



A Bayesian Approach to Linking Historical Records of America's Enslaved

Hannah Butler & Andee Kaplan

Colorado State University

Overview

Goal: Develop an effective probabilistic record linkage model that leverages the presence of alias information in historical slave records.

- Define Record Linkage
- Record Linkage for Historical Slave Records
- Model Specification & Method
- Simulation Performance
- Application & Results

What is Record Linkage?

A method used to identify **co-referent** entities from overlapping datasets.

Author	Title	Publ. Date
Stephen Jay Gould	The Mismeasure of Man	1996
Ian Hacking	Taming of Chance	1990
Stephen T. Ziliak, Deirdre N. McCloskey	The Cult of Statistical Significance	2008
Theodore Porter	Trust in Numbers	2020

Author	Title	Publ. Date
Ruha Benjamin	Race After Technology	2019
Ian Hackng	The Taming of Chance	1990
Theodore M. Porter	The Rise of Statistical Thinking	2020
Michael E. Staub	The Mismeasure of Minds	2018

- **Deterministic Record Linkage**

- Define rules to match field values
- Eg: **JOIN** in SQL

- **Probabilistic Record Linkage**

- Statistical model used to estimate matches between records
- More flexibility for error-prone data
- **Eg: Bayesian Bipartite Record Linkage (BRL)** (*Sadinle 2017*)

Comparison Data

- Calculate distances between values of a field with a chosen similarity metric
- Classify distances as **levels of disagreement** (*Fellegi & Sunter, 1969*)
- Level of disagreement | field ~ Multinomial Mixture

Assumption

- **There are no duplicate records within a single file** (*Sadinle, 2017*)

Challenges of Linking Slave Records

• Data Availability

- Comparatively few digitized data sources
- High rates of missing data
- Eg: Last names

• Uncertainty

- Data collection is unregulated/non-standardized
- Falsified/contentious/subjective values
- Transcription errors
- Eg: Ages

• Duplicate Records (Non-Erroneous)

- Many entities have multiple variable values for a field (**Aliases**)
- Potential violations of bipartite linkage structure

From: **Louisiana Kindred (Kinfolkology)**

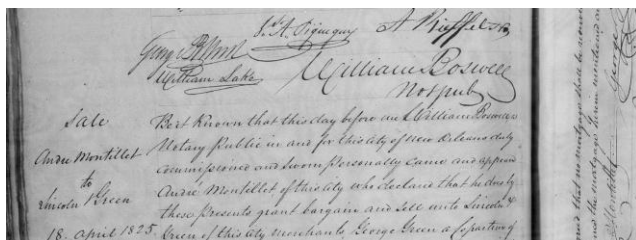
ID (LK)	Enslaved Name	Racial Descriptor	Age	Kin Relation	Kin Name
PELK100940	Cook, James	"mulatto"	25	<NA>	<NA>
PELK100498	Lucy	"negro"	37	Mother of	Harrison
PELK100498	Lucy	"negro"	37	Mother of	Arthennise
PELK100498	Lucy	"negro"	37	Mother of	[unnamed child]

1 entity, 1 alias

1 entity, 3 aliases

From: **Oceans of Kinfolk (Kinfolkology)**

ID (OOK)	Enslaved Name	Racial Descriptor	Age	Kin Relation	Kin Name
OOK104652	Cook, James	"yellow"	20	<NA>	<NA>
OOK106770	Lewis, Lucy	"black"	35	<NA>	<NA>
OOK109450	Lucy	"dark"	40	<NA>	<NA>
OOK115546	Lucy	"black"	38	Mother of	<NA>



MARY EPPS

45

I cannot say much about the place as I have ben here but a short time but so far as I have seen I like very well. you will give my Respect to your lady, & Mr & Mrs Brown. If you have not written to Petersburg you will please to write as soon as can I have nothing More to Write at present but yours Respectfully

EMMA BROWN (old name **MARY EPPS**)

- The Underground Railroad

Existing Approaches

- **Record Linkage with Naïve Pre-Processing**

- Remove alias values/records during data processing
- Perform record linkage as usual

ID (LK)	Enslaved Name	Racial Descriptor	Age	Kin Relation	Kin Name
PELK100940	Cook, James	"mulatto"	25	<NA>	<NA>
PELK100498	Lucy	"negro"	37	Mother of	Harrison
PELK100498	Lucy	"negro"	37	Mother of	Arthennise
PELK100498	Lucy	"negro"	37	Mother of	[unnamed child]

- **Record Linkage with Post-Processing**

- Keep all records and values
- Perform record linkage as usual
- Resolve conflicts after estimation

ID (LK)		ID (OOK)	
PELK100940	Cook, James	OOK104652	Cook, James
PELK100498	Lucy	OOK106770	Lewis, Lucy
PELK100498	Lucy	OOK109450	Lucy
PELK100498	Lucy	OOK115546	Lucy

Bipartite Alias Record Linkage (BARL)

1. Specify a Bayesian BRL model that accommodates alias records (BARL) without introducing conflict

Comparisons

$$\gamma_{ij}^f \mid \mathbf{m}_f, \mathbf{u}_f, \mathbf{Z} \stackrel{\text{iid}}{\sim} \text{Mult}(1, K_f, \mathbf{m}_f^{\mathbf{Z}_i=j} \mathbf{u}_f^{\mathbf{Z}_i \neq j})$$

$$\Gamma \mid \mathbf{m}, \mathbf{u}, \mathbf{Z} \sim \prod_{i=1}^{N_1} \prod_{j=1}^{N_2} \prod_{f=1}^F \prod_{k=1}^{K_f} \left[(\mathbf{m}_f(k))^{\mathbf{1}(\mathbf{Z}_i=j)} (\mathbf{u}_f(k))^{\mathbf{1}(\mathbf{Z}_i \neq j)} \right]^{\gamma_{ij}^f(k)} \mathbf{V}_{ij}$$

Level of Disagreement Probabilities

$$\mathbf{m}_f \mid \boldsymbol{\alpha}_f \sim \text{Dirichlet}(\boldsymbol{\alpha}_f) \quad \mathbf{u}_f \mid \boldsymbol{\beta}_f \sim \text{Dirichlet}(\boldsymbol{\beta}_f)$$

2. Choose a linkage prior for \mathbf{Z} that equally weights links between any two records

$$P(\mathbf{Z} \mid a, b) \propto \frac{B(a + L(\mathbf{Z}), b + N_1 - L(\mathbf{Z})) (N_2 - L(\mathbf{Z}))!}{B(a, b) N_2!} \mathbf{1}(\mathbf{Z}_i \neq \mathbf{Z}_{i'} \forall i \neq i')$$

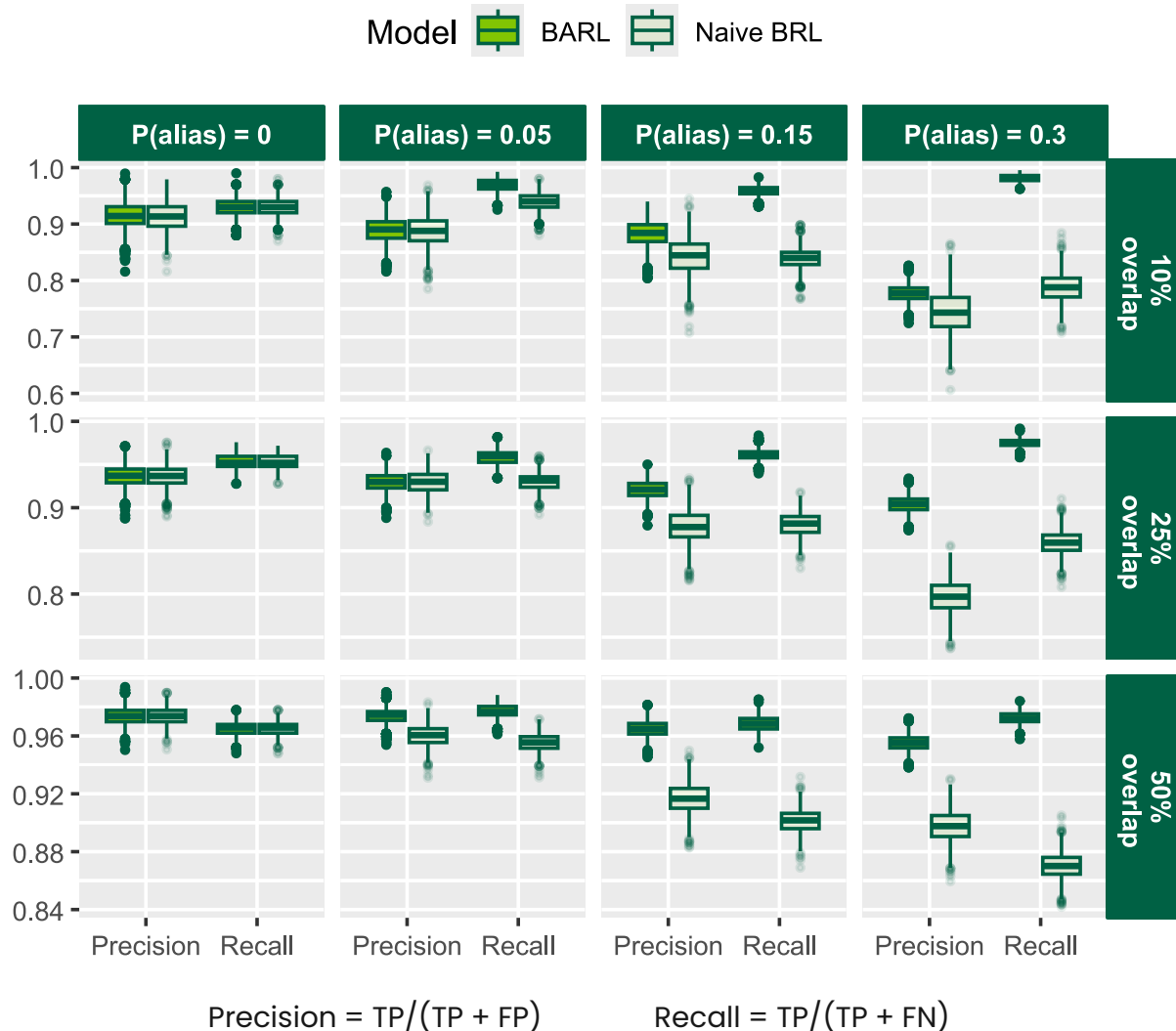
3. Use MCMC to generate posterior samples

1. Initialize agreement probabilities and links
2. For each iteration, $t = 1, \dots, T$:
 1. Sample \mathbf{Z} given data, \mathbf{m} , \mathbf{u}
 2. Sample \mathbf{m} , \mathbf{u} given data, \mathbf{Z}

4. Estimate links directly from MCMC posterior samples or consolidate records to obtain point estimate

$$\hat{\mathbf{Z}}_i = \begin{cases} j & \text{if } P(\mathbf{Z}_i = j \mid \Gamma) > 1/2 \\ N_2 + i & \text{otherwise} \end{cases}$$

Model Performance Simulation



Data Simulation

- Records generated using GeCo
- Simulated overlapping datasets with aliases
- BARL compared to BRL with additional aliases removed (Naïve BRL)
- When no aliases are present, BARL = BRL

Results

- BARL: Greater average precision & recall than Naïve BRL

Oceans of Kinfolk

Louisiana Kindred

www.kinfolkology.org/data-overview



MANIFEST OF NEGROES, MILITATES, AND PERSONS OF COLOR, taken on board the
Bay State of *Massachusetts* on the 14th of *April* 1862
 transported to the Port of *New Orleans* in Master, *hulk* *to* *be* *sent* *to* *the* *United* *States* *Army*
 or disposed of at El Paso, or to be held in service or sold.
 for the purpose of being sold

Number of Entry	NAMES	SEX.	AGE.	HEIGHT. Feet. Inches	Whether Negro, Mulatto, or Person of Color.	Inscribed or Disposed of in Name and Particular Destination.	
1	Wheeler	Male	23	5	136-58	German	Wheeler, 1862
2	Wells	Male	19	3	140	"	"
3	Smith	Male	45	6	140	"	"
4	Martha	Male	9	6	3	140	"
5	Lincol	Male	9	3	4	140	"
6	Smith	Female	1	6	11	140	"
7	Parson	Male	22	3	140	"	"
8	Johnson	"	31	5	1	140	"
9	Steele	"	44	5	1	140	"
10	Stephens	Female	16	5	1	140	"
11	Evans	"	31	5	1	140	"
12	Smith	"	5	3	0 1/2	140	"
13	May	"	6	3	2	140	"
14	Wing	Male	3	4	1	140	"
15	Kelly	"	23	5	2 1/2	140	"
16	Wells	"	26	4	1	140	"
17	Wheeler	Male	11	0	6	Master William, 1862	1862
18	Smith	Female	11	0	4	"	"
19	Evans	"	11	0	4	Smith	"
20	Smith	"	11	0	1	Smith	"
21	Marshall	"	11	0	10	Smith	"
22	Lincol	Male	13	0	9	White	"
23	Stephens	"	13	0	9	"	"
24	Smith	"	13	0	6	"	"
25	Edwards	"	16	0	1	"	"

[illegible]

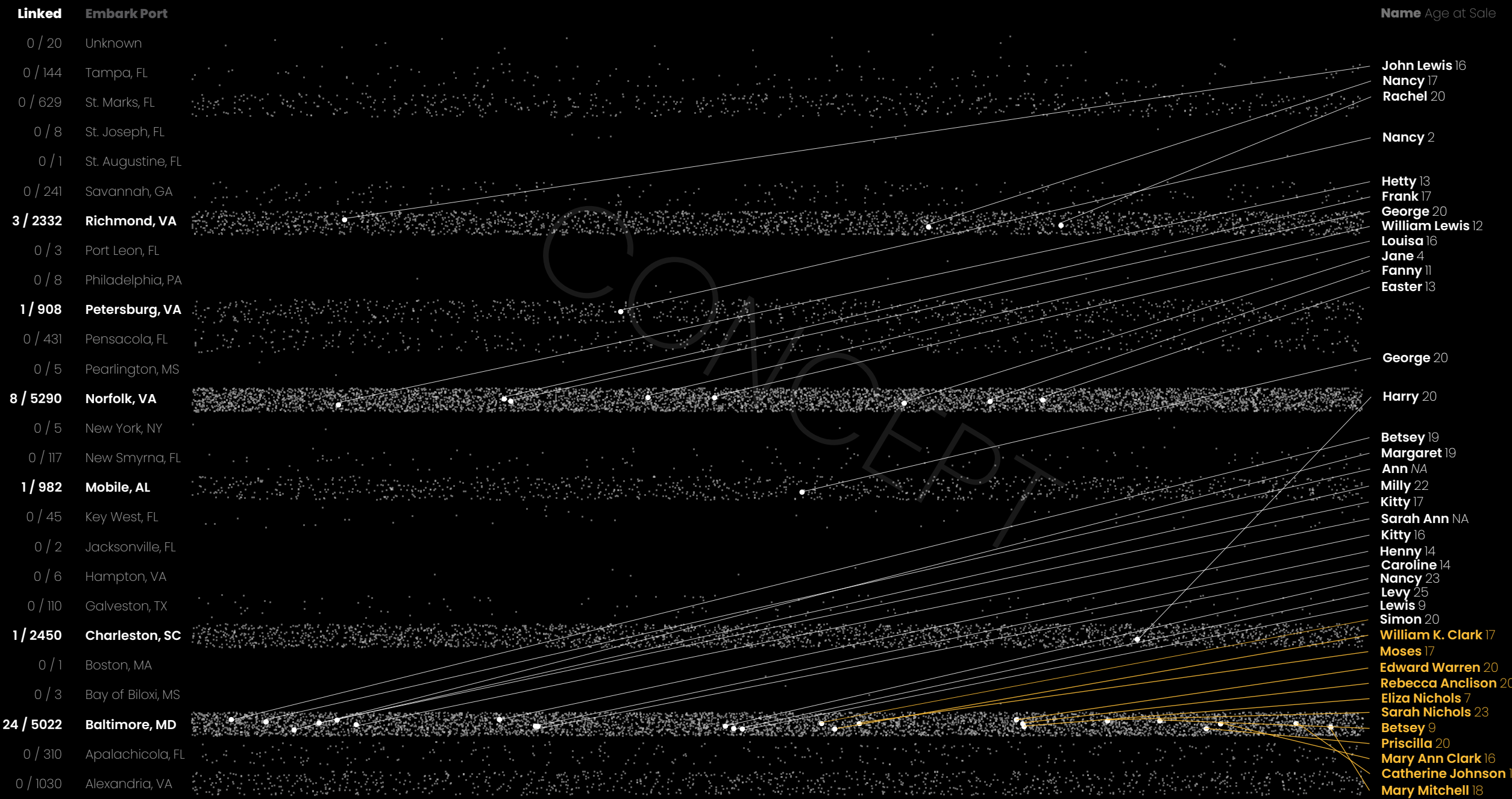
From: *Louisiana Kindred of Kinfolk on Kinfolkology*

Aligned Variables

- Enslaved Name (first, last)
- Enslaved Age
- Enslaved Gender
- Enslaved Kin Name
- Enslaved Kin Relation
- Event Date
- Enslaver Name (first, last)
- Enslaver Location (city, county, state)

Oceans of Kinfolk

Louisiana Kindred



Thank You!

Acknowledgments

- National Science Foundation¹
- Jennie K. Williams, PhD & Kinfolkology

¹ This material is based on work supported in part by the National Science Foundation under Grant No. SES-2338428 and DMS-2330089. The ideas in this work are representative of the authors and not of the NSF

Sources

- Fellegi, Ivan P., and Alan B. Sunter. "A Theory for Record Linkage." *Journal of the American Statistical Association* 64, no. 328 (1969): 1183–1210. <https://doi.org/10.2307/2286061>.
- Kinfolkology. "Kinfolkology." Accessed October 30, 2024. <https://www.kinfolkology.org>.
- National Center for Health Statistics. Division of Analysis and Epidemiology., "The Linkage of National Center for Health Statistics Survey Data to Centers for Medicare & Medicaid Services Transformed Medicaid Statistical Information System Claims Data (2014–2019): Matching Methodology and Analytic Considerations."
- Sadinle, Mauricio. "Bayesian Estimation of Bipartite Matchings for Record Linkage." *Journal of the American Statistical Association* 112, no. 518 (April 3, 2017): 600–612. <https://doi.org/10.1080/01621459.2016.1148612>.
- Still, William. *The Underground Railroad: A Record of Facts, Authentic Narratives, Letters &c., Narrating the Hardships, Hair-Breadth Escapes and Death Struggles of the Slaves in Their Efforts for Freedom*. Rev. ed. Philadelphia, Pa., Cincinnati, Ohio [etc.]: People's publishing company, 1879. <https://www.loc.gov/resource/rbc0001.2019gen24984/>.