Sampling and incomplete network data 567 Statistical analysis of social networks

Peter Hoff

Statistics, University of Washington

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

It is sometimes difficult to obtain a complete network dataset:

- the population nodeset is too large;
- gathering all relational information is too costly;
- population nodes are hard to reach.

- gather the data (i.e. design the survey);
- make inference (i.e. estimate and evaluate parameters).

Common sampling methods

- 1. node-induced subgraph sampling
- 2. edge-induced subgraph sampling
- 3. egocentric sampling
- 4. link tracing designs
- 5. censored nomination schemes

Procedure:

1. Uniformly sample a set $\mathbf{s} = \{s_1, \ldots, s_{n_s}\}$ of nodes

 $\mathbf{s} \subset \{1, \ldots, n\}.$

2. Observe relations y_s between sampled nodes

$$\mathbf{Y}_{\mathbf{s}} = \{ y_{i,j} : i \in \mathbf{s}, j \in \mathbf{s} \}.$$

Procedure:

1. Uniformly sample a set $\mathbf{s} = \{s_1, \ldots, s_{n_s}\}$ of nodes

 $\mathbf{s} \subset \{1, \ldots, n\}.$

2. Observe relations \boldsymbol{y}_{s} between sampled nodes

$$\mathbf{Y}_{\mathbf{s}} = \{ y_{i,j} : i \in \mathbf{s}, j \in \mathbf{s} \}.$$

Procedure:

1. Uniformly sample a set $\mathbf{s} = \{s_1, \ldots, s_{n_s}\}$ of nodes

 $\mathbf{s} \subset \{1, \ldots, n\}.$

2. Observe relations \mathbf{y}_s between sampled nodes

$$\mathbf{Y}_{\mathbf{s}} = \{ y_{i,j} : i \in \mathbf{s}, j \in \mathbf{s} \}.$$

Procedure:

1. Uniformly sample a set $\mathbf{s} = \{s_1, \ldots, s_{n_s}\}$ of nodes

 $\mathbf{s} \subset \{1, \ldots, n\}.$

2. Observe relations \mathbf{y}_s between sampled nodes

$$\mathbf{Y}_{\mathbf{s}} = \{ y_{i,j} : i \in \mathbf{s}, j \in \mathbf{s} \}.$$













In what ways does \mathbf{Y}_s resemble \mathbf{Y} ?

For what functions g() will $g(\mathbf{Y}_s)$ estimate $g(\mathbf{Y})$?

Consider the following setup:

- *n* × *n* sociomatrix **Y**
- *n* × *n* dyadic covariate **X**_d
- *n* × 1 nodal covariate **X**_{*n*}

$$\begin{split} \bar{y} &= \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} \ , \ \bar{x}_d = \frac{1}{n(n-1)} \sum_{i \neq j} x_{d,i,j} \ , \ \bar{x}_n = \frac{1}{n} \sum x_{n,i} \\ \overline{yx_d} &= \frac{1}{n(n-1)} \sum y_{i,j} x_{d,i,j} \ , \ \overline{yx_n} = \frac{1}{n(n-1)} \sum x_{n,i} \overline{y_i}. \end{split}$$

In what ways does \mathbf{Y}_s resemble \mathbf{Y} ? For what functions g() will $g(\mathbf{Y}_s)$ estimate $g(\mathbf{Y})$?

Consider the following setup:

- *n* × *n* sociomatrix **Y**
- *n* × *n* dyadic covariate **X**_d
- $n \times 1$ nodal covariate X_n

$$\begin{split} \bar{y} &= \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} \ , \ \bar{x}_d = \frac{1}{n(n-1)} \sum_{i \neq j} x_{d,i,j} \ , \ \bar{x}_n = \frac{1}{n} \sum x_{n,i} \\ \overline{yx_d} &= \frac{1}{n(n-1)} \sum_{i \neq i} y_{i,j} x_{d,i,j} \ , \ \overline{yx_n} = \frac{1}{n(n-1)} \sum_i x_{n,i} \bar{y}_i. \end{split}$$

In what ways does \mathbf{Y}_s resemble \mathbf{Y} ? For what functions g() will $g(\mathbf{Y}_s)$ estimate $g(\mathbf{Y})$?

Consider the following setup:

- *n* × *n* sociomatrix **Y**
- *n* × *n* dyadic covariate **X**_d
- $n \times 1$ nodal covariate X_n

$$\begin{split} \bar{y} &= \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} \ , \ \bar{x}_d = \frac{1}{n(n-1)} \sum_{i \neq j} x_{d,i,j} \ , \ \bar{x}_n = \frac{1}{n} \sum x_{n,i} \\ \overline{yx_d} &= \frac{1}{n(n-1)} \sum_{i \neq i} y_{i,j} x_{d,i,j} \ , \ \overline{yx_n} = \frac{1}{n(n-1)} \sum_i x_{n,i} \bar{y}_i. \end{split}$$

In what ways does \mathbf{Y}_{s} resemble \mathbf{Y} ?

For what functions g() will $g(\mathbf{Y}_s)$ estimate $g(\mathbf{Y})$?

Consider the following setup:

- *n* × *n* sociomatrix **Y**
- *n* × *n* dyadic covariate X_d
- *n* × 1 nodal covariate **X**_{*n*}

$$\bar{y} = \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} \ , \ \bar{x}_d = \frac{1}{n(n-1)} \sum_{i \neq j} x_{d,i,j} \ , \ \bar{x}_n = \frac{1}{n} \sum x_{n,i}$$
$$\overline{yx_d} = \frac{1}{n(n-1)} \sum y_{i,j} x_{d,i,j} \ , \ \overline{yx_n} = \frac{1}{n(n-1)} \sum x_{n,i} \bar{y}_i.$$

In what ways does \mathbf{Y}_s resemble \mathbf{Y} ?

For what functions g() will $g(\mathbf{Y}_s)$ estimate $g(\mathbf{Y})$?

Consider the following setup:

- *n* × *n* sociomatrix **Y**
- *n* × *n* dyadic covariate X_d
- *n* × 1 nodal covariate **X**_{*n*}

$$\begin{split} \bar{y} &= \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} \ , \ \bar{x}_d = \frac{1}{n(n-1)} \sum_{i \neq j} x_{d,i,j} \ , \ \bar{x}_n = \frac{1}{n} \sum_{i \neq j} x_{n,i} \\ \overline{yx_d} &= \frac{1}{n(n-1)} \sum_{i \neq j} y_{i,j} x_{d,i,j} \ , \ \overline{yx_n} = \frac{1}{n(n-1)} \sum_i x_{n,i} \overline{y}_i. \end{split}$$



$$\begin{split} g(\mathbf{Y}) &= \text{ an average of subgraphs of size } k, \text{ for } k \leq n_s \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{2}} \sum_{i < j} h(y_{i,j}, y_{j,i}) \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{3}} \sum_{i < j < k} h(y_{i,j}, y_{j,i}, y_{i,k}, y_{k,i}, y_{j,k}, y_{k,j}) \text{ if } n_s \geq 3 \end{split}$$

Why does it work?:

Each subgraph of size k appears in the sample with equal probability (although the subgraphs that appear are dependent).

- in and outdegree distributions;
- geodesics, distances, number of paths, etc.

$$\begin{split} g(\mathbf{Y}) &= \text{ an average of subgraphs of size } k, \text{ for } k \leq n_s \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{2}} \sum_{i < j} h(y_{i,j}, y_{j,i}) \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{3}} \sum_{i < j < k} h(y_{i,j}, y_{j,i}, y_{i,k}, y_{k,i}, y_{j,k}, y_{k,j}) \text{ if } n_s \geq 3 \end{split}$$

Why does it work?:

Each subgraph of size k appears in the sample with equal probability (although the subgraphs that appear are dependent).

- in and outdegree distributions;
- geodesics, distances, number of paths, etc.

$$\begin{split} g(\mathbf{Y}) &= \text{ an average of subgraphs of size } k, \text{ for } k \leq n_s \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{2}} \sum_{i < j} h(y_{i,j}, y_{j,i}) \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{3}} \sum_{i < j < k} h(y_{i,j}, y_{j,i}, y_{i,k}, y_{k,i}, y_{j,k}, y_{k,j}) \text{ if } n_s \geq 3 \end{split}$$

Why does it work?:

Each subgraph of size k appears in the sample with equal probability (although the subgraphs that appear are dependent).

- in and outdegree distributions;
- geodesics, distances, number of paths, etc.

$$\begin{split} g(\mathbf{Y}) &= \text{ an average of subgraphs of size } k, \text{ for } k \leq n_s \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{2}} \sum_{i < j} h(y_{i,j}, y_{j,i}) \\ g(\mathbf{Y}) &= \frac{1}{\binom{n}{3}} \sum_{i < j < k} h(y_{i,j}, y_{j,i}, y_{i,k}, y_{k,i}, y_{j,k}, y_{k,j}) \text{ if } n_s \geq 3 \end{split}$$

Why does it work?:

Each subgraph of size k appears in the sample with equal probability (although the subgraphs that appear are dependent).

- in and outdegree distributions;
- geodesics, distances, number of paths, etc.

Procedure:

1. Uniformly sample a set $\mathbf{e} = \{e_1, \dots, e_{n_e}\}$ of edges

 $\mathbf{e} \subset \{(i,j): y_{i,j} = 1\}$

2. Let \mathbf{Y}_s be the edge-generated subgraph of \mathbf{e} .

Procedure:

1. Uniformly sample a set $\mathbf{e} = \{e_1, \ldots, e_{n_e}\}$ of edges

$$\mathbf{e} \subset \{(i,j): y_{i,j} = 1\}$$

2. Let \mathbf{Y}_s be the edge-generated subgraph of \mathbf{e} .












Edge-induced subgraph sampling

How well do these subgraphs represent $\boldsymbol{Y}?$

Can you infer anything about **Y** from these data?

Edge-induced subgraph sampling

How well do these subgraphs represent \mathbf{Y} ?

Can you infer anything about **Y** from these data?

Edge-induced subgraph sampling



7 8 9

Procedure:

1. Uniformly sample a set $\textbf{s}_1 = \{\textbf{s}_{1,1}, \ldots, \textbf{s}_{1,\textit{n}_s}\}$ of nodes

 $\mathbf{s}_1 \subset \{1,\ldots,n\}.$

- 2. Observe the relations for each $i \in \mathbf{s}_1$, i.e. observe $\{y_{i,1}, \ldots, y_{i,n}\}$.
- 3. Let s_2 be the set of nodes having a link from anyone in $s_1.$ Observe the relations of anyone in s_2 to anyone in $s_1\cup s_2.$

$$\mathbf{Y}_{\mathbf{s}} = \{ y_{i,j} : i, j \in \mathbf{s}_1 \cup \mathbf{s}_2 \}$$

For large graphs, these data can be obtained (with high probability) by asking each $i \in \mathbf{s}_1$ the following:

- 1. Who are your friends?
- 2. Among your friends, which are friends with each other?

Procedure:

1. Uniformly sample a set $\textbf{s}_1 = \{\textbf{s}_{1,1}, \ldots, \textbf{s}_{1,\textit{n}_s}\}$ of nodes

$$\mathbf{s}_1 \subset \{1,\ldots,n\}.$$

- 2. Observe the relations for each $i \in \mathbf{s}_1$, i.e. observe $\{y_{i,1}, \ldots, y_{i,n}\}$.
- 3. Let s_2 be the set of nodes having a link from anyone in $s_1.$ Observe the relations of anyone in s_2 to anyone in $s_1\cup s_2.$

$$\mathbf{Y}_{\mathbf{s}} = \{y_{i,j} : i, j \in \mathbf{s}_1 \cup \mathbf{s}_2\}$$

For large graphs, these data can be obtained (with high probability) by asking each $i \in \mathbf{s}_1$ the following:

- 1. Who are your friends?
- 2. Among your friends, which are friends with each other?

Procedure:

1. Uniformly sample a set $\textbf{s}_1 = \{\textbf{s}_{1,1}, \ldots, \textbf{s}_{1,\textit{n}_s}\}$ of nodes

$$\mathbf{s}_1 \subset \{1,\ldots,n\}.$$

- 2. Observe the relations for each $i \in \mathbf{s}_1$, i.e. observe $\{y_{i,1}, \ldots, y_{i,n}\}$.
- 3. Let s_2 be the set of nodes having a link from anyone in $s_1.$ Observe the relations of anyone in s_2 to anyone in $s_1\cup s_2.$

$$\mathbf{Y}_{\mathbf{s}} = \{y_{i,j} : i, j \in \mathbf{s}_1 \cup \mathbf{s}_2\}$$

For large graphs, these data can be obtained (with high probability) by asking each $i \in s_1$ the following:

- 1. Who are your friends?
- 2. Among your friends, which are friends with each other?

Procedure:

1. Uniformly sample a set $s_1 = \{ \textit{s}_{1,1}, \ldots, \textit{s}_{1,\textit{n}_s} \}$ of nodes

$$\mathbf{s}_1 \subset \{1,\ldots,n\}.$$

- 2. Observe the relations for each $i \in \mathbf{s}_1$, i.e. observe $\{y_{i,1}, \ldots, y_{i,n}\}$.
- 3. Let s_2 be the set of nodes having a link from anyone in s_1 . Observe the relations of anyone in s_2 to anyone in $s_1 \cup s_2$.

$$\mathbf{Y}_{\mathbf{s}} = \{y_{i,j} : i, j \in \mathbf{s}_1 \cup \mathbf{s}_2\}$$

For large graphs, these data can be obtained (with high probability) by asking each $i \in s_1$ the following:

- 1. Who are your friends?
- 2. Among your friends, which are friends with each other?

Procedure:

1. Uniformly sample a set $s_1 = \{ \textit{s}_{1,1}, \ldots, \textit{s}_{1,\textit{n}_s} \}$ of nodes

$$\mathbf{s}_1 \subset \{1,\ldots,n\}.$$

- 2. Observe the relations for each $i \in \mathbf{s}_1$, i.e. observe $\{y_{i,1}, \ldots, y_{i,n}\}$.
- 3. Let s_2 be the set of nodes having a link from anyone in s_1 . Observe the relations of anyone in s_2 to anyone in $s_1 \cup s_2$.

$$\mathbf{Y}_{\mathbf{s}} = \{y_{i,j} : i, j \in \mathbf{s}_1 \cup \mathbf{s}_2\}$$

For large graphs, these data can be obtained (with high probability) by asking each $i \in \mathbf{s}_1$ the following:

- 1. Who are your friends?
- 2. Among your friends, which are friends with each other?

Link-tracing designs

Snowball sampling: Iteratively repeat the egocentric sampler, obtaining the stage-k nodes s_k from the links of s_{k-1} .

This is a type of link-tracing design. The links of the current nodes determine who is next to be included in the sample.

How will such subgraphs \boldsymbol{Y}_s be similar to $\boldsymbol{Y}?$ How will they differ?

Link-tracing designs

Snowball sampling: Iteratively repeat the egocentric sampler, obtaining the stage-k nodes s_k from the links of s_{k-1} .

This is a type of link-tracing design. The links of the current nodes determine who is next to be included in the sample.

How will such subgraphs \boldsymbol{Y}_s be similar to $\boldsymbol{Y}?$ How will they differ?

Link-tracing designs

Snowball sampling: Iteratively repeat the egocentric sampler, obtaining the stage-k nodes s_k from the links of s_{k-1} .

This is a type of link-tracing design. The links of the current nodes determine who is next to be included in the sample.

How will such subgraphs \boldsymbol{Y}_s be similar to $\boldsymbol{Y}?$ How will they differ?



8





















• Y_s is not generally representative of Y.

- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of **Y**_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- \bullet For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of \mathbf{Y}_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of Y_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of Y_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of Y_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - · covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of \mathbf{Y}_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of \mathbf{Y}_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

- Y_s is not generally representative of Y.
- For some statistics, weighted averages based on \mathbf{Y}_{s} can be unbiased (Horwitz-Thompson estimator).
- For many statistics, *part* of \mathbf{Y}_s can be used to obtain good estimates:
 - degree distributions can be estimated from degrees of egos;
 - covariate distributions can be estimated from those of the egos;

References

- Snijders (1992), "Estimation on the basis of snowball samples: How to Weight?"
- Kolaczyk (2009) "Sampling and Estimation in Network Graphs," chapter 5 of *Statistical Analysis of Network Data*.

Model: $Pr(\mathbf{Y} = \mathbf{y}|\theta), \theta \in \Theta$.

Complete data: Y

Observed data: Y[O], where O is a set of pairs of indices

$$\mathbf{0} = \begin{pmatrix} i_1 & j_1 \\ i_2 & j_2 \\ i_3 & j_3 \\ \vdots & \vdots \\ i_s & j_s \end{pmatrix}$$

Parameter estimation with incomplete sampled data

Model: $Pr(\mathbf{Y} = \mathbf{y}|\theta), \theta \in \Theta$.

Complete data: Y

Observed data: Y[O], where O is a set of pairs of indices

$$\mathbf{0} = \begin{pmatrix} i_1 & j_1 \\ i_2 & j_2 \\ i_3 & j_3 \\ \vdots & \vdots \\ i_s & j_s \end{pmatrix}$$

Parameter estimation with incomplete sampled data

Model: $Pr(\mathbf{Y} = \mathbf{y}|\theta), \theta \in \Theta$.

Complete data: Y

Observed data: Y[0], where **O** is a set of pairs of indices

$$\mathbf{0} = \begin{pmatrix} i_1 & j_1 \\ i_2 & j_2 \\ i_3 & j_3 \\ \vdots & \vdots \\ i_5 & j_5 \end{pmatrix}$$

Model: $Pr(\mathbf{Y} = \mathbf{y}|\theta), \theta \in \Theta$.

Complete data: Y

Observed data: Y[0], where **O** is a set of pairs of indices

$$\mathbf{0} = \begin{pmatrix} i_1 & j_1 \\ i_2 & j_2 \\ i_3 & j_3 \\ \vdots & \vdots \\ i_s & j_s \end{pmatrix}$$

Study design and missing data

Node-induced subgraph sampling

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA
Node-induced subgraph sampling

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

Node-induced subgraph sampling

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

Node-induced subgraph sampling: Observed data

	1	2	3	4	5	6
1	NA	NA	NA	NA	NA	NA
2	NA	NA	1	NA	0	NA
3	NA	0	NA	NA	0	NA
4	NA	NA	NA	NA	NA	NA
5	NA	0	1	NA	NA	NA
6	NA	NA	NA	NA	NA	NA

Edge-induced subgraph sampling

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

Edge-induced subgraph sampling

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

Edge-induced subgraph sampling: Observed data

	1	2	3	4	5	6
1	NA	NA	NA	NA	NA	1
2	NA	NA	1	NA	NA	NA
3	1	NA	NA	NA	NA	NA
4	NA	NA	NA	NA	NA	NA
5	NA	NA	1	NA	NA	NA
6	NA	NA	1	NA	NA	NA

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

	1	2	3	4	5	6
1	NA	0	1	0	0	1
2	0	NA	1	0	0	1
3	1	0	NA	1	0	0
4	0	1	0	NA	0	0
5	0	0	1	0	NA	1
6	1	0	1	0	0	NA

	1	2	3	4	5	6
1	NA	NA	NA	NA	NA	NA
2	0	NA	1	0	0	1
3	NA	0	NA	NA	NA	0
4	NA	NA	NA	NA	NA	NA
5	NA	NA	NA	NA	NA	NA
6	NA	0	1	NA	NA	NA

If the data are missing at random, i.e. the value of $\boldsymbol{o},$ what you get to observe,

- doesn't depend on $\boldsymbol{\theta}$
- doesn't depend on values of Y,

then valid likelihood and Bayesian inference can be obtained from the observed-data likelihood:

$$I_{MAR}(\theta : \mathbf{y}[\mathbf{o}]) = \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}] : \theta)$$
$$= \sum_{\mathbf{y}[\mathbf{o}^c]} \Pr(\mathbf{Y} = \mathbf{y} : \theta)$$

Inference based on $I(\theta : \mathbf{y}[\mathbf{o}])$ is provided in amen:

If the data are missing at random, i.e. the value of $\boldsymbol{o},$ what you get to observe,

- doesn't depend on $\boldsymbol{\theta}$
- doesn't depend on values of Y,

then valid likelihood and Bayesian inference can be obtained from the observed-data likelihood:

$$I_{MAR}(\theta : \mathbf{y}[\mathbf{o}]) = \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}] : \theta)$$
$$= \sum_{\mathbf{y}[\mathbf{o}^c]} \Pr(\mathbf{Y} = \mathbf{y} : \theta)$$

Inference based on $I(\theta : \mathbf{y}[\mathbf{o}])$ is provided in amen:

If the data are missing at random, i.e. the value of $\boldsymbol{o},$ what you get to observe,

- doesn't depend on $\boldsymbol{\theta}$
- doesn't depend on values of Y,

then valid likelihood and Bayesian inference can be obtained from the observed-data likelihood:

$$I_{MAR}(\theta : \mathbf{y}[\mathbf{o}]) = \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}] : \theta)$$
$$= \sum_{\mathbf{y}[\mathbf{o}^c]} \Pr(\mathbf{Y} = \mathbf{y} : \theta)$$

Inference based on $I(\theta : \mathbf{y}[\mathbf{o}])$ is provided in amen:

If the data are missing at random, i.e. the value of \mathbf{o} , what you get to observe,

- doesn't depend on $\boldsymbol{\theta}$
- doesn't depend on values of **Y**,

then valid likelihood and Bayesian inference can be obtained from the observed-data likelihood:

$$\begin{split} I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) &= \mathsf{Pr}(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]:\theta) \\ &= \sum_{\mathbf{y}[\mathbf{o}^c]} \mathsf{Pr}(\mathbf{Y} = \mathbf{y}:\theta) \end{split}$$

Inference based on $l(\theta : \mathbf{y}[\mathbf{o}])$ is provided in amen:

If the data are missing at random, i.e. the value of **o**, what you get to observe,

- doesn't depend on $\boldsymbol{\theta}$
- doesn't depend on values of **Y**,

then valid likelihood and Bayesian inference can be obtained from the observed-data likelihood:

$$\begin{split} I_{MAR}(\theta : \mathbf{y}[\mathbf{o}]) &= \mathsf{Pr}(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}] : \theta) \\ &= \sum_{\mathbf{y}[\mathbf{o}^c]} \mathsf{Pr}(\mathbf{Y} = \mathbf{y} : \theta) \end{split}$$

Inference based on $l(\theta : \mathbf{y}[\mathbf{o}])$ is provided in amen:

- Node-induced subgraph sampling?
- Edge-induced subgraph sampling?
- Egocentric sampling?

- Node-induced subgraph sampling?
- Edge-induced subgraph sampling?
- Egocentric sampling?

- Node-induced subgraph sampling?
- Edge-induced subgraph sampling?
- Egocentric sampling?

- Node-induced subgraph sampling?
- Edge-induced subgraph sampling?
- Egocentric sampling?

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- **O** = **o**, the determination of which relations you get to see;
- **Y**[**O**] = **y**[**o**], the relationship values for the observable relations.

The likelihood is then

$$\begin{split} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= l_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{split}$$

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- $\mathbf{O} = \mathbf{o}$, the determination of which relations you get to see;
- Y[O] = y[o], the relationship values for the observable relations.

The likelihood is then

$$\begin{split} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= l_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{split}$$

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- **O** = **o**, the determination of which relations you get to see;
- Y[O] = y[o], the relationship values for the observable relations.

The likelihood is then

$$\begin{split} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= l_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{split}$$

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- **O** = **o**, the determination of which relations you get to see;
- $\mathbf{Y}[\mathbf{O}] = \mathbf{y}[\mathbf{o}]$, the relationship values for the observable relations.

The likelihood is then

$$\begin{aligned} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{aligned}$$

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- **O** = **o**, the determination of which relations you get to see;
- Y[O] = y[o], the relationship values for the observable relations.

The likelihood is then

$$\begin{split} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{split}$$

While egocentric and other link-tracing designs are not MAR, they still can be analyzed as if they were. The argument is as follows:

The "data" include

- **O** = **o**, the determination of which relations you get to see;
- Y[O] = y[o], the relationship values for the observable relations.

The likelihood is then

$$\begin{split} l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}], \mathbf{O} = \mathbf{o}|\theta) \\ &= \Pr(\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]|\theta) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \\ &= I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta, \mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}]) \end{split}$$

$\textit{I}(\theta:\mathbf{o},\mathbf{y[o]}) = \textit{I}_{\textit{MAR}}(\theta:\mathbf{y[o]}) \times \Pr(\mathbf{O}=\mathbf{o}|\theta,\mathbf{Y[o]}=\mathbf{y[o]})$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that **O** equals **o**

- doesn't depend on θ
- only depends on Y through Y[o].

then the design is ignorable.

$l(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = l_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that **O** equals **o**

- doesn't depend on θ
- only depends on Y through Y[o].

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that **O** equals **o**

- doesn't depend on θ
- only depends on Y through Y[o].

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that ${\bf O}$ equals ${\bf o}$

- doesn't depend on $\boldsymbol{\theta}$
- only depends on Y through Y[o].

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that ${\bf O}$ equals ${\bf o}$

- doesn't depend on $\boldsymbol{\theta}$
- only depends on \boldsymbol{Y} through $\boldsymbol{Y}[\boldsymbol{o}].$

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that ${\bf O}$ equals ${\bf o}$

- doesn't depend on $\boldsymbol{\theta}$
- only depends on **Y** through **Y**[**o**].

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that ${\bf O}$ equals ${\bf o}$

- doesn't depend on $\boldsymbol{\theta}$
- only depends on **Y** through **Y**[**o**].

then the design is ignorable.

 $I(\theta:\mathbf{o},\mathbf{y}[\mathbf{o}]) = I_{MAR}(\theta:\mathbf{y}[\mathbf{o}]) \times \Pr(\mathbf{O} = \mathbf{o}|\theta,\mathbf{Y}[\mathbf{o}] = \mathbf{y}[\mathbf{o}])$

When is the design ignorable?

(MAR) If the probability that **O** equals **o** doesn't depend on θ or **Y** (e.g., node-induced subgraph sampling), the design is ignorable.

ID If the probability that ${\bf O}$ equals ${\bf o}$

- doesn't depend on $\boldsymbol{\theta}$
- only depends on **Y** through **Y**[**o**].

then the design is ignorable.

References

- Thompson and Frank (2000) "Model-based estimation with link-tracing sampling designs"
- Heitjan and Basu (1996) "Distinguishing 'Missing at Random' and 'Missing Completely at Random'

Simulation study - ID likelihoods

$y_{i,j} = \beta_0 + \beta_r x_{n,i} + \beta_c x_{n,j} + \beta_{d,i,j} + a_i + b_j + \epsilon_{i,j}$

fit.pop = fitted model based on complete network data

fit.**samp** = fitted model based on sampled network data

How do the parameter estimates of fit.samp compare to those of fit.pop?
Simulation study - ID likelihoods

$$y_{i,j} = \beta_0 + \beta_r x_{n,i} + \beta_c x_{n,j} + \beta_{d,i,j} + a_i + b_j + \epsilon_{i,j}$$

fit.pop = fitted model based on complete network data fit.samp = fitted model based on sampled network data

How do the parameter estimates of fit.samp compare to those of fit.pop?

Simulation study - ID likelihoods

$$y_{i,j} = \beta_0 + \beta_r x_{n,i} + \beta_c x_{n,j} + \beta_{d,i,j} + a_i + b_j + \epsilon_{i,j}$$

fit.pop = fitted model based on complete network data
fit.samp = fitted model based on sampled network data
How do the parameter estimates of fit.samp compare to those of fit.pop?

Node-induced subgraph sample

 $n_p = 32, n_s = 10$ -3.0 -2.0 -1.0 -4.0 β_0 -0.2 -0.1 0.0 0.1 0.2 0.3 βc



Egocentric sample

