

"Recovering the Structure of Sparse Markov Networks from High-Dimensional Data"

Narges Bani Asadi

Irina Rish

Dimitri Kanevsky

Katya Scheinberg

Bhuvana Ramabhadran

IBM TJ Watson Research Center

September 19, 2008

Outline

- Gaussian Markov Network

Outline

- Gaussian Markov Network
- Learning Sparse Gaussian Networks

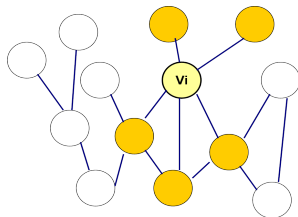
Outline

- Gaussian Markov Network
- Learning Sparse Gaussian Networks
- Learning the Sparsity of the Sparse Gaussian Networks

Outline

- Gaussian Markov Network
- Learning Sparse Gaussian Networks
- Learning the Sparsity of the Sparse Gaussian Networks
- Results

Markov Networks



$$\mathbf{X} = \{X_1, \dots, X_p\}$$

$$G = (V, E)$$

$$P(\mathbf{X}) = \frac{1}{Z} \prod_k \phi_k(\mathbf{X})$$

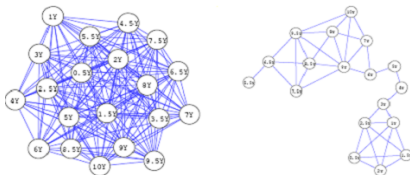
Lack of edge : conditional independence

Gaussian Markov Networks

- $P(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$
- Zeros in Σ : marginal independence
- Zeros in Σ^{-1} : conditional independence : Sparsity
- Sparsity
 - Interpretation
 - Prediction

Maximum Likelihood Estimation of the Inverse Covariance

- $\Sigma^{-1} = \arg \max_C (\log \det C - \text{Tr}(CS))$
 $\Sigma^{-1} = S^{-1}$



[borrowed from A. dAspremont presentation]

- How do we make Σ^{-1} Sparse : Penalize the likelihood

Penalized Likelihood Estimation

$$\Sigma^{-1} = \arg \max_C (\log \det C - \text{Tr}(CS) - \lambda \text{Card}(C))$$

$\text{Card}(C)$ = number of nonzero elements of C

- $\lambda = 2/(N + 1)$ for AIC
- $\lambda = \log(N + 1)/(N + 1)$ for BIC
- This is a NP-Hard combinatorial problem

L-1 Regularized Likelihood Estimation

- $\Sigma^{-1} = \arg \max_C \log(P(X))$
Subject to $\|C\|_1 \leq s$

- This is equal to solve

$$\Sigma^{-1} = \arg \max_C \log \det C - \text{tr}(SC) - \lambda \|C\|_1$$

- Convex problem with unique solution for a given λ

The Role and Choice of the Sparsity parameter: λ

- λ decides the amount of sparsity
- What is the criteria to pick the best λ ?
 - Model structure recovery
 - Prediction power on test data

Learning the Sparsity Parameter λ

- Cross Validation on Training Data
 - To maximize prediction power, ignoring the structure
 - Over fit to Model, almost No Sparsity!
- Method suggested by Banerjee et. al:

$$\lambda(\alpha) = (\max_i \sigma_i \sigma_j) \frac{t_{n-2}(\alpha/p^2)}{\sqrt{N-2+t_{n-2}^2(\alpha/p^2)}}$$

- Constructs back sparse Σ not the Σ^{-1}
- Learns a Too Sparse model!
- Weak prediction

Being Bayesian about λ

- λ as a random variable: learn its distribution
- Maximize the joint log likelihood

$$\Sigma^{-1}, \hat{\lambda} = \arg \max_{C, \lambda} \log(P(X))$$

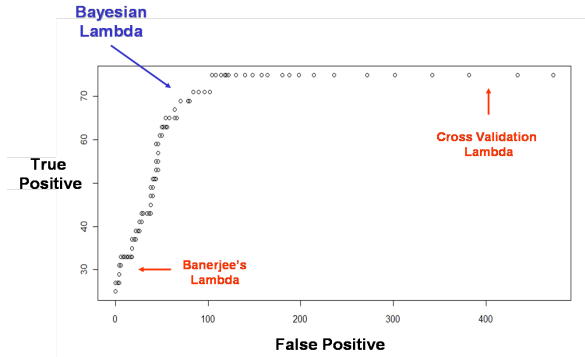
$$P(X) = P(X|C)P(C|\lambda)P(\lambda)$$

$$\begin{aligned} \log P(X) = \\ \log \det C - \text{Tr}(CS) - \lambda \|C\|_1 + P^2 \log(\lambda/2) + \log P(\lambda) \end{aligned}$$

- The choice of $P(\lambda)$

The Bayesian λ

ROC Curve

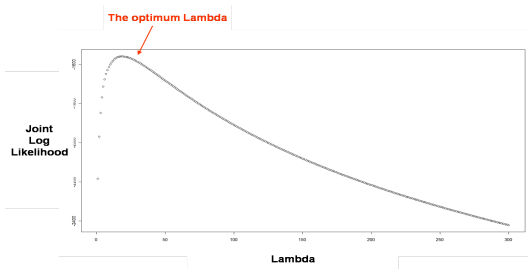


The Optimization Method

- Alternating minimization with line search
 - Estimate Σ^{-1} for an initial λ
 - Update λ in the direction of the gradient: $P^2/\lambda - \|C\|_1$
 - Iterate until convergence

The Joint likelihood with Flat prior on λ

- When $N \gg P$ we get a global maximum



- But unbounded for $N \leq P$
 - Add regularization to the objective function
 - Assume a non-flat prior for λ

The Regularized likelihood with Flat Prior on λ

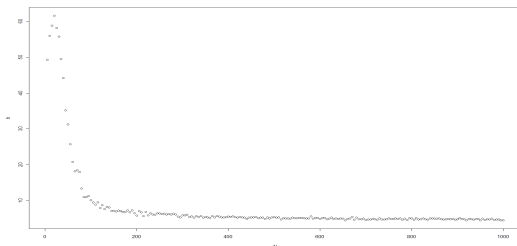
- Do not penalize the diagonal elements in the estimation of the Σ^{-1}
- Update λ as before

$$\log P(X) = \log \det C - \text{Tr}(CS) - \lambda \|C\|_1 + P^2 \log(\lambda/2) + \log P(\lambda)$$

The Joint likelihood with Exponential Prior on λ

$$P(\lambda) = -b \exp(-b\lambda) \quad E(\lambda) = 1/b$$

- How to learn b from data?
- b indicates density of the empirical inverse covariance
- approximate b with $\|S^{-1}\|_1/P^2$ multiply it by P/N



b for P=50 as a function of N

Results with Original Density = 4% and with the Prior on λ

P	N	λ	TP	FP	Prediction Error
100	30	$\lambda=34$	108	510	1.5
		$\lambda_b=603$	2	0	3.2
		$\lambda_c=2$	294	3097	0.5
100	50	$\lambda=50$	122	516	1.4
		$\lambda_b=693$	2	0	3.2
		$\lambda_c=1$	350	5087	0.6
100	500	$\lambda=62$	328	1236	0.35
		$\lambda_b=2510$	46	136	3.2
		$\lambda_c=0.1$	356	9390	0.32
100	1000	$\lambda=26$	356	2388	0.28
		$\lambda_b=3415$	60	204	3.2
		$\lambda_c=0.1$	356	9421	0.34

Results with Original Density = 52% and with the Prior on λ

P	N	λ	TP	FP	Prediction Error
100	30	$\lambda=32$	690	616	1.9
		$\lambda_b=1120$	2	2	6.5
		$\lambda_c=0.4$	2630	2157	0.62
100	50	$\lambda=47$	694	86	1.9
		$\lambda_b=2209$	6	12	6.47
		$\lambda_c=0.4$	3225	2555	0.42
100	500	$\lambda=24$	2286	1376	0.38
		$\lambda_b=5710$	116	158	5.9
		$\lambda_c=0.1$	5089	4480	0.17
100	1000	$\lambda=14$	3465	1957	0.20
		$\lambda_b=7691$	186	240	4.3
		$\lambda_c=0.1$	5102	4623	0.15

Results with Original Density = 4% and Regularized Likelihood

P	N	λ	TP	FP	Prediction Error
100	30	$\lambda=190$	22	60	3.2
		$\lambda_b=603$	2	0	3.2
		$\lambda_c=2$	294	2976	0.48
100	50	$\lambda=208$	148	156	1.4
		$\lambda_b=693$	2	0	3.2
		$\lambda_c=1$	350	4979	0.6
100	500	$\lambda=55$	336	1132	0.33
		$\lambda_b=2510$	46	136	3.2
		$\lambda_c=0.1$	356	9390	0.32
100	1000	$\lambda=27$	356	2174	0.28
		$\lambda_b=3415$	60	204	3.2
		$\lambda_c=0.1$	356	9421	0.34

Results with Original Density = 52% and Regularized Likelihood

P	N	λ	TP	FP	Prediction Error
100	30	$\lambda=500$	44	72	6.4
		$\lambda_b=1120$	2	2	6.5
		$\lambda_c=0.4$	2630	2157	0.62
100	50	$\lambda=500$	120	156	4.2
		$\lambda_b=2209$	6	12	6.3
		$\lambda_c=0.4$	3225	2555	0.42
100	500	$\lambda=24$	2183	1304	0.35
		$\lambda_b=5710$	116	158	5.9
		$\lambda_c=0.1$	5085	4480	0.17
100	1000	$\lambda=14$	3430	1899	0.20
		$\lambda_b=7691$	186	240	4.3
		$\lambda_c=0.1$	5102	4625	0.15

The Results on fMRI data

- Test on 2007 PBAIC competition
- Filtering with correlation
- comparison with the Elastic Net Results on correlation of the estimated response with the true response on the test data
 - On the 24th response $c=0.82$ with 127 pre-selected voxels
Elastic Net: $c=0.69$ with 300 pre-selected voxels
 - On the 15th response $c=0.74$ with 126 pre-selected voxels
Elastic Net: $c=0.66$ with 300 pre-selected voxels

On-going Work

- Different λ for each variable in the network
- Learning appropriate prior for λ
- Alternative penalties for the likelihood
- Application of non-convex optimization methods : EBW