, \theta_i^T}{\operatorname{argmin}} \sum_{t=1}^T \ell_{\text{avg}}(\theta_i^t) + \lambda_1 \sum_{t=1}^T \|\theta_{-i}^t\|_1 + \lambda_2 \sum_{t=2}^T \|\theta_i^t - \theta_i^{t-1}\|_2^2 \quad (i) \quad (ii)$$

$$\ell_{\text{avg}}(\theta_i^t) = \frac{1}{N^t} \sum_{d=1}^{N^t} \log P(x_{d,i}^t | x_{d,-i}^t, \theta_i^t) \quad (I)$$

- previous graph algo (KELLER); estimate param for one graph at time point  $t$  and summed over all examples.

- TESLA: simultaneously estimate params of graphs simult. from  $t=1 \dots T$ .

- loss - sum of point. likelihoods of  $X_i | X_{-i}$  with no reweighting
- sparsity  $\rightarrow$  every graph has sparse structure (i)  
 $\hookrightarrow$  differences between adjacent pair of graph structures in time also sparse (ii)
- evolution of graph structure is 'smooth' over time;
- graph structure between two time points minimally distinct

can rewrite (I) using box constraints (rewrite above constraints)

TESLA: 
$$\min_{\theta_i^1, \dots, \theta_i^T, u_i^1, \dots, u_i^T, v_i^1, \dots, v_i^T} \sum_{t=1}^T \ell(x_i^t; \theta_i^t) + \lambda_1 \sum_{t=1}^T \|u_i^t\|_1 + \lambda_2 \sum_{t=2}^T \|v_i^t\|_1$$

s.t. 
$$-u_{i,j}^t \leq \theta_{i,j}^t \leq u_{i,j}^t \quad t=1, \dots, T \quad \forall j \in V \setminus i$$

$$-v_{i,j}^t \leq \theta_{i,j}^t - \theta_{i,j}^{t-1} \leq v_{i,j}^t \quad t=2, \dots, T \quad \forall j \in V \setminus i \quad (\text{total variation differences})$$

- can run different optimisation algorithms
- consistency guarantees

(\*) theoretical analysis of KELLER, TESLA  $\rightarrow$  see Koller, Xing (2009)  
 Koller, Le Song, Amund, Xing (2010)

## Ex) Applications

- social network time-series data
- senator network  $\rightarrow$  shows predictive power on Chafee's political movements!
- breast cancer  $\rightarrow$  use network inference
- TV distance is significant if you know dependency structure across networks (through time)

## (\*) Estimating time varying networks; Isler, U. Song, Ahmed Xing (2010)

- paper contains rigorous formulation
  - $\hookrightarrow$  time-varying graphical structure estimation of GGM for discrete r.v.s. (temporally smoothed L-reg. log. regression)
- Assumptions required on structure of parameter vector.
- optimisation algorithms presented.
- Hyperparam selection  $\rightarrow$  (graph sparsity, BIC and a heuristic)
- Real-world datasets  $\rightarrow$  particularly interesting
- voting: - Time varying characteristics of voting bill records 109<sup>th</sup> congress (e.g. senators)
  - senators with swaying political stance can be discovered (not possible for time-invariant network global estimation)

## (\*) High-dimensional graphs and variable selection with the LASSO, Meinshausen and Bühlmann (2006)

- Pattern of zero-entries in precision matrix of MVG  $\rightarrow$  C.I. restrictions between variables
- neighbourhood selection via LASSO is computationally attr. alt. to standard covariance selection for sparse high-dimensional graphs
- A equivalent to variable selection for Gaussian linear models.

(\*) gives results on consistency

- consistent estimation of the full-edge set in a sparse high-dim graph is possible (asymptotically, prob. of estimating correct neighbours converges exponentially to 1).