

GLasso: regularized maximum likelihood estimation

data X_1, \dots, X_n i.i.d. $\sim \mathcal{N}_p(\mu, \Sigma)$

goal: estimate $K = \Sigma^{-1}$ (precision matrix)

approach, called GLasso (Friedman, Hastie and Tibshirani, 2008):

$$\hat{K}, \hat{\mu} = \operatorname{argmin}_{K \succ 0, \mu} (-\log\text{-likelihood}(K, \mu; X_1, \dots, X_n) + \lambda \|K\|_1)$$

$$\hat{\mu} = n^{-1} \sum_{i=1}^n X_i \text{ decouples}$$

$$\hat{K} = \operatorname{argmin}_{K \succ 0} \underbrace{(-\log\text{-likelihood}(K, \hat{\mu}; X_1, \dots, X_n) + \lambda \|K\|_1)}_{\propto -\log(\det K) + \operatorname{trace}(\hat{\Sigma}_{\text{MLE}} K)}$$

$$\|K\|_1 = \sum_{j,k} |K_{j,k}| \text{ or } \sum_{j \neq k} |K_{j,k}|$$

$$\hat{\Sigma}_{\text{MLE}} = n^{-1} \sum_{i=1}^n (X_i - \hat{\mu})(X_i - \hat{\mu})^T$$

- ▶ GLasso is computationally (much) slower than nodewise regression
 $O(np^3)$ computational complexity (for potentially dense problems)
- ▶ GLasso provides estimates of Σ^{-1} and also of Σ by inversion
- ▶ one can run a hybrid approach:
nodewise selection first with estimated edge set \hat{E}
GLasso **restricted** to \hat{E} with $\lambda = 0$:
that is, unpenalized MLE restricted to \hat{E}

fast and accurate!

analogous to Lasso-OLS hybrid in regression

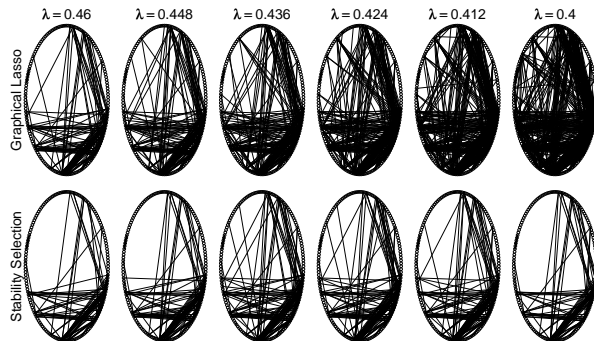
Tuning of the methods

cross-validation of the (nodewise) likelihood

and/or Stability Selection

$p = 160$ gene expressions, $n = 115$

GLasso estimator, selecting among the $\binom{p}{2} = 12'720$ features stability selection with $\mathbb{E}[V] \leq v_0 = 30$



The nonparanormal graphical model

(Liu, Lafferty and Wasserman, 2009)

motivating question: are there other “interesting” distributions, besides the Gaussian, where conditional independence between two rv.’s is encoded as zero entries in a matrix?

nonparanormal graphical model:

X has a nonparanormal distribution if there exist functions f_j ($j = 1, \dots, p$) such that

$$Z = f(X) = (f_1(X^{(1)}), \dots, f_p(X^{(p)})) \sim \mathcal{N}_p(\mu, \Sigma)$$

w.l.o.g. $\mu = 0$ and $\Sigma_{jj} = 1$

$\leadsto Z_j = f_j(X^{(j)}) \sim \mathcal{N}(0, 1)$ and therefore:

$f_j(\cdot) = \Phi^{-1}(F_j(\cdot))$ where $F_j(u) = \mathbb{P}[X^{(j)} \leq u]$: **monotone**

\leadsto a semiparametric Gaussian copula model

Lemma

Assume that (G, P) is a nonparanormal graphical model with f_j being differentiable for all $j = 1, \dots, p$. Then:

$$(j, k) \in E \iff X^{(j)} \not\perp X^{(k)} | X^{(V \setminus \{j, k\})} \iff \Sigma_{j,k}^{-1} \neq 0$$

Proof: the density of X is

$$p(x) = \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2} (f(x) - \mu)^T \Sigma^{-1} (f(x) - \mu)\right) \prod_{j=1}^p |f'_j(x_j)|$$

\rightsquigarrow the density factorizes exactly as in the Gaussian case according to Σ^{-1} □

we only have to estimate the non-zeroes of Σ^{-1}
but Σ is **not** the covariance matrix of $X = (X^{(1)}, \dots, X^{(p)})$
 Σ is the covariance matrix of the unknown $f_1(X^{(1)}), \dots, f_p(X^{(p)})$

the “best” proposal (Lue and Zhou, 2012): **rank-based!**
compute empirical rank correlation of $X^{(1)}, \dots, X^{(p)}$ with a bias correction from Kendall (1948)
denote this empirical rank correlation matrix as \hat{R} (invariant under monotone f_j 's)

stick it into GLasso:

$$\hat{K} = \operatorname{argmin}_{K \succ 0} -\log(\det K) + \operatorname{trace}(\hat{R}K) + \lambda \|K\|_1$$

this has provable guarantees in the case of a nonparanormal graphical model for estimating Σ^{-1}

as an important implication:

the rank-based version of GLasso exhibits some **robustness** for estimating the conditional independence pattern of $X \sim P$ that is: if the distribution is nonparanormal, it still works well and properly!

this is different and much better than:

GLasso works for estimating $\text{Cov}(X)^{-1}$ even if $X \sim P$ is non-Gaussian

although this is true, if sufficient amount of moments exist for non-Gaussian P : zeroes of $\text{Cov}(X)^{-1}$ do **not** encode conditional independencies!

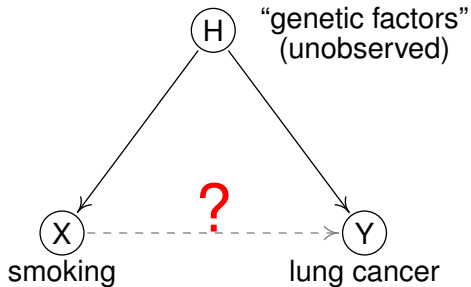
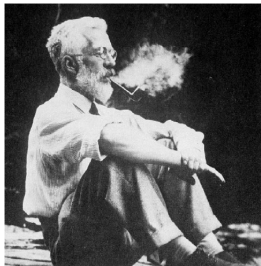
The danger of hidden confounding!

Lasso, Group Lasso, neural networks, neighborhood selection, GLasso,...

for (generalized) linear models, nonlinear models, undirected graphical models, ...

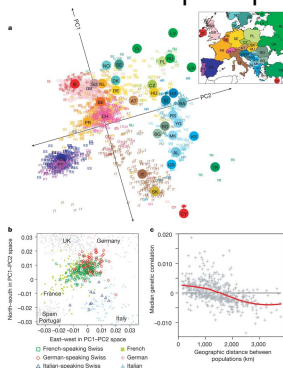
they all give “wrong” answers in presence of hidden confounding

Does smoking cause lung cancer?

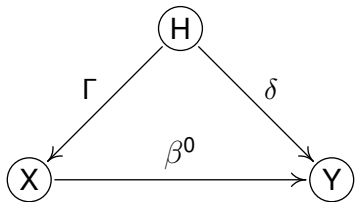


Genes mirror geography within Europe (Novembre et al., 2008)

SNP data plotted on first 2 principal components



confounding effects about geographical origin of data are found on the first principal components



$$Y \leftarrow X_{n \times p} \beta^0 + H \delta + \eta$$

$$X \leftarrow H_{n \times q} \Gamma + E$$

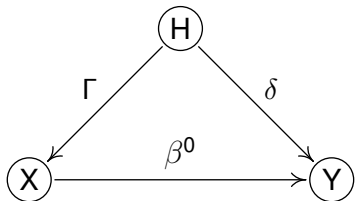
goal: infer β^0 from observations $(X_1, Y_1), \dots, (X_n, Y_n)$

the population least squares principle leads to the parameter

$$\beta^* = \operatorname{argmin}_u \mathbb{E}[(Y - Xu)^2], \quad \beta^* = \beta^0 + \underbrace{b}_{\text{"bias"}}$$

$$\|b\|_2 \leq \frac{\|\delta\|_2}{\sqrt{\text{"number of } X\text{-components affected by } H}}$$

small "bias" if confounder has dense effects!
 blessing of high dimensionality!



$$Y \leftarrow X_{n \times p} \beta^0 + H \delta + \eta$$

$$X \leftarrow H_{n \times q} \Gamma + E$$

goal: infer β^0 from observations $(X_1, Y_1), \dots, (X_n, Y_n)$

the population least squares principle leads to the parameter

$$\beta^* = \operatorname{argmin}_u \mathbb{E}[(Y - Xu)^2], \quad \beta^* = \beta^0 + \underbrace{b}_{\text{"bias"}}$$

$$\|b\|_2 \leq \frac{\|\delta\|_2}{\sqrt{\text{"number of } X\text{-components affected by } H}}$$

small “bias” if confounder has dense effects!
blessing of high dimensionality!

perhaps more importantly: view this as

$$Y = X\beta^* + \varepsilon = X \underbrace{(\beta^0 + b)}_{\text{sparse + dense}} + \varepsilon,$$

$$\varepsilon = Y - \mathbb{E}[Y|x]$$

~> we should use high-dimensional methods for “sparse + dense” regression parameter vector

- ▶ Lava (Chernozhukov, Hansen & Liao, 2017)
- ▶ Spectral Deconfounding (Ćevic, Bühlmann & Meinshausen, 2020, Guo, Ćevic & Bühlmann, 2022)

similarly for undirected graphical modeling:

$$\text{Cov}(X)^{-1} = \text{sparse matrix} + \text{low rank matrix}$$

~> use Gaussian likelihood for $\text{Cov}(X)^{-1}$ but with penalty enforcing **sparsity + low rank**

(Chandrasekaran, Parrilo & Willsky, 2012)

still lots of things to do!