

Algorithmic approaches to fitting ERG models

Ruth Hummel, Penn State University
Mark Handcock, University of Washington
David Hunter, Penn State University

Research funded by Office of Naval Research Award No. N00014-08-1-1015

MURI meeting, April 24, 2009

Outline

- 1 Introduction
- 2 Partial stepping
- 3 Biological network example
- 4 References

The class of Exponential-family Random Graph Models (ERGMs):

Definition

$$P_{\eta}(\mathbf{Y} = \mathbf{y}) = \frac{\exp\{\eta^t \mathbf{g}(\mathbf{y})\}}{\kappa(\eta)}$$

where

- \mathbf{Y} is a random network written as an adjacency matrix so that Y_{ij} is the indicator of an edge from i to j ;
- $\mathbf{g}(\mathbf{y})$ is a vector of the network statistics of interest;
- η is a vector of parameters corresponding to the vector $\mathbf{g}(\mathbf{y})$;
- $\kappa(\eta)$ is the constant of proportionality which makes the probabilities sum to one; **intractable**

Loglikelihood

The loglikelihood for this class of models is

$$l(\boldsymbol{\eta}) = \boldsymbol{\eta}^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log \sum_{\mathbf{z} \in \mathcal{Y}} \exp(\boldsymbol{\eta}^t \mathbf{g}(\mathbf{z})). \quad (1)$$

which can be written:

$$l(\boldsymbol{\eta}) - l(\boldsymbol{\eta}_0) = (\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log E_{\boldsymbol{\eta}_0} [\exp \{(\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{Y})\}], \quad (2)$$

Maximum Pseudolikelihood Estimation

Notation: For a network y and a pair (i, j) of nodes,

- $y_{ij} = 0$ or 1 , depending on whether there is an edge
- y_{ij}^c denotes the status of all pairs in y other than (i, j)
- y_{ij}^+ denotes the same network as y but with $y_{ij} = 1$
- y_{ij}^- denotes the same network as y but with $y_{ij} = 0$

- Assume no dependence among the Y_{ij} .
- In other words, assume $P(Y_{ij} = 1) = P(Y_{ij} = 1 | Y_{ij}^c = y_{ij}^c)$.
- Then some algebra gives

$$\log \frac{P(Y_{ij} = 1)}{P(Y_{ij} = 0)} = \theta^t \left[g(y_{ij}^+) - g(y_{ij}^-) \right],$$

so θ is estimated by straightforward logistic regression.

- Result: The **maximum pseudolikelihood estimate**.

Maximum Pseudolikelihood Estimation

Notation: For a network y and a pair (i, j) of nodes,

- $y_{ij} = 0$ or 1 , depending on whether there is an edge
 - y_{ij}^c denotes the status of all pairs in y other than (i, j)
 - y_{ij}^+ denotes the same network as y but with $y_{ij} = 1$
 - y_{ij}^- denotes the same network as y but with $y_{ij} = 0$
- Assume no dependence among the Y_{ij} .
 - In other words, assume $P(Y_{ij} = 1) = P(Y_{ij} = 1 | Y_{ij}^c = y_{ij}^c)$.
 - Then some algebra gives

$$\log \frac{P(Y_{ij} = 1)}{P(Y_{ij} = 0)} = \theta^t \left[g(y_{ij}^+) - g(y_{ij}^-) \right],$$

so θ is estimated by straightforward logistic regression.

- Result: The **maximum pseudolikelihood estimate**.

MPLE's behavior

- **MLE (maximum likelihood estimation)**: Well-established method but very hard because the normalizing constant $\kappa(\alpha)$ is difficult (usually impossible) to evaluate, so we approximate it instead.
- **MPLE (maximum pseudo-likelihood estimation)**: Easy to do using logistic regression, but based on an independence assumption that is often not justified.

Several authors, notably van Duijn et al. (2009), argue forcefully against the use of MPLE.

Back to the loglikelihood

- Remember that the loglikelihood

$$l(\boldsymbol{\eta}) = \boldsymbol{\eta}^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log \sum_{z \in \mathcal{Y}} \exp(\boldsymbol{\eta}^t \mathbf{g}(z)). \quad (3)$$

can be written:

$$l(\boldsymbol{\eta}) - l(\boldsymbol{\eta}_0) = (\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log E_{\boldsymbol{\eta}_0} [\exp \{(\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{Y})\}],$$

This leads to our first approximation for the MLE:

- MCMC MLE idea: Pick $\boldsymbol{\theta}_0$, draw Y_1, \dots, Y_m from this model using MCMC, then approximate the population mean above by a sample mean.
- We can take $\boldsymbol{\theta}_0$ to be, for example, the MPLE.

Back to the loglikelihood

- Remember that the loglikelihood

$$l(\boldsymbol{\eta}) = \boldsymbol{\eta}^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log \sum_{z \in \mathcal{Y}} \exp(\boldsymbol{\eta}^t \mathbf{g}(z)). \quad (3)$$

can be written:

$$l(\boldsymbol{\eta}) - l(\boldsymbol{\eta}_0) = (\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log E_{\boldsymbol{\eta}_0} [\exp \{(\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{Y})\}],$$

This leads to our first approximation for the MLE:

- MCMC MLE idea: Pick θ_0 , draw Y_1, \dots, Y_m from this model using MCMC, then approximate the population mean above by a sample mean.
- We can take θ_0 to be, for example, the MPLE.

Back to the loglikelihood

- Remember that the loglikelihood

$$l(\boldsymbol{\eta}) = \boldsymbol{\eta}^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log \sum_{\mathbf{z} \in \mathcal{Y}} \exp(\boldsymbol{\eta}^t \mathbf{g}(\mathbf{z})). \quad (3)$$

can be written:

$$l(\boldsymbol{\eta}) - l(\boldsymbol{\eta}_0) = (\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{y}_{\text{obs}}) - \log E_{\boldsymbol{\eta}_0} [\exp \{(\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{Y})\}],$$

This leads to our first approximation for the MLE:

- MCMC MLE idea: Pick $\boldsymbol{\theta}_0$, draw Y_1, \dots, Y_m from this model using MCMC, then approximate the population mean above by a sample mean.
- We can take $\boldsymbol{\theta}_0$ to be, for example, the MPLE.

Existence of MLE?

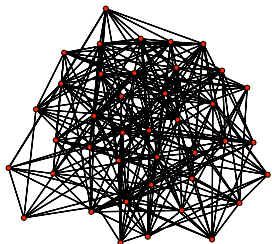
There are two issues we need to be concerned with when using MCMC MLE:

- First, if the observed $\mathbf{g}(\mathbf{y}_{obs})$ is not in the interior of the convex hull of the sampled (MCMC generated) $\mathbf{g}(\mathbf{y}_i)$, then no maximizer of the approximate loglikelihood ratio exists.
- Also, the approximation of the likelihood surface,

$$l(\boldsymbol{\eta}) - l(\boldsymbol{\eta}_0) \approx (\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{y}_{obs}) - \log \frac{1}{m} \sum_{i=1}^m \exp((\boldsymbol{\eta} - \boldsymbol{\eta}_0)^t \mathbf{g}(\mathbf{Y}_i)),$$

is not very good when we get far from $\boldsymbol{\eta}_0$.

Quick Erdős-Rényi illustration



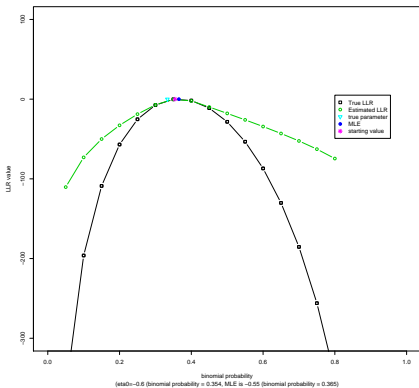
The Erdős-Rényi model can be written as an ERGM, if $g(\mathbf{y})$ is the number of edges of \mathbf{y} .

Since each edge exists independently with probability b ,

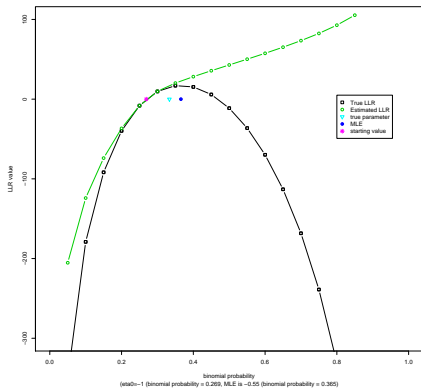
$$P(\mathbf{Y} = \mathbf{y}) = (b)^{g(\mathbf{y})} (1 - b)^{n^* - g(\mathbf{y})} = \left(\frac{b}{1 - b} \right)^{g(\mathbf{y})} (1 - b)^{n^*}.$$

Comparison of the true loglikelihood ratio and its estimation for an Erdos-Renyi graph

True and estimated loglikelihood ratios plotted for $\eta = -0.69$ (binomial probability = 1/3)



True and estimated loglikelihood ratios plotted for $\eta = -0.69$ (binomial probability = 1/3)



Outline

- 1 Introduction
- 2 Partial stepping**
- 3 Biological network example
- 4 References

Mean value parametrization

- 1 We would like to move our arbitrary initial value η_0 close enough to the MLE that samples (of the statistics) generated from η_0 cover the observed statistics, $\mathbf{g}(\mathbf{y}_{obs})$.
- 2 New Idea: Intuitively, then, we might try to find a way to make the journey from η_0 to the MLE in pieces.

In order to take steps toward the MLE, we need to have some idea where we are going. We obviously don't know where the MLE is, but we do know that the MLE has the following property:

Definition

The MLE, if it exists, is the unique parameter vector η satisfying $E_{\eta} \mathbf{g}(\mathbf{Y}) = \mathbf{g}(\mathbf{y}_{obs}) = \hat{\xi}$.

This is true for any exponential family (Barndorff-Nielsen, 1978, and Brown, 1986).

Mean value parametrization

- 1 We would like to move our arbitrary initial value η_0 close enough to the MLE that samples (of the statistics) generated from η_0 cover the observed statistics, $\mathbf{g}(\mathbf{y}_{obs})$.
- 2 New Idea: Intuitively, then, we might try to find a way to make the journey from η_0 to the MLE in pieces.

In order to take steps toward the MLE, we need to have some idea where we are going. We obviously don't know where the MLE is, but we do know that the MLE has the following property:

Definition

The MLE, if it exists, is the unique parameter vector η satisfying $E_{\eta} \mathbf{g}(\mathbf{Y}) = \mathbf{g}(\mathbf{y}_{obs}) = \hat{\xi}$.

This is true for any exponential family (Barndorff-Nielsen, 1978, and Brown, 1986).

Mean value parametrization

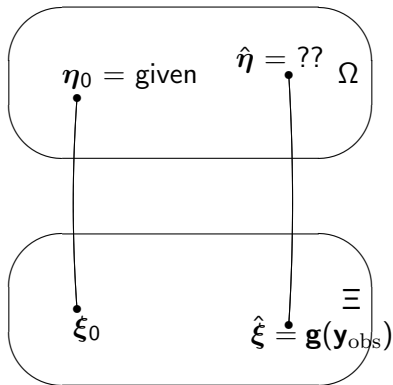
- 1 We would like to move our arbitrary initial value η_0 close enough to the MLE that samples (of the statistics) generated from η_0 cover the observed statistics, $\mathbf{g}(\mathbf{y}_{obs})$.
- 2 New Idea: Intuitively, then, we might try to find a way to make the journey from η_0 to the MLE in pieces.

In order to take steps toward the MLE, we need to have some idea where we are going. We obviously don't know where the MLE is, but we do know that the MLE has the following property:

Definition

The MLE, if it exists, is the unique parameter vector η satisfying $E_{\eta}\mathbf{g}(\mathbf{Y}) = \mathbf{g}(\mathbf{y}_{obs}) = \hat{\xi}$.

This is true for any exponential family (Barndorff-Nielson, 1978, and Brown, 1986).



- Ω is the original η parameter space (defined for the change statistics)
- Ξ is the corresponding mean value parameter space
- η_0 is the (given) initial value of η in the Markov Chain
- ξ_0 is the vector of mean statistics corresponding to η_0
- $\hat{\eta}$ is the (unknown) MLE
- $\hat{\xi}$ is the observed statistics, $\mathbf{g}(\mathbf{y}_{\text{obs}})$, which is the corresponding mean statistics vector for $\hat{\eta}$.

Taking steps

- We now want to move toward $\hat{\xi} = \mathbf{g}(\mathbf{y}_{\text{obs}})$ in mean value parameter space. Here we specify a step length, $0 \leq \gamma \leq 1$, as a fraction of the distance toward $\hat{\xi}$ that we want to traverse. We move the fraction γ toward $\hat{\xi}$ and call this point $\xi_1 = \gamma\hat{\xi} + (1 - \gamma)\xi_0$.
- At each step, we choose this fraction to be the biggest move toward $\mathbf{g}(\mathbf{y}_{\text{obs}})$ that does not leave the convex hull of the current sample.

Taking steps

- We now want to move toward $\hat{\xi} = \mathbf{g}(\mathbf{y}_{\text{obs}})$ in mean value parameter space. Here we specify a step length, $0 \leq \gamma \leq 1$, as a fraction of the distance toward $\hat{\xi}$ that we want to traverse. We move the fraction γ toward $\hat{\xi}$ and call this point $\xi_1 = \gamma \hat{\xi} + (1 - \gamma) \xi_0$.
- At each step, we choose this fraction to be the biggest move toward $\mathbf{g}(\mathbf{y}_{\text{obs}})$ that does not leave the convex hull of the current sample.

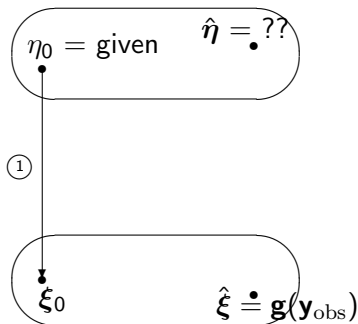
Finding the MCMC MLE

Next, using the sample from the model defined by η_0 , we maximize the approximate loglikelihood

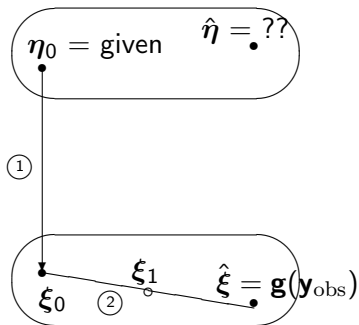
$$l(\eta) - l(\eta_0) \approx (\eta - \eta_0)^t \mathbf{g}(\mathbf{y}_{obs}) - \log \frac{1}{m} \sum_{i=1}^m \exp((\eta - \eta_0)^t \mathbf{g}(\mathbf{Y}_i)), \quad (6)$$

but with ξ_1 substituted in place of $\mathbf{g}(\mathbf{y}_{obs})$. The resulting maximizer will be called η_1 . The process then repeats, with η_1 taking the place of η_0 .

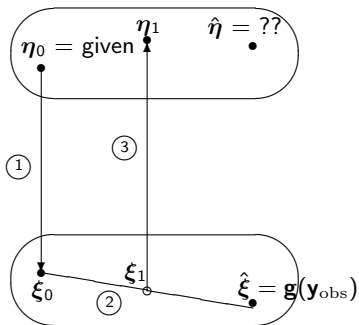
Partial stepping in Mean Value Parameter space



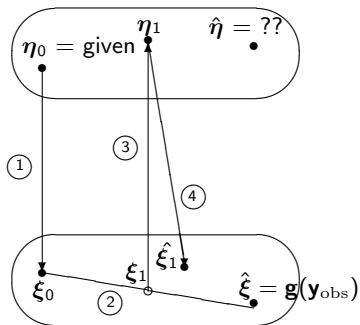
0. Set η_0 .
1. Take an MCMC sample from the model defined by $\eta = \eta_0$ to get ξ_0 .



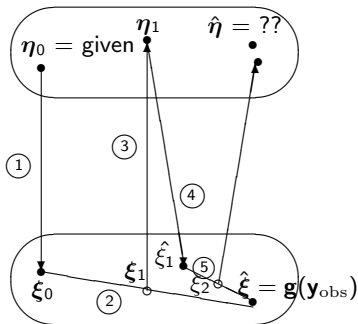
2. Go $\gamma\%$ toward $\hat{\xi} = \mathbf{g}(\mathbf{y}_{\text{obs}})$ in mean value parameter space. Call this ξ_1 .



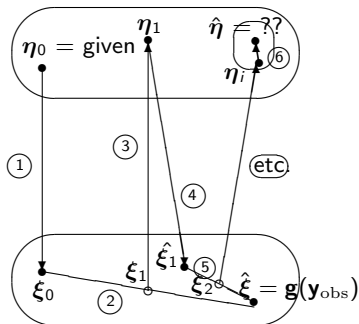
3. Use η_0 in MCMC MLE to find the η that corresponds to ξ_1 . Call this η_1 .



4. Re-estimate ξ_1 from an MCMC sample from the model defined by $\eta = \eta_1$. Call this $\hat{\xi}_1$.



- Repeat step (2) by going $\gamma\%$ toward $\hat{\xi} = \mathbf{g}(\mathbf{y}_{obs})$ from $\hat{\xi}_1$. Call this ξ_2 . Also repeat steps (3) and (4) to obtain $\hat{\xi}_2$. Keep going until $\mathbf{g}(\mathbf{y}_{obs})$ is in the convex hull of the new sample.



6. Use MCMC MLE with η_i as the initial value to find $\hat{\eta}$. If we have made it into the appropriate neighborhood around $\hat{\eta}$, this will now be possible.

In review:

In other words, for each $t \geq 0$, we first use MCMC to draw a random sample $\mathbf{Y}_1, \dots, \mathbf{Y}_m$ from the model determined by η_t , then we set

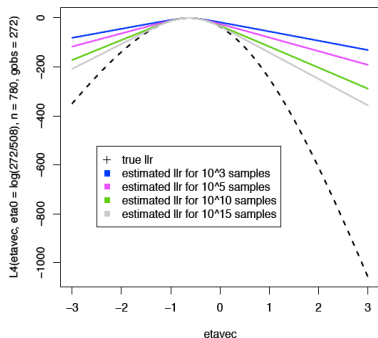
$$\hat{\xi}_t = \frac{1}{m} \sum_{i=1}^m \mathbf{g}(\mathbf{Y}_i);$$

$$\xi_{t+1} = \gamma_t \hat{\xi} + (1 - \gamma_t) \xi_t;$$

$$\eta_{t+1} = \arg \max_{\eta} \left\{ (\eta - \eta_0)^t \xi_{t+1} - \log \left[\frac{1}{m} \sum_{i=1}^m \exp \{ (\eta - \eta_0)^t \mathbf{g}(\mathbf{y}_i) \} \right] \right\}$$

We iterate until $\hat{\xi} = \mathbf{g}(\mathbf{y}_{\text{obs}})$ is in the convex hull of the statistics generated from $\hat{\xi}_t$.

A second approximation: the lognormal



- Here, the ERGM is very simple, $g(y) = \text{edges}$.
- In other words, the model is binomial.

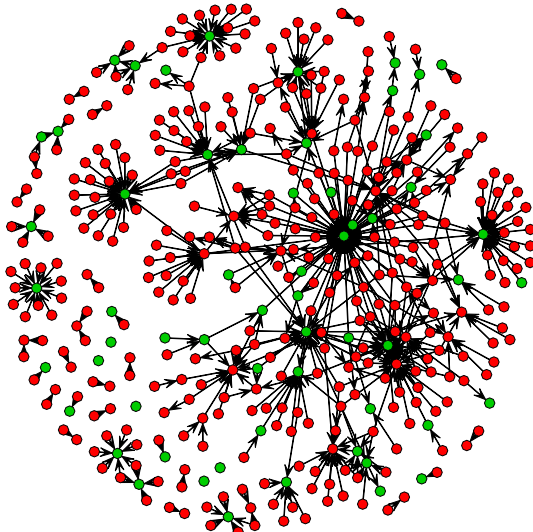
- Nevertheless, naive MCMC approximation of the log likelihood ratio is not good far from θ_0 , even for gigantic samples.
- One possible remedy: Assume $(\theta - \theta_0)^t g(Y)$ is normally distributed in

$$\ell(\theta) - \ell(\theta_0) = (\theta - \theta_0)^t g(y^{\text{obs}}) - \log \mathbb{E}_{\theta_0} \left[\exp \left\{ (\theta - \theta_0)^t g(Y) \right\} \right].$$

Outline

- 1 Introduction
- 2 Partial stepping
- 3 Biological network example**
- 4 References

Example Bionet: E. Coli (Salgado et al 2001)



- A node is an operon
- Edge $A \rightarrow B$ means A encodes a transcription factor that regulates B .
- Green indicates self-regulation

Another ERGM for the E. Coli network

We fit a model similar to that of Saul and Filkov (2007)

Term(s)	Description:
Edges	Number of Edges
2-Deg, ..., 5-Deg	Nodes with degree 2, ..., 5
GWDeg	Single statistic: Weighted sum of 1-Deg, ..., $(n - 1)$ -Deg with weights tending to 1 at a geometric rate

```
model <- ergm(ecoli2 ~ edges + degree(2:5) +
  gwdegree(0.25, fixed=TRUE), MPLEonly=TRUE)
```

MPLE

```
edges    degree2    degree3    degree4    degree5    gwdegree
-5.35    -2.58      -3.06      -2.39      -1.85      8.13
```


Another ERGM for the E. Coli network

We fit a model similar to that of Saul and Filkov (2007)

Term(s)	Description:
Edges	Number of Edges
2-Deg, ..., 5-Deg	Nodes with degree 2, ..., 5
GWDeg	Single statistic: Weighted sum of 1-Deg, ..., $(n - 1)$ -Deg with weights tending to 1 at a geometric rate

```
model <- ergm(ecoli2 ~ edges + degree(2:5) +
  gwdegree(0.25, fixed=TRUE), MPLEonly=TRUE)
```

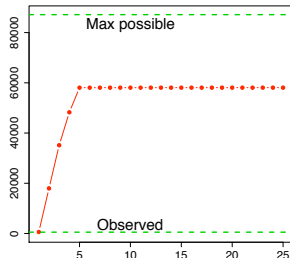
MPLE

edges	degree2	degree3	degree4	degree5	gwdegree
-5.35	-2.58	-3.06	-2.39	-1.85	8.13

MPLE fit is degenerate

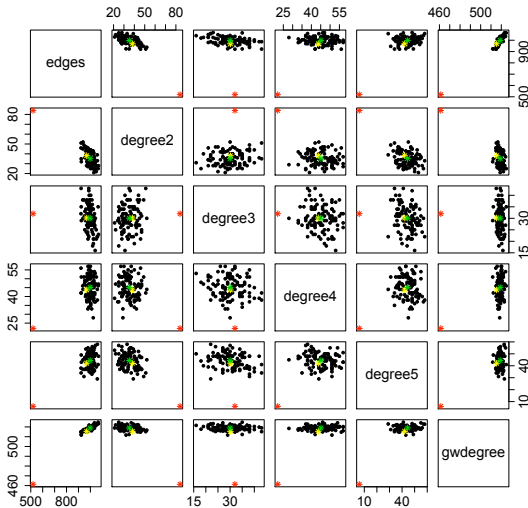
With the MPLE we encounter problems:

Here is a time-series plot of the edge-count of the 25 networks generated from the MPLE:



A sample from the MPLE-fitted model

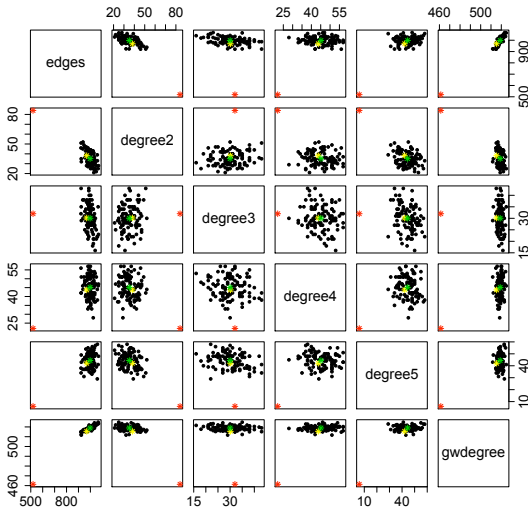
Iteration 1: Yellow = 0.06(Red) + 0.94(Green)



- Red: $g(y^{obs})$
- Green: Sample mean
- Theory: No maximizer of the approximated likelihood exists because $g(y^{obs})$ is not in the interior of the convex hull of the sampled points.
- However, the likelihood depends on the data only through $g(y^{obs})$
- Idea: What if we pretend $g(y^{obs})$ is the yellow point?

A sample from the MPLE-fitted model

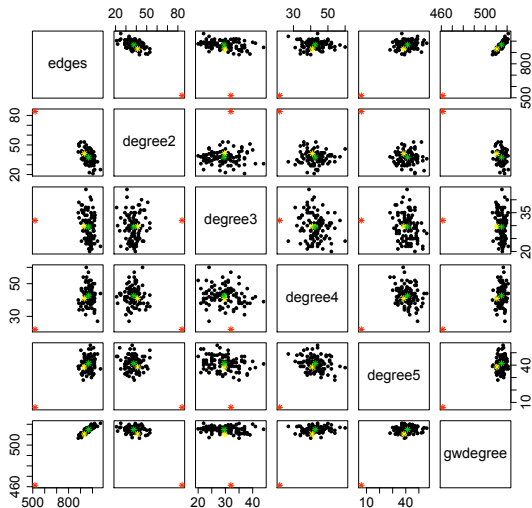
Iteration 1: Yellow = 0.06(Red) + 0.94(Green)



- Red: $g(y^{obs})$
- Green: Sample mean
- Theory: No maximizer of the approximated likelihood exists because $g(y^{obs})$ is not in the interior of the convex hull of the sampled points.
- However, the likelihood depends on the data only through $g(y^{obs})$
- Idea: What if we pretend $g(y^{obs})$ is the yellow point?

A sample from a model with a better θ_0

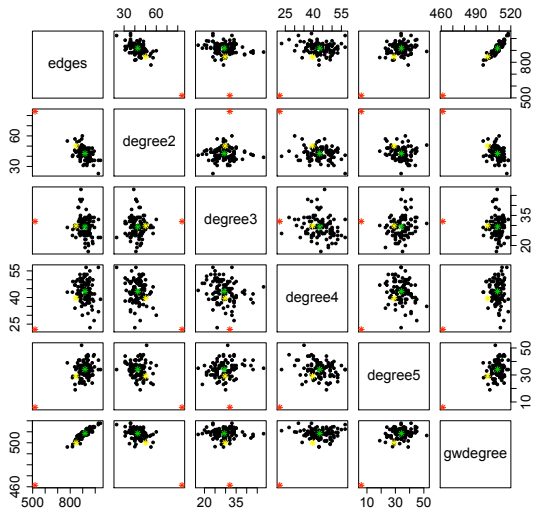
Iteration 2: Yellow = 0.08(Red) + 0.92(Green)



- For θ_0 we have replaced the MPLE by the MCMC “MLE” obtained by pretending that $g(y^{\text{obs}})$ was the yellow point on the previous slide.

A sample from a model with a better θ_0

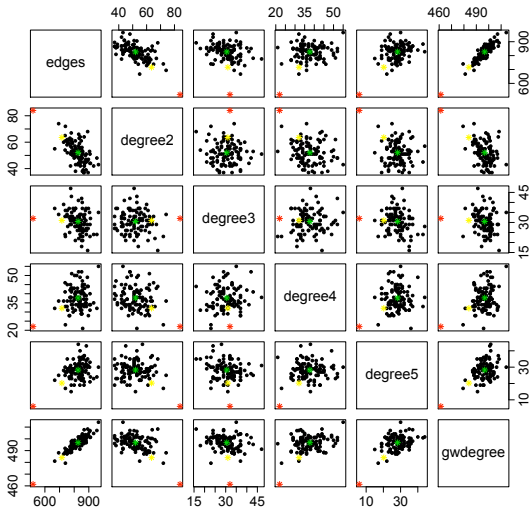
Iteration 3: Yellow = 0.18(Red) + 0.82(Green)



• Continue to iterate...

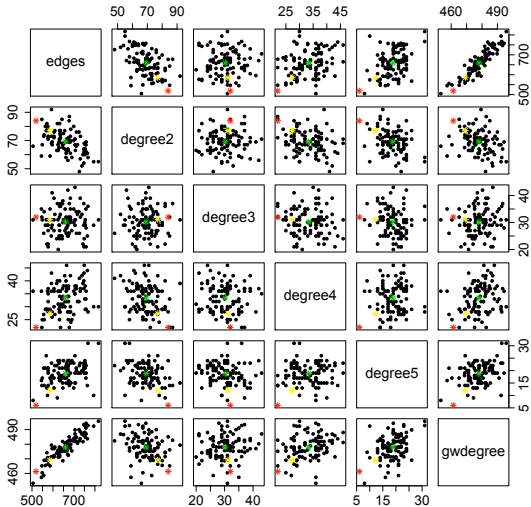
A sample from a model with a better θ_0

Iteration 4: Yellow = 0.36(Red) + 0.64(Green)



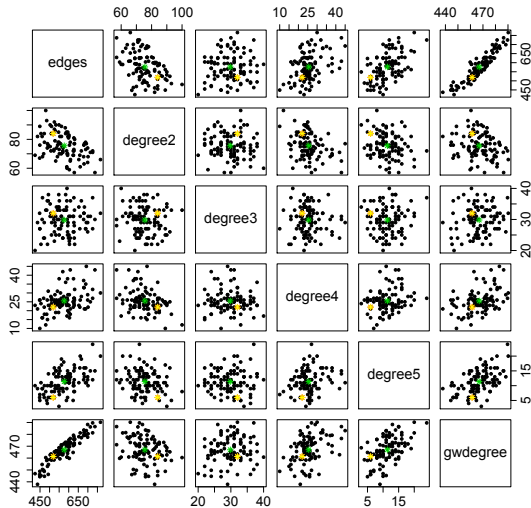
A sample from a model with a better θ_0

Iteration 5: Yellow = 0.53(Red) + 0.47(Green)



A sample from a model with a better θ_0

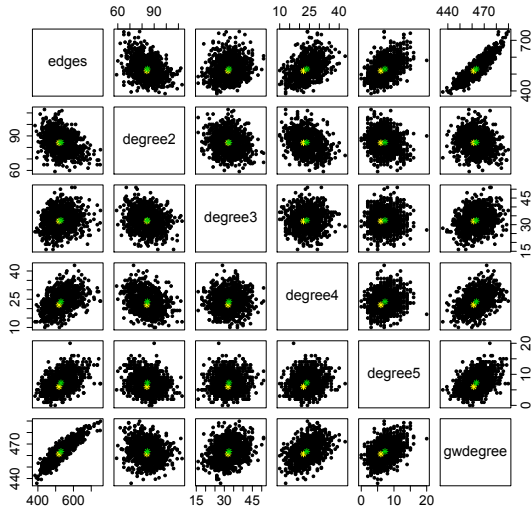
Iteration 6: Yellow = 1(Red) + 0(Green)



- Finally, we don't need to pretend; the true $g(y^{\text{obs}})$ is actually interior to the convex hull of sampled points. . .

A sample from a model with a better θ_0

Final Iteration (#7): Green = mean, Yellow=observed



- ... so now we can take a larger sample and get a better final estimate of θ .

Finally, an MLE

Original MPLE:

edges	degree2	degree3	degree4	degree5	gwdegree
-5.35	-2.58	-3.06	-2.39	-1.85	8.13

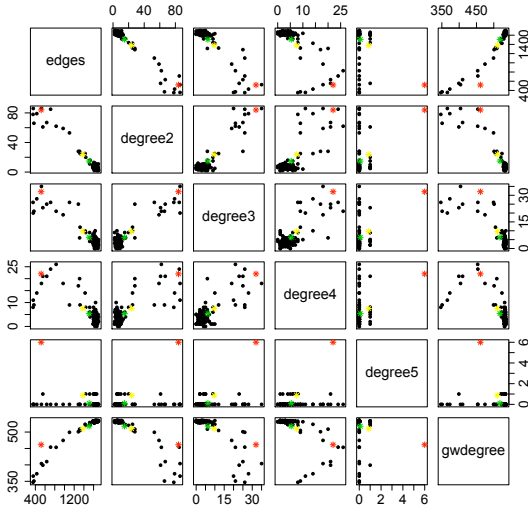
Final (approximated) MLE:

edges	degree2	degree3	degree4	degree5	gwdegree
-5.06	-1.45	-2.35	-2.28	-2.91	1.77

We know this is close to the true MLE $\hat{\theta}$ because the true MLE uniquely gives $E_{\hat{\theta}}g(Y) = g(y^{\text{obs}})$.

Failure of first approximation in this example

Iteration 5: Yellow = 0.13(Red) + 0.87(Green)



- A steplength of 0.01 is too small in this case.

Conclusions

- MPLE looks increasingly dangerous; it can mask problems when they exist and miss badly when they don't
- Naive MCMC MLE may not perform well even in very simple problems, but it may be modified. Here, we had success in a hard problem using two ideas:
 - (a) Partial stepping toward the MLE in mean-value parameter space;
 - (b) A log-normal approximation to the normalizing constant.
- By making MLE more and more automatic, we hope that scientists will be able to focus on modeling, not programming.

Outline

- 1 Introduction
- 2 Partial stepping
- 3 Biological network example
- 4 References

Cited References

- Alon, U. (2007, *Nature Reviews Genetics*), Network Motifs: Theory and Experimental Approaches.
- Barndorff-Nielsen, O. (1978). "Information and exponential families in statistical theory," New York: Wiley series in probability and mathematical statistics.
- Brown, L. D. (1986). "Fundamentals of statistical Exponential Families," Hayward: Institute of Mathematical Statistics.
- Saul ZM and Filkov V (2007, *Bioinformatics*), Exploring biological network structure using exponential random graph models.
- Salgado H et al. (2001, *Nucleic Acids Res.*), Regulondb (version 3.2): Transcriptional regulation and operon organization in Escherichia Coli k-12.
- van Duijn MAJ, Gile K, and Handcock MS (2009, *Social Networks*), A Framework for the Comparison of Maximum Pseudo-Likelihood and Maximum Likelihood Estimation of Exponential Family Random Graph Models.

Thank You!