Extended Bayesian scores for reconstructing gene regulatory networks.

Jimmy Vandel

Simon de Givry & Brigitte Mangin



ECCS'10 September 16th, 2010



OUTLINES :

- \rightarrow Biological background
- \rightarrow State of the art
- \rightarrow Score-based structure learning
- \rightarrow Model with extra biological knowledge
- \rightarrow Results
- \rightarrow Conclusion and perspectives





 \rightarrow gene activities (mRNA concentrations)



- \rightarrow gene activities (mRNA concentrations)
- \rightarrow gene regulations



- \rightarrow gene activities (mRNA concentrations)
- \rightarrow gene regulations

<u>Goal</u> : Reconstruction of gene regulatory network.

Escherichia coli (423 genes 578 regulations)













 \rightarrow DNA mutations in genes - in promoter region (impact on gene activity)





→ DNA mutations in genes - in promoter region (impact on gene activity)
 - in coding region (modify protein structure)





→ DNA mutations in genes - in promoter region (impact on gene activity)
 - in coding region (modify protein structure)

 \rightarrow observable through one genetic marker for each gene

Vandel Jimmy

1.Biological background

Pairwise correlation model

Correlations measurement + threshold

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

Pairwise correlation model

Correlations measurement + threshold

- → Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)
- \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

 \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Linear model

Graphical Gaussian Models (GGM)

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

 \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Linear model

Graphical Gaussian Models (GGM)

→ Global search (GeneNet Schäfer 2005, SIMoNe Chiquet 2008)

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

 \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Linear model

Graphical Gaussian Models (GGM)

→ Global search (GeneNet Schäfer 2005, SIMoNe Chiquet 2008)

 \rightarrow Local regressions (SEM Lasso Liu 2008)

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

 \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Linear model

Graphical Gaussian Models (GGM)

- → Global search (GeneNet Schäfer 2005, SIMoNe Chiquet 2008)
- \rightarrow Local regressions (SEM Lasso Liu 2008)
- → Boosting strategy (GGMSelect Giraud 2008)

(Spearman correlation+GGM)

Pairwise correlation model

Correlations measurement + threshold

→ Mutual Information (ARACNE Margolin 2006, CLR Faith 2007)

 \rightarrow Pearson / Spearman partial correlation (ParCorA de la Fuente 2004)

Linear model

Graphical Gaussian Models (GGM)

- → Global search (GeneNet Schäfer 2005, SIMoNe Chiquet 2008)
- \rightarrow Local regressions (SEM Lasso Liu 2008)
- → Boosting strategy (GGMSelect Giraud 2008)

(Spearman correlation+GGM)

Probabilistic discrete graphical model

 \rightarrow Bayesian networks on discrete data (*Friedman* 2000, *Zhu* 2007)

Bayesian network

* Directed acyclic graph G composed of n variables X_i with domain size r_i

* Conditional distribution for variable X_i , given its parents Pa_i in G:

$$P_G(X_i/Pa_i^j) = \theta_{ij}$$

* Representation of a joint probability distribution :

$$P_G(X) = \prod_{i=1}^n P_G(X_i / Pa_i)$$

* We note $Dim(G) = \sum_{i=1}^{n} (r_i - 1) * q_i$ the dimension of the network with $q_i = \prod_{Pa_i^j} r_j$

Vandel Jimmy

3.Score based learning

6/18

We look for the graph $G_{score} = argmax_{G_i} P(G_i/D)$ with dataset D.

We look for the graph $G_{score} = argmax_{G_i} P(G_i/D)$ with dataset D.

$$P(G_i/D) = \frac{P(D/G_i)P(G_i)}{P(D)}$$

We look for the graph $G_{score} = argmax_{G_i} P(G_i/D)$ with dataset D.

$$P(G_i/D) = \frac{P(D/G_i)P(G_i)}{P(D)}$$

\$\approx P(D/G_i)P(G_i)\$

We look for the graph $G_{score} = argmax_{G_i} P(G_i/D)$ with dataset D.

$$P(G_i/D) = \frac{P(D/G_i)P(G_i)}{P(D)}$$

 $\propto P(D/G_i)P(G_i)$

 P(D/G_i):marginal likelihood of Gi

 • exact under → Bayesian

 hypothesis

 Dirichlet score

 • estimation → Bayesian

 Information

 Criterion score

We look for the graph $G_{score} = argmax_{G_i} P(G_i/D)$ with dataset D.

$$P(G_i/D) = \frac{P(D/G_i)P(G_i)}{P(D)}$$

 $\propto P(D/G_i)P(G_i)$

P(D/G_i):marginal likelihood of Gi
 exact under → Bayesian
 hypothesis Dirichlet score
 estimation → Bayesian
 Information
 Criterion score

▷ $P(G_i)$:prior probability of the graph Gi
→ assumed to be uniform

> BD Score (D.Heckerman Machine learning 1995)

· A priori on conditional probability θ_{ij} following a Dirichlet distribution with parameter α_{ijk}

$$BD(G) = \prod_{i}^{n} \prod_{j}^{q_{i}} \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} \prod_{k}^{r_{i}} \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

- > BD Score (D.Heckerman Machine learning 1995)
 - · A priori on conditional probability θ_{ij} following a Dirichlet distribution with parameter α_{ijk}

$$BD(G) = \prod_{i}^{n} \prod_{j}^{q_{i}} \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} \prod_{k}^{r_{i}} \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

- BIC score (G.Schwartz Annals of statistics 1978)
 - · Laplace approximation

$$BIC(G) = \log(P(D/G, \hat{\theta})) - \frac{1}{2}Dim(G)\log(nb_{sample})$$

- > BD Score (**D.Heckerman** Machine learning 1995)
 - · A priori on conditional probability θ_{ij} following a Dirichlet distribution with parameter α_{ijk}

$$BD(G) = \prod_{i}^{n} \prod_{j}^{q_{i}} \frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} \prod_{k}^{r_{i}} \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})}$$

- BIC score (G.Schwartz Annals of statistics 1978)
 - · Laplace approximation

$$BIC(G) = \log(P(D/G, \hat{\theta})) - \frac{1}{2}Dim(G)\log(nb_{sample})$$

- fNML score (*T.Silander* International journal of approximate reasoning 2010) (factorized Normalized Maximum Likelihood)
 - · How much G explain D compare to other possible data sets D'

 $fNML(G) = \log(P(D/G, \hat{\theta})) - \sum_{i=1}^{n} \sum_{j=1}^{q_i} \log C_{N_{ij}}^{r_i} \quad with \quad C_{N_{ij}}^{r_i} the normalizing sum$

- > What about $P(G_i)$?
 - · Uniform for all G_i

- > What about $P(G_i)$?
 - · Uniform for all G_i

- > What about $P(G_i)$?
 - · Uniform for all G_i

How many restricted graph to node *j* with $|Pa_j^{G_i^j}| = l$ exist ?

 $\tau_l(G_i^j) = 1 \qquad \text{for } l = 0$

- > What about $P(G_i)$?
 - · Uniform for all G_i

$$\tau_l(G_i^j) = 1 \qquad \text{for } l = 0$$
$$= (n-1) \qquad l = 1$$

- > What about $P(G_i)$?
 - · Uniform for all G_i

$$\tau_l(G_i^J) = 1 \qquad \text{for } l = 0$$

= $(n-1)$ $l = 1$
= $(n-1)*(\frac{n-2}{2})$ $l = 2$

- > What about $P(G_i)$?
 - · Uniform for all G_i

$$\tau_l(G_i^J) = 1 \qquad \text{for } l = 0$$

$$= (n-1) \qquad l = 1$$

$$= \binom{n-1}{p} \qquad l = 2$$

$$\vdots$$

$$l = p$$

- > What about $P(G_i)$?
 - · Uniform for all G_i

$$\begin{aligned} \tau_l(G_i^j) &= 1 & \text{for } l = 0 \\ &= (n-1) & l = 1 \\ &= (n-1) \ast \left(\frac{n-2}{2}\right) & l = 2 \\ &= \begin{bmatrix} n-1 \\ p \end{bmatrix} & l = p \end{aligned}$$
 For $l < \frac{n-1}{2} \\ \longrightarrow \tau_{l-1}(G_i^j) < \tau_l(G_i^j) \end{aligned}$

- > What about $P(G_i)$?
 - · Uniform for all G_i

How many restricted graph to node *j* with $|Pa_j^{G_i^j}| = l$ exist ?

$$\begin{aligned} \tau_l(G_i^j) &= 1 & \text{for } l = 0 \\ &= (n-1) & l = 1 \\ &= (n-1) * \left(\frac{n-2}{2}\right) & l = 2 \\ &= \binom{n-1}{p} & l = p \end{aligned}$$
 For $l < \frac{n-1}{2} \\ & \longrightarrow \tau_{l-1}(G_i^j) < \tau_l(G_i^j) \\ P_{l-1}(G_i^j) < P_l(G_i^j) \end{aligned}$

Against principle of parsimony?

3.Score based learning

- > What about $P(G_i)$?
 - · Uniform for all G_i

How many restricted graph to node *j* with $|Pa_j^{G_i^j}| = l$ exist?

$$\tau_l(G_i^J) = 1 \qquad \text{for } l = 0$$

$$= (n-1) \qquad l = 1$$

$$= \binom{n-1}{p} \qquad l = 2$$

$$\vdots$$

$$l = p$$

For
$$l < \frac{n-1}{2}$$

 $\rightarrow \tau_{l-1}(G_i^j) < \tau_l(G_i^j)$
 $P_{l-1}(G_i^j) < P_l(G_i^j)$

Against principle of parsimony ?

· Extend BIC idea (Chen Biometrika 2008)

$$P(G_i^j) \propto \tau_l^{-\gamma}(G_i^j) \quad with \quad \gamma \in [0,1]$$
$$\log (P(G_i)) \simeq -\gamma \log (\prod_{j=1}^n \tau_l(G_i^j))$$

 \rightarrow applicable for all scores seen previously

Full observed data



Full observed data One marker per gene



Full observed data One marker per gene Discrete expression data



Full observed data One marker per gene Discrete expression data SNP markers



Full observed data One marker per gene Discrete expression data SNP markers



 \rightarrow Genetic linkage between markers

Full observed data One marker per gene Discrete expression data SNP markers



- \rightarrow Genetic linkage between markers
- \rightarrow Mutation in promoter region of G_i (ex G₁ and G₃)
 - · Force link $M_i \rightarrow G_i$
 - Forbid links $M_i \rightarrow G_j \qquad \forall j \neq i$

Full observed data One marker per gene Discrete expression data SNP markers



 \rightarrow Genetic linkage between markers

- \rightarrow Mutation in promoter region of G_i (ex G₁ and G₃)
 - Force link $M_i \rightarrow G_i$
 - Forbid links $M_i \rightarrow G_j \qquad \forall j \neq i$

Full observed data One marker per gene Discrete expression data SNP markers



 \rightarrow Genetic linkage between markers

- \rightarrow Mutation in promoter region of G_i (ex G₁ and G₃)
 - · Force link $M_i \rightarrow G_i$
 - Forbid links $M_i \to G_j \qquad \forall j \neq i$
- \rightarrow Mutation in coding region for G_i (ex G₂)
 - · Forbid link $M_i \rightarrow G_i$

Full observed data One marker per gene Discrete expression data SNP markers



 \rightarrow Genetic linkage between markers

- \rightarrow Mutation in promoter region of G_i (ex G₁ and G₃)
 - · Force link $M_i \rightarrow G_i$
 - Forbid links $M_i \rightarrow G_j \qquad \forall j \neq i$
- \rightarrow Mutation in coding region for G_i (ex G₂)
 - · Forbid link $M_i \rightarrow G_i$

Full observed data One marker per gene Discrete expression data SNP markers



 \rightarrow Genetic linkage between markers

- \rightarrow Mutation in promoter region of G_i (ex G₁ and G₃)
 - · Force link $M_i \rightarrow G_i$
 - Forbid links $M_i \to G_j \qquad \forall j \neq i$
- \rightarrow Mutation in coding region for G_i (ex G₂)
 - · Forbid link $M_i \rightarrow G_i$

Network studied



- *similar structure of known networks few hubs*
- × 50 nodes / 50 edges



Network evaluation

- > Mean over 50 artificial networks
- Sample size between 50 and 500 individuals
- > In our model we project all $M_i \rightarrow G_j$ as a $G_i \rightarrow G_j$ $\forall j \neq i$
- We use Greedy search (implemented in Banjo Hartemink 2005)
- > Evaluation doesn't takes into account edges orientation
- > 2 common metrics Predictive value / Sensitivity

$$Predictive value: \frac{TP}{TP + FP}$$

Sensitivity: $\frac{TP}{TP+FN}$

TP : number of correct learned edges*FP* : number of wrong learned edges*FN* : number of missed edges

Extended scores impact

fNML score 1.0 1.0 0.8 0.8 0.6 0.6 Sensitivity Predictive 0.4 0.4 0.2 0.2 0.0 0.0 100 300 400 100 200 500 200 300 400 500 Sample size Sample size

- · We compare BIC and fNML scores with several \mathcal{Y} values
 - $\cdot \mathcal{Y} = 0 \rightarrow \text{classic scores (solid line)}$
 - $\cdot \mathcal{Y}$ = 0.5 (dashed line)

Legend: BIC score

 $\cdot \mathcal{Y} = 1 \rightarrow$ same probability over connectivity classes (dotted line)

Vandel Jimmy

5.Results

13/18

Biological knowledge impact

fNML score 1.0 1.0 0.8 0.8 0.6 0.6 Predictive Sensitivity 0.4 0.4 0.2 0.2 0.0 0.0 100 300 400 200 500 100 200 300 400 500 Sample size Sample size

 \cdot We fix / ban relations in regulation network inference with $\mathcal{Y}\text{=}1$

- without extra biological knowledge (solid line)
- with extra biological knowledge (dashed line)

Vandel Jimmy

Legend: BIC score

5.Results

14/18

Comparative results



Conclusion

- Scoring criteria study with uniform node in-degree prior
- Model description taking into account specific Biological knowledge
- Comparison with several regulatory network inference methods
- > Robustness of Bayesian networks for sparse graphs

Perspectives

- Improve learning algorithm performance
- Study causality
- Try on real data

Bibliography

- * **D.M. Chickering** 'Efficient Approximations for the Marginal Likelihood of Incomplete Data Given a Bayesian Network' 1997
- * J. Chen 'Extended Bayesian information criteria for model selection with large model spaces' 2008
- * **T.Silander** 'Learning locally minimax optimal Bayesian Networks' 2010
- * **N.Meinshausen** 'High dimensional graphs and variable selection with the lasso' 2006
- * **M.Bansal** 'How to infer gene networks from expression profiles' 2006
- * **B.Liu** 'Gene network inference via structural equation modeling in genetical genomics experiments' 2008
- *× J.Zhu* 'Increasing the power to detect causal associations by combining genotypic and expression data in segregating populations' 2007