

Model selection in penalized Gaussian graphical models

Ernst Wit
University of Groningen

e.c.wit@rug.nl
<http://www.math.rug.nl/~ernst>

January 2014



Penalized likelihood generates a PATH of solutions

- Consider an experiment: $|\Gamma|$ genes measured across $|T|$ time points.
- Assume n iid samples $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}$, where $\mathbf{y}^{(i)} = (y_1^{(i)}, \dots, y_{|\Gamma|}^{(i)})$.
- Assume $\mathbf{Y}^{(i)} \sim N(\mathbf{0}, \mathbf{K}^{-1})$, then

Likelihood:

$$l(\mathbf{K}|\mathbf{y}) \propto \frac{n}{2} \{\log(|K|) - \text{tr}(\mathbf{S}\mathbf{K})\}.$$

AIM: Optimization of penalized likelihood:

$$\hat{\mathbf{K}}_\lambda := \operatorname{argmax}_{\mathbf{K}} \{l(\mathbf{K}|\mathbf{y})\}$$

subject to

- $\mathbf{K} \succeq 0$;
- $\|\mathbf{K}\|_1 \leq 1/\lambda$... **for λ in some range!!**;
- some factorial colouring F .

Between proximity and truth

True process for data Y :

$$Y \sim Q.$$

A *statistical model* is a collection of measures:

$$\mathcal{M}_i = \{\mathcal{P}_\theta \mid \theta \in \Theta_i\}$$

... and typically we consider several: $\mathbb{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_k\}$.

What is the best model?

- **“Proximity”**: the model that is closest to the truth:

$$\min_{i \leq k} KL(\mathcal{P}_{\hat{\theta}_i}; Q).$$

- **“Truth”**: the model that is most likely to be the truth:

$$\max_{i \leq k} P(\mathcal{M}_i | Y).$$

With flat prior on $\Theta_{\mathcal{M}}$ and \mathcal{M} , model probability for $\mathcal{M} \in \mathbb{M}$ is:

$$\begin{aligned}
 p(\mathcal{M}|y) &\propto p(y|\mathcal{M})/p(y) \\
 &\propto \int_{\Theta_{\mathcal{M}}} e^{n\bar{\ell}(\theta)} d\theta \\
 &\approx e^{\ell(\hat{\theta})} \int_{\Theta_{\mathcal{M}}} e^{-\frac{1}{2}(\theta-\hat{\theta})^t n \frac{\partial^2}{\partial \theta^2} \bar{\ell}(\hat{\theta})(\theta-\hat{\theta})} d\theta \\
 &\approx e^{\ell(\hat{\theta})} (2\pi)^{-p/2} n^{-p/2} \left| \frac{\partial^2}{\partial \theta^2} \bar{\ell}(\hat{\theta}) \right|^{-1/2},
 \end{aligned}$$

where $p = \dim(\Theta_{\mathcal{M}})$ and $\bar{\ell}(\hat{\theta}) = \ell(\hat{\theta})/n$.

Schwarz (1978) ignored terms not depending on n :

$$p(\mathcal{M}|y) \approx e^{\ell(\hat{\theta})} n^{-p/2}, \quad \text{for large } n.$$

BIC: Bayesian information criterion

By applying the $-\log$ and ignoring constant terms c :

$$BIC(\mathcal{M}) = -2l(\hat{\theta}) + p \log(n).$$

Minimizing BIC corresponds to maximizing posterior model probability.

Define model weights $W(\mathcal{M})$,

$$W(\mathcal{M}) = e^{-BIC(\mathcal{M})/2},$$

which rescaled correspond to posterior model probabilities,

$$p(\mathcal{M}_i|y) = \frac{W(\mathcal{M}_i)}{\sum_{\mathcal{M} \in \mathbb{M}} W(\mathcal{M})}.$$

BIC for Gaussian graphical models

The BIC for an estimated Gaussian graphical model \hat{K} :

$$BIC(\hat{K}) = n(-\log |\hat{K}| + \text{Tr}(S\hat{K})) + p_{\hat{K}} \log(n),$$

where

$$p_K = |\{\text{unique non-zeroes in } K\}|,$$

which is less than $p(p-1)/2$.



The BIC for an estimated Gaussian graphical model \hat{K} :

$$BIC(\hat{K}) = n(-\log |\hat{K}| + \text{Tr}(S\hat{K})) + p_{\hat{K}} \log(n),$$

where

$$p_K = |\{\text{unique non-zeroes in } K\}|,$$

which is less than $p(p-1)/2$.

But...

- BIC is an asymptotic approximation. What is n ??
- For penalized estimate \hat{K} are number of parameters not smaller??

eBIC: extended Bayesian information criterion

The BIC for an estimated Gaussian graphical model \hat{K} :

$$eBIC_{\gamma}(\hat{K}) = n(-\log |\hat{K}| + \text{Tr}(S\hat{K})) + p_{\hat{K}} \log(n) + 4\gamma p_{\hat{K}} \log(p),$$

where

$$p_K = |\{\text{unique non-zeroes in } K\}|$$

$$p = \text{number of nodes}$$

$$\gamma = \text{tuning parameter } (0 \leq \gamma \leq 1)$$

If

$$\gamma = 0 \Rightarrow \text{ordinary BIC}$$

$$\gamma = 1 \Rightarrow \text{additional sparsity}$$

$$\gamma = 0.5 \Rightarrow \text{good trade-off (Foygel \& Drton, 2010)}$$



Proximity

Let the true density of the data Y be:

$$Y \sim q = dQ.$$

The Kullback-Leibler (KL) divergence between fitted and true model

$$KL(\hat{\theta}_i) = \int q(y) \log q(y) dy - \int q(y) \log p(y; \hat{\theta}_i) dy$$



Proximity

Let the true density of the data Y be:

$$Y \sim q = dQ.$$

The Kullback-Leibler (KL) divergence between fitted and true model

$$\begin{aligned} KL(\hat{\theta}_i) &= \int q(y) \log q(y) dy - \int q(y) \log p(y; \hat{\theta}_i) dy \\ &= C - E_q(\ell(\hat{\theta}_i)) \end{aligned}$$



Let the true density of the data Y be:

$$Y \sim q = dQ.$$

The Kullback-Leibler (KL) divergence between fitted and true model

$$\begin{aligned} KL(\hat{\theta}_i) &= \int q(y) \log q(y) dy - \int q(y) \log p(y; \hat{\theta}_i) dy \\ &= C - E_q(\ell(\hat{\theta}_i)) \\ &\approx C - \ell(\hat{\theta}_i) + p_i^* \end{aligned}$$

where

- $p(\cdot; \theta)$ is density associated with \mathcal{P}_θ .
- $p_i^* = \text{Trace} \left(J_i^{-1} K_i \right) \approx \dim(\Theta_i)$,

where J_i and K_i are Fisher informations using model \mathcal{M}_i :

$$J_i = E_g \left(\frac{\partial^2 \log f(Y, \hat{\theta}_i)}{\partial \theta \partial \theta^t} \right), \quad K_i = V_g \left(\frac{\partial \log(f(Y, \hat{\theta}_i))}{\partial \theta} \right).$$

AIC: Akaike's information criterion

By applying the $-\log$ and ignoring constant terms c :

$$AIC(\mathcal{M}) = -2l(\hat{\theta}) + 2p.$$

Minimizing AIC corresponds to maximizing posterior model probability.

Define *Akaike weights* $W(\mathcal{M})$,

$$W(\mathcal{M}) = e^{-AIC(\mathcal{M})/2},$$

which rescaled correspond to probability weights that add up to one,

$$p(\mathcal{M}_i) = \frac{W(\mathcal{M}_i)}{\sum_{\mathcal{M} \in \mathbb{M}} W(\mathcal{M})}.$$

AIC for Gaussian graphical models

The AIC for an estimated Gaussian graphical model \hat{K} :

$$AIC(\hat{K}) = n(-\log |\hat{K}| + \text{Tr}(S\hat{K})) + 2p_{\hat{K}},$$

where

$$p_K = |\{\text{unique non-zeroes in } K\}|,$$

which is less than $p(p-1)/2$.



AIC for Gaussian graphical models

The AIC for an estimated Gaussian graphical model \hat{K} :

$$AIC(\hat{K}) = n(-\log |\hat{K}| + \text{Tr}(S\hat{K})) + 2p_{\hat{K}},$$

where

$$p_K = |\{\text{unique non-zeros in } K\}|,$$

which is less than $p(p-1)/2$.

But...

- AIC is an asymptotic approximation. What is n ??
- For penalized estimate \hat{K} are number of parameters not smaller??

In the next slides, we propose three alternatives.

1. Exact AIC

$$KL(K_0 || \hat{K}) = \frac{1}{2} \{ \text{Tr}(\hat{K}\Sigma_0) - \log |\hat{K}\Sigma_0| - p \}$$

Scaling this by 2 and ignoring a constant

$$2KL(K_0 || \hat{K}) \cong -\{ \log |\hat{K}| - \text{Tr}(\hat{K}S) \} + 2 \cdot \frac{1}{2} \{ \text{Tr}(\hat{K}(\Sigma_0 - S)) \}.$$

We can write this as

$$2KL(K_0 || \hat{K}) \cong -2\ell(\hat{K}) + 2 \times \frac{1}{2} \{ \text{Tr}(\hat{K}(\Sigma_0 - S)) \}.$$

Definition (Degrees of freedom in Gaussian graphical model)

Let $Y \sim N(0, K_0^{-1})$ and \hat{K} an estimate of K_0 :

$$\text{df}_{\hat{K}} = \frac{1}{2} \{ \text{Tr}(\hat{K}(\Sigma_0 - S)) \}.$$

Approximate Exact AIC

We obviously don't know Σ_0 , but we can estimate it:

$$\hat{\Sigma}_0 = \pi S - (1 - \pi) \text{diag}\{\sigma_{11}^2, \dots, \sigma_{pp}^2\},$$

for some tuning parameter π .

Definition (Approximate Exact AIC)

Let $Y \sim N(0, K_0^{-1})$ and \hat{K} an estimate of K_0 :

$$AIC(\hat{K}) = -2\ell(\hat{K}) + 2\hat{\text{df}},$$

where

$$\hat{\text{df}} = \frac{1}{2} \{ \text{Tr}(\hat{K}(\hat{\Sigma}_0 - S)) \}.$$



2. Exact GIC

Problem of AIC: \hat{K}_λ is *not* a MLE.

GIC for M -estimator \hat{K} derived by Konishi & Kitagawa (1996):

$$\text{GIC} = -2 \sum_{k=1}^n l_k(\hat{K}; x_k) + 2\text{tr}\{R^{-1}Q\}, \quad (1)$$

where R and Q are square matrices of order p^2 given by

$$R = -\frac{1}{n} \sum_{k=1}^n \{D\psi(x_k, K)\}^T \Big|_{K=\hat{K}},$$

$$Q = \frac{1}{n} \sum_{k=1}^n \psi(x_k, K) D l_k(K) \Big|_{K=\hat{K}}.$$

Problem: Bias term in (1) requires inversion of $d^2 \times d^2$ matrix  R.

Approximate GIC (Abbruzzo, Vujacic, Wit)

We derived an explicit estimator of KL that avoids matrix inversion:

$$\widehat{\text{GIC}}(\lambda) = -2l(\hat{K}_\lambda) + 2\widehat{\text{df}}_{\text{GIC}},$$

where

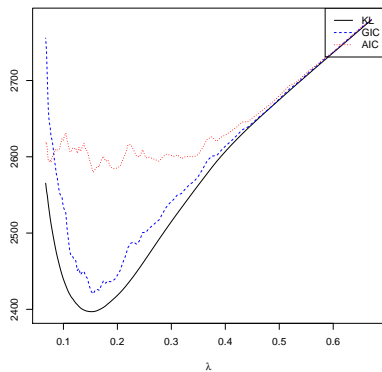
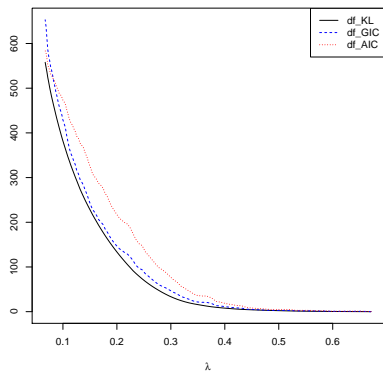
$$\widehat{\text{df}}_{\text{GIC}} = \frac{\sum_{k=1}^n c(S_k \circ I_\lambda)^\top c(\hat{K}_\lambda(S_k \circ I_\lambda) \hat{K}_\lambda)}{2n} - \frac{c(S \circ I_\lambda)^\top c(\hat{K}_\lambda(S \circ I_\lambda) \hat{K}_\lambda)}{2},$$

where

- $c()$ is the vectorize operator,
- $S_k = x_k x_k^\top$,
- $I_\lambda = 1 * (\hat{K}_\lambda \neq 0)$, an indicator matrix.



Approximate GIC: Simulations



Bias term (left) and **KL** divergence (right) estimates for $n = 100, d = 50$.

3. Exact Cross-validation (CV)

Ignoring an additive constant, the KL divergence is:

$$KL(K_0 || \hat{K}) = -\frac{1}{2} \log |\hat{K}| + E_{\Sigma_0} \frac{1}{2} \{Tr(\hat{K}XX^t)\}$$

The idea of cross-validation is to replace the final expectation by

$$KL(K_0 || \hat{K}) \approx -\frac{1}{2n} \sum_{k=1}^n \{ \log |\hat{K}^{-(k)}| + Tr(\hat{K}^{-(k)} x_k x_k^t) \}$$

Clearly, this would require recalculating the estimate $\hat{K}^{-(k)}$ many times.

Approximate Cross-validation (Vujacic, Abbruzzo, Wit)

We write the KL divergence as follows,

$$KL(K_0|\hat{K}_\lambda) = -\frac{1}{n}l(\hat{K}_\lambda) + \text{bias},$$

where $l(K) = n\{\log |K| - \text{tr}(KS)\}/2$ and $\text{bias} = \text{tr}(\hat{K}_\lambda(K_0^{-1} - S))/2$.

Definition

LOOCV-inspired estimate of Kullback-Leibler divergence

$$KL_{LCV}(\lambda) = -\frac{1}{n}l(\hat{K}_\lambda) + \frac{\sum_{i=1}^n c[(\hat{K}_\lambda^{-1} - S_i) \circ I_\lambda]^T (\hat{K}_\lambda \otimes \hat{K}_\lambda) c[(S - S_i) \circ I_\lambda]}{n(n-1)},$$

where

- $c()$ is the vectorize operator,
- $S_k = x_k x_k^T$,
- $I_\lambda = 1 * (\hat{K}_\lambda \neq 0)$, an indicator matrix.

$$\begin{aligned}
LOOCV &= -\frac{1}{2n} \sum_{i=1}^n f(S_i, \hat{K}^{(-i)}) \\
&= -\frac{1}{2} f(S, \hat{K}) - \frac{1}{2n} \sum_{i=1}^n [f(S_i, \hat{K}^{(-i)}) - f(S_i, \hat{K})] \\
&\approx -\frac{1}{n} l(\hat{K}) - \frac{1}{2n} \sum_{i=1}^n \left[\frac{df(S_i, \hat{K})}{d\Omega} \right]^\top c(\hat{K}^{(-i)} - \hat{K}).
\end{aligned}$$

Matrix differential calculus: $df(S_i, \hat{K})/d\Omega = c(\hat{K}^{-1} - S_i)$.

The term $c(\hat{K}^{(-i)} - \hat{K})$ is obtained via the Taylor expansion

$$0 \approx \frac{df(S, \hat{K})}{d\Omega} + \frac{d^2f(S, \hat{K})}{d\Omega^2} c(\hat{K}^{(-i)} - \hat{K}) + \frac{d^2f(S, \hat{K})}{d\Omega dS} c(S^{(-i)} - S).$$

Inserting: $df(S, \hat{K})/d\Omega = c(\hat{K}^{-1} - S)$, $d^2f(S, \hat{K})/d\Omega dS = -I_{p^2}$,
 $d^2f(S, \hat{K})/d\Omega^2 = -\hat{K}^{-1} \otimes \hat{K}^{-1}$

and consequently

$$c(\hat{K}^{(-i)} - \hat{K}) = -(\hat{K} \otimes \hat{K})c(S^{(-i)} - S).$$



KLCV simulation results

d=100	KL ORACLE	KLCV	AIC	GACV
n=20	8.06 (0.37)	8.60 (0.45)	12.24 (0.28)	28.59 (19.94)
n=30	6.87 (0.34)	7.29 (0.39)	10.59 (0.41)	32.07 (2.77)
n=50	5.24 (0.27)	5.63 (0.33)	7.33 (0.81)	16.93 (1.40)
n=100	3.34 (0.19)	3.57 (0.23)	3.63 (0.48)	6.81 (0.52)
n=400	1.13 (0.07)	1.20 (0.08)	1.17 (0.08)	1.24 (0.07)



Conclusions

- BIC aims to find true model
- AIC aims to come closest to the truth
- BIC gives sparser model than AIC (typically)
- AIC/BIC have asymptotic issues.
- What are number of parameters for penalized inference?
- Improved versions are available!