



FACULTAD DE
INGENIERÍA



UNIVERSIDAD
DE LA REPÚBLICA
URUGUAY

Aprendizaje Automático para Datos en Grafos

Aprendiendo Grafos desde Datos

Marcelo Fiori

Muy basado en transparencias de **Gonzalo Mateos**

`m fiori@fing.edu.uy`

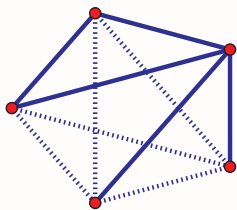
`http://www.fing.edu.uy/~mfiori/`

Octubre 2022

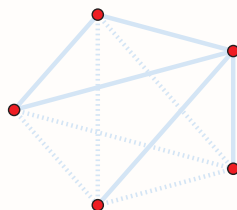
Recordemos

- Tres problemas canónicos de **inferencia de topología del grafo** [Kolaczyk'09]
 - (i) Predicción de enlaces
 - (ii) Association network inference ← Énfasis de estas clases
 - (iii) Tomographic network topology inference

Association network inference



Original graph



Association network inference

- Supongamos que solo observamos la señal en el grafo $\mathbf{x} = [x_1, \dots, x_N]^\top$; y
- Asumimos (i, j) definido por un ‘nivel de asociación’ no trivial entre x_i, x_j
- **Objetivo:** predecir estado de aristas para todos los pares de vértices $\mathcal{V}^{(2)}$

Recordemos

- Redes de correlación
 - ⇒ Test de hipótesis usando la correlación empírica
- Correlaciones parciales
 - ⇒ Similar, pero condicionando a otros nodos para evitar influencias indirectas

Undirected Gaussian graphical models

- Supongamos que las variables $\{x_i\}_{i \in \mathcal{V}}$ tienen distribución gaussiana multivariada
⇒ Consideramos $\rho_{ij|\mathcal{V} \setminus \{i,j\}}$ condicionada respecto a todos los otros nodos ($m = N - 2$)

Teorema

Bajo este supuesto gaussiano, los nodos $i, j \in \mathcal{V}$ tienen correlación parcial

$$\rho_{ij|\mathcal{V} \setminus \{i,j\}} = 0$$

si y solo si x_i y x_j son *condicionalmente independientes* dados $\{x_k\}_{k \in \mathcal{V} \setminus \{i,j\}}$

- **Def:** el *grafo de independencia condicional* $G(\mathcal{V}, \mathcal{E})$ tiene conjunto de aristas

$$\mathcal{E} = \left\{ (i, j) \in \mathcal{V}^{(2)} : \rho_{ij|\mathcal{V} \setminus \{i,j\}} \neq 0 \right\}$$

⇒ Es un caso especial y popular de redes de correlaciones parciales

- También conocido como *Gaussian Markov random field (GMRF)*

Covariance selection

- Sea Σ la matriz de covarianza de $\mathbf{x} = [x_1, \dots, x_N]^\top$

Def: la **matriz de precisión** es $\Theta := \Sigma^{-1}$ con entradas θ_{ij}

- **Resultado clave:** Para GMRFs, las correlaciones parciales pueden ser expresadas como

$$\rho_{ij|V \setminus \{i,j\}} = -\frac{\theta_{ij}}{\sqrt{\theta_{ii}\theta_{jj}}}$$

\Rightarrow Entradas no nulas en $\Theta \Leftrightarrow$ Aristas en el grafo G

- El problema de inferir G desde \mathcal{X} es conocido como **covariance selection** [Dempster'74]

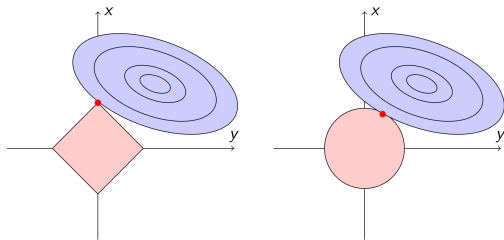
\Rightarrow Métodos clásicos son 'network-agnostic,' y testean

$$H_0 : \rho_{ij|V \setminus \{i,j\}} = 0 \quad \text{versus} \quad H_1 : \rho_{ij|V \setminus \{i,j\}} \neq 0$$

- En general no escalable, y $P \ll N$ entonces la estimación de $\hat{\Sigma}$ es desafiante

Digresión: Sparsity y la norma ℓ_1

- Consideremos la minimización de una función cuadrática de θ como en mínimos cuadrados
- Q: ¿Cuál es el efecto de una restricción de norma ℓ_1 , i.e., $\|\theta\|_1 = \sum_i |\theta_i| \leq \tau$?



⇒ Si conjuntos de nivel intersecan restricciones en un vértice → **soluciones esparsas**

- Estimador Lasso permite hacer **selección de variables** [Tibshirani'94]

$$\hat{\theta}_{Lasso} = \arg \min_{\theta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta)^2, \text{ s. to } \|\theta\|_1 \leq \tau$$

Graphical Lasso

- Estimador de máxima verosimilitud de Θ con regularización que promueve esparsidad [Yuan-Lin'07]

$$\hat{\Theta} \in \arg \max_{\Theta \succeq 0} \left\{ \log \det \Theta - \text{trace}(\hat{\Sigma} \Theta) - \lambda \|\Theta\|_1 \right\}$$

⇒ Efectivo cuando $P \ll N$, da lugar a modelos interpretables

⇒ Métodos escalables usan coordinate-descent [Friedman et al'08]

- **Garantía de performance:** Graphical lasso con $\lambda = 2\sqrt{\frac{\log N}{P}}$ satisface

$$\|\hat{\Theta} - \Theta_0\|_2 \leq \sqrt{\frac{d_{\max}^2 \log N}{P}} \quad \text{w.h.p.}$$

⇒ donde Θ_0 es el ground-truth, y grado máximo d_{\max}

- **Consistencia del soporte** para $P = \Omega(d_{\max}^2 \log N)$ [Ravikumar et al'11]

GMRFs con restricciones de Laplaciano

- Graphical model selection with **Laplacian constraints** $\Theta = \mathbf{L}$
 - Entradas fuera de diagonal $\theta_{ij} = L_{ij} = -A_{ij} \leq 0 \Rightarrow$ GMRF atractivo
 - El Laplaciano es singular ($\mathbf{L}\mathbf{1} = \mathbf{0}$) \Rightarrow GMRF impropio
- Estimar un GMRF propio agregando algo en la diagonal [Lake-Tenembaum'07]

$$\begin{aligned} & \max_{\Theta \geq \mathbf{0}, \gamma \geq 0} \left\{ \log \det \Theta - \text{trace}(\hat{\Sigma}\Theta) - \lambda \|\Theta\|_1 \right\} \\ & \text{s. to } \Theta = \mathbf{L} + \gamma \mathbf{I} \\ & \quad \mathbf{L}\mathbf{1} = \mathbf{0}, L_{ij} \leq 0, i \neq j \end{aligned}$$

\Rightarrow Se interpreta γ^{-1} como la varianza de fluctuaciones isotrópicas gaussianas

- Favorece grafos sobre los cuales las señales son suaves (más sobre esto la clase siguiente)

$$\text{trace}(\hat{\Sigma}\mathbf{L}) \propto \sum_{p=1}^P \mathbf{x}_p^\top \mathbf{L} \mathbf{x}_p = \sum_{p=1}^P \text{TV}(\mathbf{x}_p)$$

Covariance selection meets linear regression

- **Idea:** Estimar de forma independiente vecindades $\mathcal{N}_i := \{j : (i, j) \in \mathcal{E}\}$, $i \in \mathcal{V}$
- El valor esperado condicional de x_i dado $\mathbf{x}_{\setminus i} := [x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N]^\top$ es

$$\mathbb{E}[x_i \mid \mathbf{x}_{\setminus i}] = \mathbf{x}_{\setminus i}^\top \boldsymbol{\beta}^{(i)} \quad (\text{es lineal})$$

- Las entradas de $\boldsymbol{\beta}^{(i)}$ se pueden expresar en términos de $\Theta = \Sigma^{-1}$, a saber

$$\beta_j^{(i)} = -\frac{\theta_{ij}}{\theta_{ii}}$$

\Rightarrow valor no nulo $\beta_j^{(i)} \Leftrightarrow$ valor no nulo θ_{ij} en $\Theta \Leftrightarrow$ Arista (i, j) in G

\Rightarrow En otras palabras, $\text{supp}(\boldsymbol{\beta}^{(i)}) := \{j : \beta_j^{(i)} \neq 0\} \equiv \mathcal{N}_i$

- Sugiere inferencia de G via regresión de mínimos cuadrados (LS), pues

$$\boldsymbol{\beta}^{(i)} = \arg \min_{\boldsymbol{\beta}} \mathbb{E} \left[(x_i - \mathbf{x}_{\setminus i}^\top \boldsymbol{\beta})^2 \right], \quad i \in \mathcal{V}$$

Neighborhood-based sparse regression

- Recorrer los nodos $i \in \mathcal{V}$ y estimar $\hat{\mathcal{N}}_i = \text{supp}(\hat{\boldsymbol{\beta}}^{(i)})$, donde

$$\hat{\boldsymbol{\beta}}^{(i)} \in \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{N-1}} \left\{ \sum_{p=1}^P (x_{pi} - \mathbf{x}_{p, \setminus i}^\top \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}$$

⇒ Problemas lasso separados por cada nodo

- No hay garantía que $\hat{\beta}_j^{(i)} \neq 0$ implica $\hat{\beta}_i^{(j)} \neq 0$ y vice versa

⇒ Combina información en $\hat{\mathcal{N}}_i$ y $\hat{\mathcal{N}}_j$ para forzar simetría

⇒ **Regla OR:** $(i, j) \in \mathcal{E}$ si $\beta_j^{(i)} \neq 0$ o $\beta_i^{(j)} \neq 0$. De igual forma, **regla AND**

- **Consistencia del soporte** para ambas reglas [Meinshausen-Bühlmann'06]

- con cierta elección de λ , esparsidad de Θ_0 , y relación de tamaño y muestras $P \ll N$

Mapa conceptual para GMRF model selection

Testing partial correlations

For each $(i, j) \in \mathcal{V} \times \mathcal{V}$, test the hypothesis

$$H_0 : \rho_{ij|\mathcal{V}\setminus ij} = 0 \quad \text{versus} \quad H_1 : \rho_{ij|\mathcal{V}\setminus ij} \neq 0$$

Covariance selection

$$\rho_{ij|\mathcal{V}\setminus ij} = -\frac{\theta_{ij}}{\sqrt{\theta_{ii}\theta_{jj}}} \rightarrow \rho_{ij|\mathcal{V}\setminus ij} \neq 0 \Leftrightarrow \theta_{ij} \neq 0$$

Infer non-zero entries $\theta_{ij} \neq 0$ of the precision matrix

$$\Theta := \Sigma^{-1}$$

Neighborhood-based regression

$$\beta_j^{(i)} = -\frac{\theta_{ij}}{\theta_{ii}} \rightarrow \beta_j^{(i)} \neq 0 \Leftrightarrow \theta_{ij} \neq 0$$

For each $i \in \mathcal{V}$, infer non-zero regression coefficients $\beta_j^{(i)} \neq 0$ in

$$\beta^{(i)} = \arg \min_{\beta} \mathbb{E} \left[(x_i - \mathbf{x}_{\setminus i}^T \beta)^2 \right]$$

Resumen comparativo

- **Paralelizable** neighborhood-based regression (NBR)

- ⇒ Verosimilitud condicionales por nodo $i \in \mathcal{V}$

- ⇒ no fuerza $\Theta \succeq \mathbf{0}$

- ⇒ Tiende a ser más rápido computacionalmente

- Graphical Lasso minimiza una **verosimilitud global** (regularizada)

$$\mathcal{L}(\Theta; \mathcal{X}) = \log \det \Theta - \text{trace}(\hat{\Sigma} \Theta)$$

- ⇒ Tiende a ser estadísticamente más eficiente

- Método NBR se puede adaptar para graphical models discretos o mixtos

- ⇒ *Ising-model selection* para $\mathbf{x} \in \{-1, +1\}^N$ [Ravikumar'10]