# Record Linkage: Theory and Practice

William E. Yancey

`william.e.yancey@census.gov`

U.S. Census Bureau

# Introduction

- Definition

# Introduction

- Definition
- Terminology

# Introduction

- Definition
- Terminology
- Uses

# Introduction

- Definition

- Terminology

- Uses

- Context

# Definition

- A procedure to find pairs of records in two files that represent the same entity

US CENSUS BUREAU

# Definition

- A procedure to find pairs of records in two files that represent the same entity

- When both files are the same file, the procedure is to find duplicate records

# Terminology

- Matched records:  both records represent the same entity in truth

U S C E N S U S B U R E A U

# Terminology

- Matched records: both records represent the same entity in truth

- Linked records: Both records are identified by record linkage procedure as probably representing the same entity

USCENSUSBUREAU

# Uses

- Updating and deduplicating a survey frame

# Uses

- Updating and deduplicating a survey frame
- Merging two files for microdata anlysis

USCENSUSBUREAU

# Uses

- Updating and deduplicating a survey frame
- Merging two files for microdata anlysis
- Determine confidentiality of microdata

US CENSUS BUREAU

# Uses

- Updating and deduplicating a survey frame

- Merging two files for microdata anlysis

- Determine confidentiality of microdata

- Measure a population by capture-recapture

USCENSUSBUREAU

# Capture-Recapture

- Let $A, B$ be independent random samples of sample space $S$

$$x_{11} = |A \cap B| \quad x_{10} = |A \setminus B|$$
$$x_{01} = |B \setminus A| \quad x_{00} = |S \setminus (A \cup B)|$$

USCENSUSBUREAU

# Capture-Recapture

- Let $A, B$ be independent random samples of sample space $S$

$$x_{11} = |A \cap B| \quad x_{10} = |A \setminus B|$$
$$x_{01} = |B \setminus A| \quad x_{00} = |S \setminus (A \cup B)|$$

- Then

$$\hat{x}_{00} = E\left[x_{00}\right] = \frac{x_{1+}x_{+1}}{x_{11}}$$

# Capture-Recapture, Cont.

- Take two independent surveys of a region and estimate the number of people missed.

USCENSUSBUREAU

# Capture-Recapture, Cont.

- Take two independent surveys of a region and estimate the number of people missed.

- Note accuracy of estimate depends on accuracy of $x_{11}$, as determined by record linkage

## USCENSUSBUREAU

# Capture-Recapture, Cont.

- Take two independent surveys of a region and estimate the number of people missed.

- Note accuracy of estimate depends on accuracy of $x_{11}$, as determined by record linkage

- Y.M.M. Bishop, S.E. Fienberg, P.W. Holland, *Discrete Multivariate Analysis, Theory and Priactice*, Chapter 6. MIT Press, 1975

USCENSUSBUREAU

# Record Linkage Basics

- Context

# Record Linkage Basics

- Context

- Deterministic Record Linkage

# Record Linkage Basics

- Context

- Deterministic Record Linkage

- Probabilistic Record Linkage

U S C E N S U S B U R E A U

# Record Linkage Basics

- Context

- Deterministic Record Linkage

- Probabilistic Record Linkage

- Not Statistical Matching

USCENSUSBUREAU

# Record Linkage Basics

- Context

- Deterministic Record Linkage

- Probabilistic Record Linkage

- Not Statistical Matching

- Need for Automated Record Linkage

# Context

- Files have records of fixed length with fields of fixed length and position (or in a database with retrievable individual fields)

# Context

- Files have records of fixed length with fields of fixed length and position (or in a database with retrievable individual fields)

- Not a search algorithm

# Deterministic Record Linkage

- Records are linked when

U S C E N S U S B U R E A U

# Deterministic Record Linkage

- Records are linked when
  - They agree exactly on all matching fields

USCENSUSBUREAU

# Deterministic Record Linkage

- Records are linked when
  - They agree exactly on all matching fields
  - Or on predetermined portion of fields

# Probabilistic Record Linkage

- Assign a probabilistic weighting to record pairs

# Probabilistic Record Linkage

- Assign a probabilistic weighting to record pairs

- Accepts record pairs with sufficiently high weights as linked pairs

USCENSUSBUREAU

# Not Statistical Matching

- Statistical matching: Bring together pairs of records with statistically similar characteristics, not necessarily representing the same entity

# Not Statistical Matching

- Statistical matching: Bring together pairs of records with statistically similar characteristics, not necessarily representing the same entity

- Usually for two files that represent different samples of a population

# Not Statistical Matching

- Statistical matching: Bring together pairs of records with statistically similar characteristics, not necessarily representing the same entity

- Usually for two files that represent different samples of a population

- Older practice than exact matching (deterministic or probabilistic)

# Need for Automated Record Linkage

Clerical matching is:

# Need for Automated Record Linkage

Clerical matching is:

- expensive

# Need for Automated Record Linkage

Clerical matching is:

- expensive

- slow

# Need for Automated Record Linkage

Clerical matching is:

- expensive

- slow

- error prone

# Need for Automated Record Linkage

Clerical matching is:

- expensive

- slow

- error prone

|                           | Clerical | 1988 | 1990 |
|---------------------------|----------|------|------|
| Computer match proportion | 0%       | 70%  | 75%  |
| # clerks                  | 3000     | 600  | 200  |
| #months                   | 6        | 1.5  | 1.5  |
| False match rate          | 5%       | 0.5% | 0.2% |

U S C E N S U S B U R E A U

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

- Agreement Patterns

USCENSUSBUREAU

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

- Agreement Patterns

- Example Comparison Space

USCENSUSBUREAU

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

- Agreement Patterns

- Example Comparison Space

- Conditional Probabilities

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

- Agreement Patterns

- Example Comparison Space

- Conditional Probabilities

- Linkage Rule

USCENSUSBUREAU

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation

- Agreement Patterns

- Example Comparison Space

- Conditional Probabilities

- Linkage Rule

- Error Rates

USCENSUSBUREAU

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation
- Agreement Patterns
- Example Comparison Space
- Conditional Probabilities
- Linkage Rule
- Error Rates
- Clerical Region

USCENSUSBUREAU

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation
- Agreement Patterns
- Example Comparison Space
- Conditional Probabilities
- Linkage Rule
- Error Rates
- Clerical Region
- Fundamental Theorem

U S C E N S U S B U R E A U

# Rec. Link. Theory: Fellegi & Sunter

- Basic Definitions and Notation
- Agreement Patterns
- Example Comparison Space
- Conditional Probabilities
- Linkage Rule
- Error Rates
- Clerical Region
- Fundamental Theorem
- Conditional Independence Assumption

# Basic Definitions and Notation

- Sets of entities $A, B$

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha\left(A\right), \beta\left(B\right)$

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha(A), \beta(B)$

- Sample space $\alpha(A) \times \beta(B)$

USCENSUSBUREAU

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha\left(A\right), \beta\left(B\right)$

- Sample space $\alpha\left(A\right) \times \beta\left(B\right)$

- Matches

$$M = \left\{ \left(\alpha\left(a\right), \beta\left(b\right)\right) \mid a = b \right\}$$

USCENSUSBUREAU

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha\left(A\right), \beta\left(B\right)$

- Sample space $\alpha\left(A\right) \times \beta\left(B\right)$

- Matches

$$M = \left\{\left(\alpha\left(a\right), \beta\left(b\right)\right) \mid a = b\right\}$$

- Nonmatches

$$U = \left\{\left(\alpha\left(a\right), \beta\left(b\right)\right) \mid a \neq b\right\}$$

# Basic Definitions and Notation

- Sets of entities $A, B$

U S C E N S U S B U R E A U

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha(A), \beta(B)$

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha\left(A\right), \beta\left(B\right)$

- Sample space $\alpha\left(A\right) \times \beta\left(B\right)$

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha(A), \beta(B)$

- Sample space $\alpha(A) \times \beta(B)$

- Matches

$$M = \{(\alpha(a), \beta(b)) \mid a = b\}$$

USCENSUSBUREAU

# Basic Definitions and Notation

- Sets of entities $A, B$

- Corresponding files of records $\alpha(A), \beta(B)$

- Sample space $\alpha(A) \times \beta(B)$

- Matches

$$M = \{(\alpha(a), \beta(b)) \mid a = b\}$$

- Nonmatches

$$U = \{(\alpha(a), \beta(b)) \mid a \neq b\}$$

USCENSUSBUREAU

# Agreement Patterns

- Comparison space

$$\alpha\left(A\right) \times \beta\left(B\right) \rightarrow \Gamma$$

USCENSUSBUREAU

# Agreement Patterns

- Comparison space

$$\alpha\left(A\right) \times \beta\left(B\right) \rightarrow \Gamma$$

- Comparison vector

$$\gamma \in \Gamma$$

# Agreement Patterns

- Comparison space

$$\alpha\left(A\right) \times \beta\left(B\right) \rightarrow \Gamma$$

- Comparison vector

$$\gamma \in \Gamma$$

- Each component of comparison vector can take on finitely many values, typically two

# Agreement Patterns

- Comparison space

$$\alpha\left(A\right) \times \beta\left(B\right) \to \Gamma$$

- Comparison vector

$$\gamma \in \Gamma$$

- Each component of comparison vector can take on finitely many values, typically two

$$\gamma = \left(\gamma_1, \gamma_2, \dots, \gamma_n\right)$$
$$\gamma_i \in \{0, 1\}$$

US CENSUS BUREAU

# Example Comparison Space

- Consider 3 binary comparisons

# Example Comparison Space

- Consider 3 binary comparisons
  - $\gamma_1$      pair agrees on last name

# Example Comparison Space

- Consider 3 binary comparisons
  - $\gamma_1$      pair agrees on last name
  - $\gamma_2$      pair agrees on first name

USCENSUSBUREAU

# Example Comparison Space

- Consider 3 binary comparisons
  - $\gamma_1$       pair agrees on last name
  - $\gamma_2$       pair agrees on first name
  - $\gamma_3$       pair agrees on street name

U S C E N S U S B U R E A U

# Example Comparison Space

- Consider 3 binary comparisons
  - $\gamma_1$      pair agrees on last name
  - $\gamma_2$      pair agrees on first name
  - $\gamma_3$      pair agrees on street name
- Sample agreement pattern

$$\gamma = (1, 0, 1)$$

# Conditional Probabilities

- Probability that a record pair has agreement pattern $\gamma$, given that it is a match/nonmatch

$$\Pr\left(\gamma|M\right)$$
$$\Pr\left(\gamma|U\right)$$

USCENSUSBUREAU

# Conditional Probabilities

- Probability that a record pair has agreement pattern $\gamma$, given that it is a match/nonmatch

$$\Pr\left(\gamma|M\right)$$
$$\Pr\left(\gamma|U\right)$$

- Agreement ratio

$$R\left(\gamma\right) = \frac{\Pr\left(\gamma|M\right)}{\Pr\left(\gamma|U\right)}$$

USCENSUSBUREAU

# Conditional Probabilities

- Probability that a record pair has agreement pattern $\gamma$, given that it is a match/nonmatch

$$\Pr(\gamma|M)$$
$$\Pr(\gamma|U)$$

- Agreement ratio

$$R(\gamma) = \frac{\Pr(\gamma|M)}{\Pr(\gamma|U)}$$

- Conditioned on the unobservable truth

U S C E N S U S B U R E A U

# Linkage Rule

Designate a record pair's status based on its agreement pattern

# Linkage Rule

Designate a record pair's status based on its agreement pattern

- 🔴 Link

US CENSUS BUREAU

# Linkage Rule

Designate a record pair's status based on its agreement pattern

- Link
- Non-link

# Linkage Rule

Designate a record pair's status based on its agreement pattern

- Link

- Non-link

- Undecided

# Linkage Rule

Designate a record pair's status based on its agreement pattern

- Link

- Non-link

- Undecided

$$L : \Gamma \rightarrow \{L, N, C\}$$

USCENSUSBUREAU

# Error Rates

- 🔴 False match–a linked pair that is not a match

# Error Rates

- False match–a linked pair that is not a match

- False nonmatch–a nonlinked pair that is a match

# Error Rates

- 🔴 False match–a linked pair that is not a match

- 🔴 False nonmatch–a nonlinked pair that is a match

- 🔴 False match rate–probability that a designated link is a nonmatch

U S C E N S U S B U R E A U

# Error Rates

- False match–a linked pair that is not a match

- False nonmatch–a nonlinked pair that is a match

- False match rate–probability that a designated link is a nonmatch

$$\mu = \Pr\left(L|U\right)$$

U S C E N S U S B U R E A U

# Error Rates

- False match–a linked pair that is not a match

- False nonmatch–a nonlinked pair that is a match

- False match rate–probability that a designated link is a nonmatch

$$\mu = \Pr\left(L|U\right)$$

- False nonmatch rate–probability that a designated nonlink is a match

# Error Rates

- False match–a linked pair that is not a match

- False nonmatch–a nonlinked pair that is a match

- False match rate–probability that a designated link is a nonmatch

$$\mu = \Pr\left(L|U\right)$$

- False nonmatch rate–probability that a designated nonlink is a match

$$\lambda = \Pr\left(N|M\right)$$

USCENSUSBUREAU

# Error Rates, Cont.

If all pairs of records are designated link or nonlink

# Error Rates, Cont.

If all pairs of records are designated link or nonlink

|  | Match | Nonmatch |
|---|---|---|
| Link | $1 - \lambda$ | $\mu = \Pr\left(L|U\right)$ |
| Nonlink | $\lambda = \Pr\left(N|M\right)$ | $1 - \mu$ |

U S C E N S U S B U R E A U

# Clerical Region

- The set $C$ of record pairs which are designated neither probable link nor probable nonlink by the linkage rule

# Clerical Region

- The set $C$ of record pairs which are designated neither probable link nor probable nonlink by the linkage rule

- The match status of these pairs is left to clerical review

# Fundamental Theorem

- Fellegi & Sunter ("*A Theory for Record Linkage*", JASA, December,1969)

# Fundamental Theorem

- Fellegi & Sunter ("*A Theory for Record Linkage*", JASA, December,1969)

- Order the comparison vectors $\left\{\gamma^j\right\}$ by their agreement ratios $R\left(\gamma^j\right)$

# Fundamental Theorem

- Fellegi & Sunter ("*A Theory for Record Linkage*", JASA, December,1969)

- Order the comparison vectors $\left\{ \gamma^j \right\}$ by their agreement ratios $R\left(\gamma^j\right)$

- Choose upper $T_\mu$ and lower $T_\lambda$ cutoff values for $R\left(\gamma\right)$

USCENSUSBUREAU

# Fundamental Theorem, Cont.

Linkage rule:

# Fundamental Theorem, Cont.

Linkage rule:

- Pairs with $R\left(\gamma^j\right) \geq T_\mu$ are designated links

# Fundamental Theorem, Cont.

Linkage rule:

- Pairs with $R\left(\gamma^j\right) \geq T_\mu$ are designated links

- Pairs with $R\left(\gamma^j\right) \leq T_\lambda$ are designated nonlinks

# Fundamental Theorem, Cont.

Linkage rule:

- Pairs with $R\left(\gamma^j\right) \geq T_\mu$ are designated links

- Pairs with $R\left(\gamma^j\right) \leq T_\lambda$ are designated nonlinks

- Pairs with $T_\lambda < R\left(\gamma^j\right) < T_\mu$ are in the clerical region

## USCENSUSBUREAU

# Fundamental Theorem, Cont.

The error rates for this linkage rule are

# Fundamental Theorem, Cont.

The error rates for this linkage rule are

$$\mu = \sum_{R(\gamma^j) \geq T_\mu} \mathrm{Pr}\left(\gamma^j | U\right)$$

# Fundamental Theorem, Cont.

The error rates for this linkage rule are

$$\mu = \sum_{R(\gamma^j) \geq T_\mu} \Pr\left(\gamma^j | U\right)$$
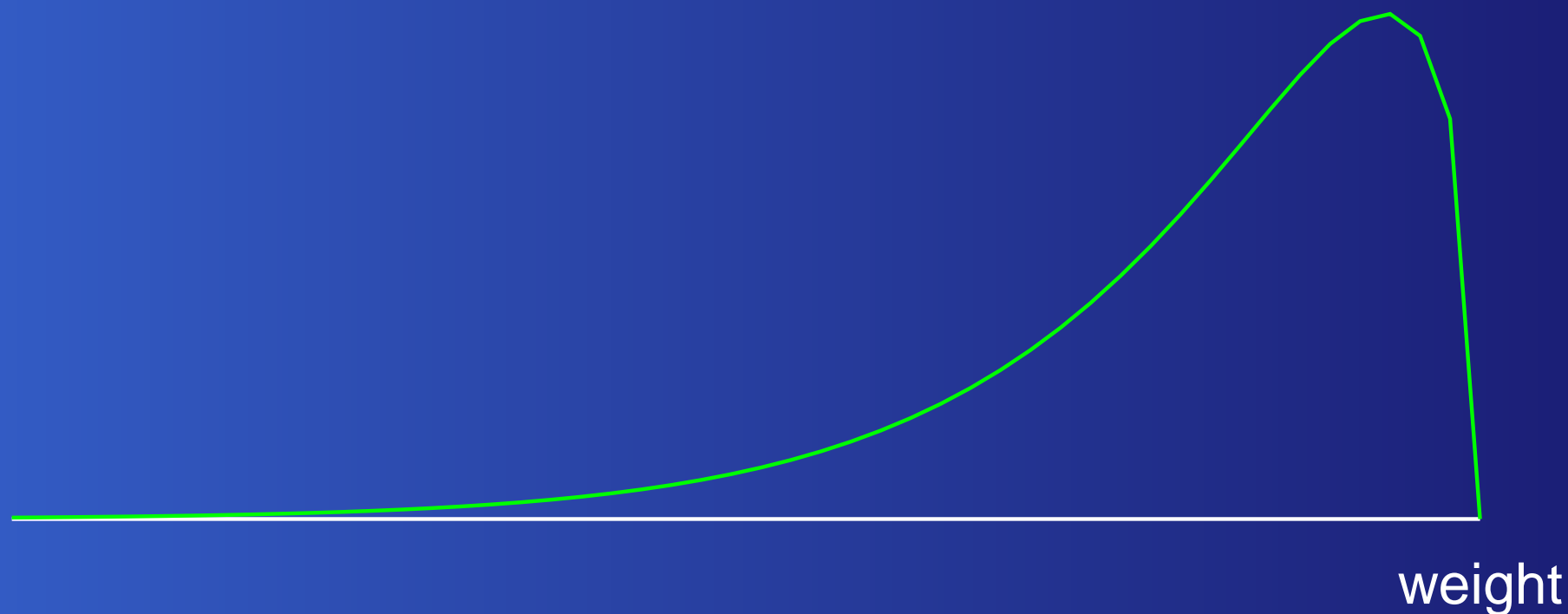
$$\lambda = \sum_{R(\gamma^j) \leq T_\lambda} \Pr\left(\gamma^j | M\right)$$

USCENSUSBUREAU

# Fundamental Theorem, Cont.

- Theorem: For these error rates $\mu, \lambda$, this is the optimal linkage rule, in the sense of producing the minimum size critical region

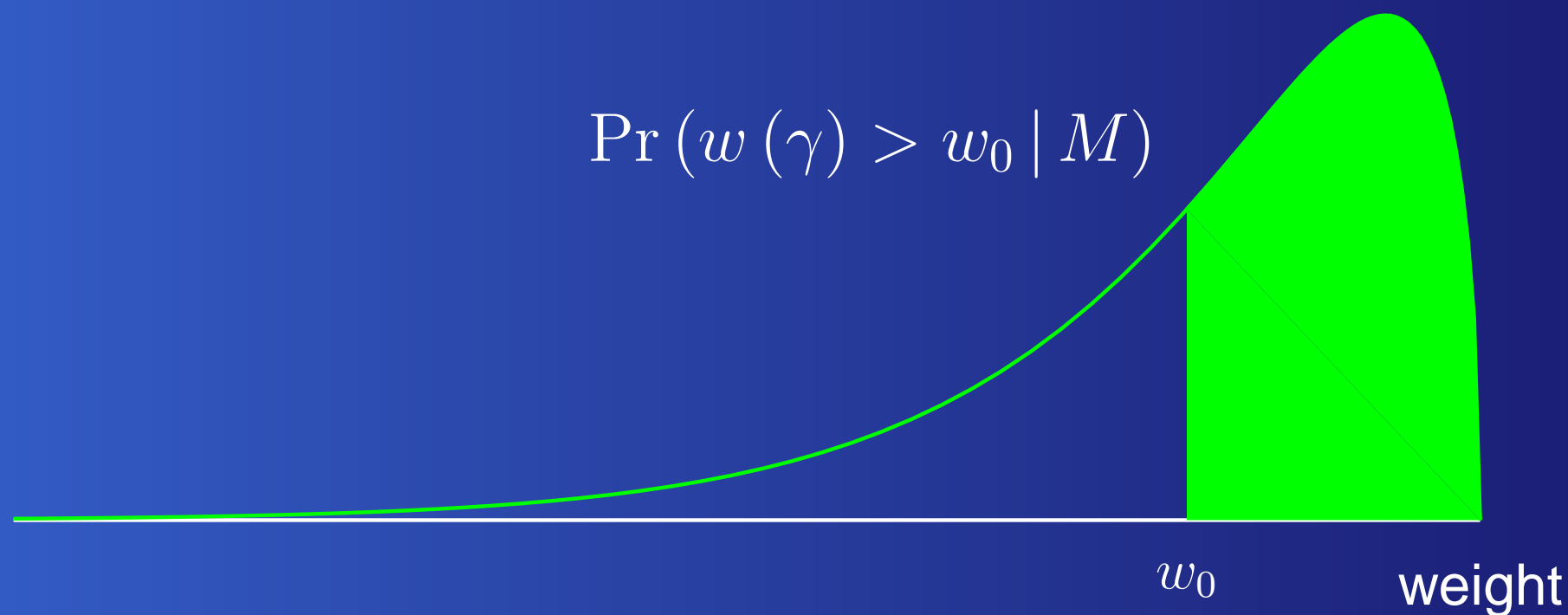U S C E N S U S B U R E A U

# Fundamental Theorem, Cont.

- Theorem: For these error rates $\mu, \lambda$, this is the optimal linkage rule, in the sense of producing the minimum size critical region

- In other words, for a given error bound tolerance, this rule make as many linkage decisions as possible
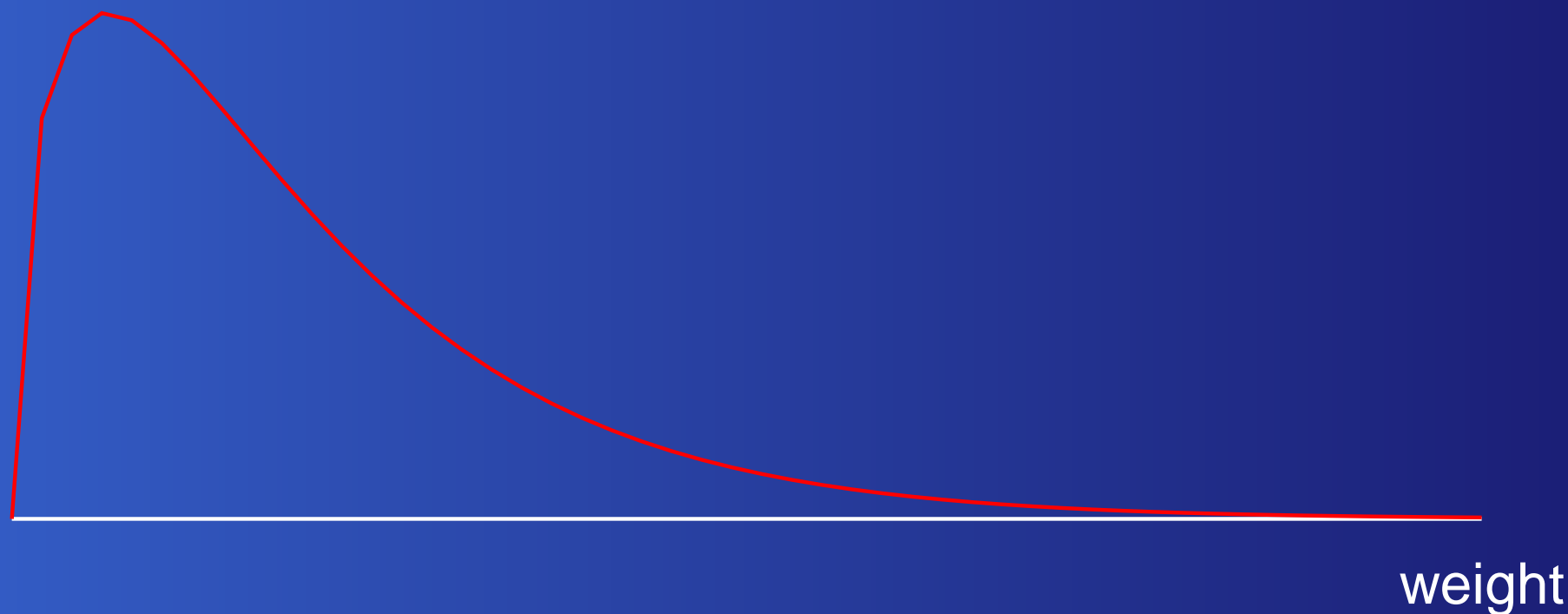
# Weight Distribution for Matches
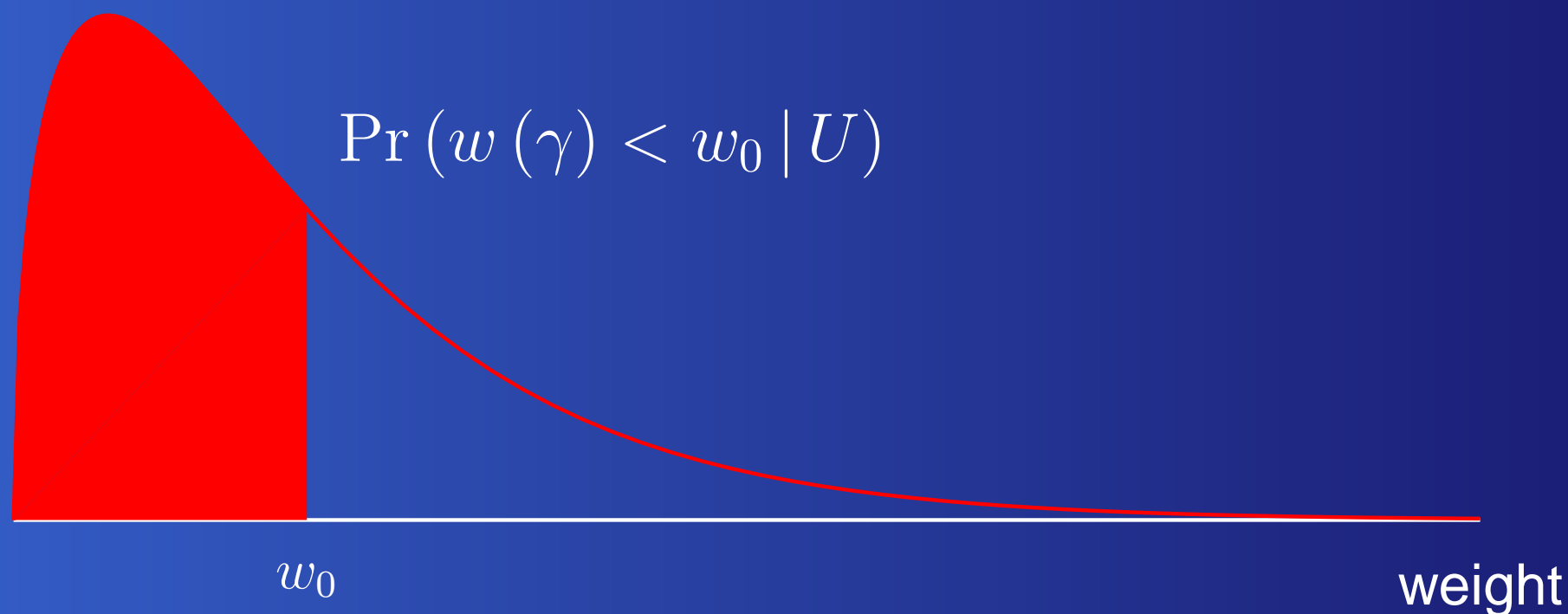
weight

# Weight Distribution for Matches



$$\Pr\left(w\left(\gamma\right) > w_0 \mid M\right)$$

$w_0$

weight

# Weight Distribution for Non-Matches

weight

U S C E N S U S B U R E A U

# Weight Distribution for Non-Matches

$$\Pr\left(w\left(\gamma\right) < w_0 \,|\, U\right)$$

$w_0$

weight

# Idealized Distributions

weight

# Idealized Distributions



$T_\lambda$    $T_\mu$    weight

# Idealized Distributions

Non-Links   Clerical   Links

$T_\lambda$   $T_\mu$   weight

# Error Rates, Clerical Review Region

weight

U S C E N S U S B U R E A U

# Error Rates, Clerical Review Region



$T_\lambda$ $\quad$ $T_\mu$ $\quad$ weight

U S C E N S U S B U R E A U

# Error Rates, Clerical Review Region

$$\lambda = \Pr\left(w\left(\gamma\right) < T_\lambda \mid M\right)$$

$$\mu = \Pr\left(w\left(\gamma\right) > T_\mu \mid U\right)$$

$T_\lambda$

$T_\mu$

weight

U S C E N S U S B U R E A U

# Conditional Independence Assumption

- To facilitate computation of conditional probabilities, Fellegi & Sunter assume conditional independence of comparison vector components

USCENSUSBUREAU

# Conditional Independence Assumption

- To facilitate computation of conditional probabilities, Fellegi & Sunter assume conditional independence of comparison vector components

- For $\gamma = (\gamma_1, \gamma_2, \ldots, \gamma_n)$, assume

$$\Pr(\gamma|M) = \prod_{i=1}^{n} \Pr(\gamma_i|M)$$

$$\Pr(\gamma|U) = \prod_{i=1}^{n} \Pr(\gamma_i|U)$$

USCENSUSBUREAU

# Cond. Indep. Assumption, Cont.

- The factors $\Pr\left(\gamma_i | M\right), \Pr\left(\gamma_i | U\right)$ are called *marginal probabilities*

USCENSUSBUREAU

# Cond. Indep. Assumption, Cont.

- The factors $\Pr(\gamma_i|M), \Pr(\gamma_i|U)$ are called *marginal probabilities*

- The ratio

$$\frac{\Pr(\gamma_i|M)}{\Pr(\gamma_i|U)}$$

determines the *distinguishing power* of the comparison $\gamma_i$

USCENSUSBUREAU

# Cond. Indep. Assumption, Cont.

Under conditional independence assumption, it is convenient to compute the *weight* of the comparison vector

# Cond. Indep. Assumption, Cont.

Under conditional independence assumption, it is convenient to compute the *weight* of the comparison vector

$$
\begin{aligned}
w\left(\gamma\right) &= \log R\left(\gamma\right) \\
&= \sum_{i=1}^{n} \frac{\log \Pr\left(\gamma_i | M\right)}{\log \Pr\left(\gamma_i | U\right)} \\
&= \sum_{i=1}^{n} w\left(\gamma_i\right)
\end{aligned}
$$

USCENSUSBUREAU

# Cond. Indep. Assumption, Cont.

- Motivation: Reduce the number of parameters

# Cond. Indep. Assumption, Cont.

- Motivation: Reduce the number of parameters

- For $n$ binary comparisons and two conditional classes $M, U$, there are $2^{n+1}$ parameters

# Cond. Indep. Assumption, Cont.

- Motivation: Reduce the number of parameters

- For $n$ binary comparisons and two conditional classes $M, U$, there are $2^{n+1}$ parameters

  - $2^n$ comparison vectors

# Cond. Indep. Assumption, Cont.

- 🔴 Motivation: Reduce the number of parameters

- 🔴 For $n$ binary comparisons and two conditional classes $M, U$, there are $2^{n+1}$ parameters

  - 🟢 $2^n$ comparison vectors

  - 🟢 $2$ conditional probabilities for each vector

**U S C E N S U S B U R E A U**

# Cond. Indep. Assumption, Cont.

- Motivation: Reduce the number of parameters

- For $n$ binary comparisons and two conditional classes $M, U$, there are $2^{n+1}$ parameters

  - $2^n$ comparison vectors

  - $2$ conditional probabilities for each vector

- Under conditional independence assumption, there are $2n$ parameters

USCENSUSBUREAU

# Cond. Indep. Assumption, Cont.

- 🔴 Motivation: Reduce the number of parameters

- 🔴 For $n$ binary comparisons and two conditional classes $M, U$, there are $2^{n+1}$ parameters

  - 🟢 $2^n$ comparison vectors

  - 🟢 $2$ conditional probabilities for each vector

- 🔴 Under conditional independence assumption, there are $2n$ parameters

- 🔴 Rationale: Given $M$, errors producing disagreement should be random

## U S C E N S U S B U R E A U

# Cond. Indep. Assumption, Cont.

- Often computable in closed form

# Cond. Indep. Assumption, Cont.

- Often computable in closed form

- Can produce good decision rules even if model inaccurate

# Cond. Indep. Assumption, Cont.

- Often computable in closed form

- Can produce good decision rules even if model inaccurate

- Refered to as *naive Bayes* in machine learning

USCENSUSBUREAU

# Conditional Independence Example

- Suppose

$$\Pr\left(\gamma_1 = 1 | M\right) = 0.9$$
$$\Pr\left(\gamma_2 = 1 | M\right) = 0.8$$
$$\Pr\left(\gamma_3 = 1 | M\right) = 0.7$$

USCENSUSBUREAU

# Conditional Independence Example

- Suppose

$$\Pr\left(\gamma_1 = 1 | M\right) = 0.9$$
$$\Pr\left(\gamma_2 = 1 | M\right) = 0.8$$
$$\Pr\left(\gamma_3 = 1 | M\right) = 0.7$$

- Then for $\gamma = (1, 0, 1)$,

$$\Pr\left(\gamma | M\right) = 0.9 * 0.2 * 0.7 = 0.126$$

## USCENSUSBUREAU

# Conditional Independence Example

- Suppose

$$\Pr\left(\gamma_1 | M\right) = 0.8 \quad \Pr\left(\gamma_1 | U\right) = 0.1$$
$$\Pr\left(\gamma_2 | M\right) = 0.9 \quad \Pr\left(\gamma_2 | U\right) = 0.3$$

# Conditional Independence Example

- Suppose

$$\Pr\left(\gamma_1|M\right) = 0.8 \quad \Pr\left(\gamma_1|U\right) = 0.1$$
$$\Pr\left(\gamma_2|M\right) = 0.9 \quad \Pr\left(\gamma_2|U\right) = 0.3$$

- Then

$$\frac{\Pr\left(\gamma_1|M\right)}{\Pr\left(\gamma_1|U\right)} = 8.0$$

$$\frac{\Pr\left(\gamma_2|M\right)}{\Pr\left(\gamma_2|U\right)} = 3.0$$

U S C E N S U S B U R E A U

# Fellegi-Sunter Summary

- Choose conditional probability parameters

US CENSUS BUREAU

# Fellegi-Sunter Summary

- Choose conditional probability parameters

- Conduct field comparisons on record pairs

USCENSUSBUREAU

# Fellegi-Sunter Summary

- Choose conditional probability parameters

- Conduct field comparisons on record pairs

- Classify record pairs based on weight of comparison vector

US CENSUS BUREAU

# Record Linkage Methodology

- Parameter estimation

# Record Linkage Methodology

- Parameter estimation
  - EM Algorithm

U S C E N S U S B U R E A U

# Record Linkage Methodology

- Parameter estimation
  - EM Algorithm
- Blocking

# Choosing Parameters

- Informal

# Choosing Parameters

- Informal

- EM Algorithm

U S C E N S U S B U R E A U

# Choosing Parameters

- Informal

- EM Algorithm

- Other Methods

# Informal Methods

- Guess

# Informal Methods

- Guess

$$0 < \Pr\left(\gamma | U\right) < \Pr\left(\gamma | M\right) < 1$$

# Informal Methods

- Guess

$$0 < \Pr\left(\gamma|U\right) < \Pr\left(\gamma|M\right) < 1$$

- Approximate

# Informal Methods

🔴 Guess

$$0 < \Pr\left(\gamma | U\right) < \Pr\left(\gamma | M\right) < 1$$

🔴 Approximate

$$\Pr\left(\gamma | U\right) \approx \Pr\left(\gamma | S\right)$$

US CENSUS BUREAU

# Informal Methods

- Guess

$$0 < \Pr\left(\gamma | U\right) < \Pr\left(\gamma | M\right) < 1$$

- Approximate

$$\Pr\left(\gamma | U\right) \approx \Pr\left(\gamma | S\right)$$

- Iterate

USCENSUSBUREAU

# Informal Methods

🔴 Guess

$$0 < \Pr\left(\gamma|U\right) < \Pr\left(\gamma|M\right) < 1$$

🔴 Approximate

$$\Pr\left(\gamma|U\right) \approx \Pr\left(\gamma|S\right)$$

🔴 Iterate

🟢 Perform matching with current parameters

USCENSUSBUREAU

# Informal Methods

- 🔴 Guess

$$0 < \Pr\left(\gamma|U\right) < \Pr\left(\gamma|M\right) < 1$$

- 🔴 Approximate

$$\Pr\left(\gamma|U\right) \approx \Pr\left(\gamma|S\right)$$

- 🔴 Iterate

  - 🟢 Perform matching with current parameters
  - 🟢 Review results

## USCENSUSBUREAU

# Informal Methods

- 🔴 Guess

$$0 < \Pr\left(\gamma | U\right) < \Pr\left(\gamma | M\right) < 1$$

- 🔴 Approximate

$$\Pr\left(\gamma | U\right) \approx \Pr\left(\gamma | S\right)$$

- 🔴 Iterate

  - 🟢 Perform matching with current parameters

  - 🟢 Review results

  - 🟢 Adjust parameters based on observation

U S C E N S U S B U R E A U

# EM Algorithm

- Dempster, Laird, Rubin. "Maximum likelihood from incomplete data via the EM algorithm". *Journal of the Royal Statistical Society. SeriesB.* **39**. pp. 1–39. 1977.

# EM Algorithm

- Dempster, Laird, Rubin. "Maximum likelihood from incomplete data via the EM algorithm". *Journal of the Royal Statistical Society. SeriesB.* **39**. pp. 1–39. 1977.

- McLachlan, Krishnan. *The EM Algorithm and Extensions.* Wiley-Interscience. 2nd Ed. 2007.

U S C E N S U S B U R E A U

# EM Algorithm

- Dempster, Laird, Rubin. "Maximum likelihood from incomplete data via the EM algorithm". *Journal of the Royal Statistical Society. SeriesB.* **39**. pp. 1–39. 1977.

- McLachlan, Krishnan. *The EM Algorithm and Extensions.* Wiley-Interscience. 2nd Ed. 2007.

- Maximum likelihood method

# EM Algorithm

- Dempster, Laird, Rubin. "Maximum likelihood from incomplete data via the EM algorithm". *Journal of the Royal Statistical Society.* SeriesB. **39**. pp. 1–39. 1977.

- McLachlan, Krishnan. *The EM Algorithm and Extensions.* Wiley-Interscience. 2$^{nd}$ Ed. 2007.

- Maximum likelihood method

- Latent class

## USCENSUSBUREAU

# EM Algorithm

- Dempster, Laird, Rubin. "Maximum likelihood from incomplete data via the EM algorithm". *Journal of the Royal Statistical Society.* SeriesB. **39**. pp. 1–39. 1977.

- McLachlan, Krishnan. *The EM Algorithm and Extensions.* Wiley-Interscience. 2nd Ed. 2007.

- Maximum likelihood method

- Latent class

- Mixture model

# Likelihood Function

- $$L = \prod_{(a,b) \in S} \Pr\left(\gamma\left(a, b\right)\right)$$

$$= \prod_{j} \left(\Pr\left(\gamma^j | M\right) \Pr\left(M\right) + \Pr\left(\gamma^j | U\right) \Pr\left(U\right)\right)^{n_j}$$

USCENSUSBUREAU

# Likelihood Function

$$L = \prod_{(a,b) \in S} \mathrm{Pr}\left(\gamma\left(a,b\right)\right)$$

$$= \prod_{j} \left(\mathrm{Pr}\left(\gamma^j | M\right) \mathrm{Pr}\left(M\right) + \mathrm{Pr}\left(\gamma^j | U\right) \mathrm{Pr}\left(U\right)\right)^{n_j}$$

$$n_j = \left|\left\{(a,b) \in S \mid \gamma\left(a,b\right) = \gamma^j\right\}\right|$$

USCENSUSBUREAU

# Complete-data Likelihood Function

🔴 Consider

$$\chi_j = \begin{cases} 1 & \text{if } (a,b)^j \in M \\ 0 & \text{if } (a,b)^j \in U \end{cases}$$

$$X_j = \sum_{\gamma(a,b)=\gamma^j} \chi_j\,(a,b)$$

USCENSUSBUREAU

# Complete-data Likelihood Function

- Consider

$$\chi_j = \begin{cases} 1 & \text{if } (a,b)^j \in M \\ 0 & \text{if } (a,b)^j \in U \end{cases}$$

$$X_j = \sum_{\gamma(a,b)=\gamma^j} \chi_j(a,b)$$

- Then

$$L = \prod_j \left( \left( \Pr\left(\gamma^j | M\right) \Pr(M) \right)^{\overline{X}_j} \left( \Pr\left(\gamma^j | U\right) \Pr(U) \right)^{1-\overline{X}_j} \right)^{n_j}$$

USCENSUSBUREAU

# Expectation Step

- Given current estimates of conditional probabilities and $\Pr(M), \Pr(U)$, compute

# Expectation Step

- Given current estimates of conditional probabilities and $\Pr(M), \Pr(U)$, compute

$$E\left(\overline{X}^j\right) = \Pr\left(M|\gamma^j\right)$$

$$= \frac{\Pr\left(\gamma^j|M\right)\Pr(M)}{\Pr\left(\gamma^j|M\right)\Pr(M) + \Pr\left(\gamma^j|U\right)\Pr(U)}$$

$$= \hat{X}^j$$

U S C E N S U S B U R E A U

# Maximization Step

- Given unobserved data estimates $\hat{X}^j$, compute probabilities $\Pr\left(\gamma^j|M\right)$, $\Pr\left(\gamma^j|U\right)$, $\Pr\left(M\right)$, $\Pr\left(U\right)$ maximizing

# Maximization Step

- Given unobserved data estimates $\hat{X}^j$, compute probabilities $\mathrm{Pr}\left(\gamma^j|M\right)$, $\mathrm{Pr}\left(\gamma^j|U\right)$, $\mathrm{Pr}\left(M\right)$, $\mathrm{Pr}\left(U\right)$ maximizing

$$\log L =$$

$$\sum_j n_j \left( \hat{X}^j \left( \log \mathrm{Pr}\left(\gamma^j|M\right) + \log \mathrm{Pr}\left(M\right)\right) \right.$$

$$\left. + \left(1 - \hat{X}^j\right) \left( \log \mathrm{Pr}\left(\gamma^j|U\right) + \log \mathrm{Pr}\left(U\right)\right)\right)$$

# Max Step, Cont.

- Under conditional independence

# Max Step, Cont.

- Under conditional independence

$$\log L =$$

$$\sum_j n_j \left( \sum_i \hat{X}^j \left( \log \Pr\left(\gamma_i^j | M\right) + \log \Pr\left(M\right) \right) \right.$$

$$\left. + \left(1 - \hat{X}^j\right) \left( \sum_i \log \Pr\left(\gamma_i^j | U\right) + \log \Pr\left(U\right) \right) \right)$$

# Max Step, Cont.

For

$$n = \sum_j n_j$$

estimate

$$\Pr\left(M\right) = \frac{1}{n} \sum_j n_j \bar{X}^j$$

# Max Step, Cont.

- Let

$$
k_i^j = \begin{cases} 1 & \text{if } \gamma_i^j = 1 \\ 0 & \text{if } \gamma_i^j = 0 \end{cases}
$$

# Max Step, Cont.

- Let

$$k_i^j = \begin{cases} 1 & \text{if } \gamma_i^j = 1 \\ 0 & \text{if } \gamma_i^j = 0 \end{cases}$$

- and estimate

$$\Pr\left(\gamma_i | M\right) = \frac{1}{n} \sum_j n_j \bar{X}^j k_i^j$$

# EM Algorithm

1. Initialize with probability values

# EM Algorithm

1. Initialize with probability values

2. Iterate

U S C E N S U S B U R E A U

# EM Algorithm

1. Initialize with probability values

2. Iterate

   (a) Expectation Step

# EM Algorithm

1. Initialize with probability values

2. Iterate
   (a) Expectation Step
   (b) Maximization Step

USCENSUSBUREAU

# EM Algorithm

1. Initialize with probability values

2. Iterate
   (a) Expectation Step
   (b) Maximization Step

3. Until convergence of likelihood function

USCENSUSBUREAU

# EM Algorithm Remarks

- Each EM iteration increases likelihood, so algorithm converges to a (local) maximum

# EM Algorithm Remarks

- Each EM iteration increases likelihood, so algorithm converges to a (local) maximum

- For this conditional independence model, convergence is efficient and generally insensitive to initial data

US CENSUS BUREAU

# EM Algorithm Remarks

- Each EM iteration increases likelihood, so algorithm converges to a (local) maximum

- For this conditional independence model, convergence is efficient and generally insensitive to initial data

- For latent class to be numerically detected, it must be represented by about 5% of the total record pair data

U S C E N S U S B U R E A U

# EM Algorithm Remarks

- Each EM iteration increases likelihood, so algorithm converges to a (local) maximum

- For this conditional independence model, convergence is efficient and generally insensitive to initial data

- For latent class to be numerically detected, it must be represented by about 5% of the total record pair data

- Check: Do $\Pr(M), \Pr(U)$ seem reasonable?

USCENSUSBUREAU

# EM Remarks, Cont.

- If $\mathrm{Pr}\left(M\right), \mathrm{Pr}\left(U\right)$ are off, everything is off

# EM Remarks, Cont.

- If $\Pr(M), \Pr(U)$ are off, everything is off

- We can extend problem to comparisons taking on more that 2 values

# EM Remarks, Cont.

- If $\Pr(M), \Pr(U)$ are off, everything is off

- We can extend problem to comparisons taking on more that 2 values

  - Creates more pattern types and probability parameters

# EM Remarks, Cont.

- If $\Pr(M), \Pr(U)$ are off, everything is off

- We can extend problem to comparisons taking on more that 2 values
  - Creates more pattern types and probability parameters

- Can extend algorithm to more that 2 classes

# EM Remarks, Cont.

- If $\Pr(M), \Pr(U)$ are off, everything is off

- We can extend problem to comparisons taking on more that 2 values

  - Creates more pattern types and probability parameters

- Can extend algorithm to more that 2 classes

  - Increases number of parameters to be estimated

U S C E N S U S B U R E A U

# Blocking

- If set $A$ contains $m$ records and set $B$ contains $n$ records then $A \times B$ contains $mn$ record pairs

# Blocking

- If set $A$ contains $m$ records and set $B$ contains $n$ records then $A \times B$ contains $mn$ record pairs

- It is computationally inefficient to compare all record pairs

USCENSUSBUREAU

# Blocking

- If set $A$ contains $m$ records and set $B$ contains $n$ records then $A \times B$ contains $mn$ record pairs

- It is computationally inefficient to compare all record pairs

- In practice, just bring together record pairs that agree on some chosen features (blocking criterion)

U S C E N S U S B U R E A U

# Blocking

- If set $A$ contains $m$ records and set $B$ contains $n$ records then $A \times B$ contains $mn$ record pairs

- It is computationally inefficient to compare all record pairs

- In practice, just bring together record pairs that agree on some chosen features (blocking criterion)

- Generally repeat record linkage procedure for several different blocking criteria

# Blocking Criteria

- Geographic codes

# Blocking Criteria

- Geographic codes

- Postal or phone codes

US CENSUS BUREAU

# Blocking Criteria

- Geographic codes

- Postal or phone codes

- Name prefix

# Blocking Criteria

- Geographic codes

- Postal or phone codes

- Name prefix

- Phonetic name codes

# Blocking Criteria

- Geographic codes

- Postal or phone codes

- Name prefix

- Phonetic name codes
  - Soundex

# Blocking Criteria

- Geographic codes

- Postal or phone codes

- Name prefix

- Phonetic name codes
  - Soundex
  - NYSIIS

# Blocking Criteria

- Geographic codes
- Postal or phone codes
- Name prefix
- Phonetic name codes
  - Soundex
  - NYSIIS
- Combinations

USCENSUSBUREAU

# Record Linkage Refinements

- String comparator

# Record Linkage Refinements

- String comparator
- Third latent class

U S C E N S U S B U R E A U

# Record Linkage Refinements

- String comparator
- Third latent class
- Third comparison value

# Record Linkage Refinements

- String comparator

- Third latent class

- Third comparison value

- One-to-one matching

# String Comparator

- For some comparisons (*e.g.* categorical variables), it is sufficient to assign agree/disagree

# String Comparator

- For some comparisons (*e.g.* categorical variables), it is sufficient to assign agree/disagree

- For string variables (*e.g.* first names, last names, street names) this is probably too restrictive

# String Comparator

- For some comparisons (*e.g.* categorical variables), it is sufficient to assign agree/disagree

- For string variables (*e.g.* first names, last names, street names) this is probably too restrictive

- A string comparator allows us to assign comparison values between full agreement and full disagreement

# String Comparator Context

- Binary comparison $\gamma \in \{0, 1\}$

US CENSUS BUREAU

# String Comparator Context

- Binary comparison $\gamma \in \{0, 1\}$

- Weight assignment

$$a_w = \log \frac{\Pr(\gamma = 1 | M)}{\Pr(\gamma = 1 | U)}$$

$$d_w = \log \frac{\Pr(\gamma = 0 | M)}{\Pr(\gamma = 0 | U)}$$

$$d_w < 0 < a_w$$

USCENSUSBUREAU

# String Comparator Context, Cont.

- For alphabet $\Sigma$, our string comparator is a *similarity function*

$$\gamma : \Sigma^* \times \Sigma^* \to [0, 1]$$

$$\gamma\left(\alpha, \beta\right) = 1 \text{ if } \alpha = \beta$$

USCENSUSBUREAU

# String Comparator Context, Cont.

- For alphabet $\Sigma$, our string comparator is a *similarity function*

$$\gamma : \Sigma^* \times \Sigma^* \rightarrow [0, 1]$$

$$\gamma(\alpha, \beta) = 1 \text{ if } \alpha = \beta$$

- Weight assignment function $w$ is an increasing interpolation function

$$w : [0, 1] \rightarrow [d_w, a_w]$$

$$w(1) = a_w$$

U S C E N S U S B U R E A U

# Some String Comparator Types

- Bigram, *n*-gram

# Some String Comparator Types

- Bigram, *n*-gram

- Jaro-Winkler

U S C E N S U S B U R E A U

# Some String Comparator Types

- Bigram, *n*-gram

- Jaro-Winkler

- Edit distance

USCENSUSBUREAU

# Bigrams

- Decompose string into a set of 2-character (contiguous) substrings

$$alphabet \rightarrow \{al, lp, ph, ha, ab, be, et\}$$

# Bigrams

- Decompose string into a set of 2-character (contiguous) substrings

$$alphabet \rightarrow \{al, lp, ph, ha, ab, be, et\}$$

- For alphabet of $s = |\Sigma|$ characters, record bigram counts in a vector of dimension $s^2$

# Bigrams, Cont.

- Two strings can be compared by computing the "angle" between their bigram vectors $a, b$

$$\cos \theta = \frac{a \cdot b}{|a|\,|b|}$$

USCENSUSBUREAU

# Bigrams, Cont.

- Two strings can be compared by computing the "angle" between their bigram vectors $a, b$

$$\cos \theta = \frac{a \cdot b}{|a|\,|b|}$$

- Obvious generalization to $n$-grams

U S C E N S U S B U R E A U

# Bigrams, Cont.

- Two strings can be compared by computing the "angle" between their bigram vectors $a, b$

$$\cos \theta = \frac{a \cdot b}{|a| \, |b|}$$

- Obvious generalization to $n$-grams
- Vector for $n$-gram is in $s^n$ dimensional space

USCENSUSBUREAU

# Bigrams, Cont.

- Computation algorithm is fast (linear)

# Bigrams, Cont.

- Computation algorithm is fast (linear)
- Don't work very well for record linkage

# Bigrams, Cont.

- Computation algorithm is fast (linear)
- Don't work very well for record linkage
  - Ignores order of bigram occurrence

$$abcba \approx bcbab$$

# Bigrams, Cont.

- Computation algorithm is fast (linear)
- Don't work very well for record linkage
  - Ignores order of bigram occurrence

$$abcba \approx bcbab$$

  - Works better for small alphabet, long strings than *vice versa*

USCENSUSBUREAU

# Jaro-Winkler Comparator

- In the following, let
  $\alpha = (a_1, a_2, \ldots a_m)$, $\beta = (b_1, b_2, \ldots, b_n)$ be
  strings of lengths $m, n$ respectively with
  $m \leq n$

# Jaro-Winkler Comparator

- In the following, let $\alpha = (a_1, a_2, \ldots a_m), \beta = (b_1, b_2, \ldots, b_n)$ be strings of lengths $m, n$ respectively with $m \le n$

- Comparator value depends on number of common characters and character "transpositions"

USCENSUSBUREAU

# Jaro-Winkler Comparator, Cont.

- Strings $\alpha, \beta$ have common characters $a_i, b_j$ iff

$$
\begin{aligned}
a_i &= b_j \\
|i - j| &< \left\lfloor \frac{n}{2} \right\rfloor
\end{aligned}
$$

# Jaro-Winkler Comparator, Cont.

- Strings $\alpha, \beta$ have common characters $a_i, b_j$ iff

$$
\begin{aligned}
a_i &= b_j \\
|i - j| &< \left\lfloor \frac{n}{2} \right\rfloor
\end{aligned}
$$

- The number of transpositions is computed as the greatest integer of half of the number of out-of-order common character pairs

USCENSUSBUREAU

# Jaro-Winkler Comparator, Cont.

- For string pair with $c$ common characters and $t$ transpositions, basis similarity score is

$$x = \frac{1}{3}\left(\frac{c}{m} + \frac{c}{n} + \frac{c-t}{c}\right)$$

USCENSUSBUREAU

# Jaro-Winkler Example

- Consider the strings $(b,a,r,n,e,s)$ and $(a,n,d,e,r,s,o,n)$

# Jaro-Winkler Example

- Consider the strings (*b,a,r,n,e,s*) and (*a,n,d,e,r,s,o,n*)

- Search range $d$

$$
\begin{aligned}
n &= 8 \\
d &= \left\lfloor \frac{8}{2} \right\rfloor - 1 = 3
\end{aligned}
$$

# Jaro-Winkler Example

- Consider the strings (*b,a,r,n,e,s*) and (*a,n,d,e,r,s,o,n*)

- Search range $d$

$$
\begin{aligned}
n &= 8 \\
d &= \left\lfloor \frac{8}{2} \right\rfloor - 1 = 3
\end{aligned}
$$

- Common characters

$$(a, r, n, e, s)$$
$$(a, n, e, r, s)$$

US CENSUS BUREAU

# Jaro-Winkler Example, Cont.

- Five common characters with 3 out of order, so $c = 5, t = 1$

# Jaro-Winkler Example, Cont.

- Five common characters with 3 out of order, so $c = 5, t = 1$

- Score

$$x = \frac{1}{3} \left( \frac{5}{6} + \frac{5}{8} + \frac{4}{5} \right) = \frac{271}{360} \doteq 0.75280$$

USCENSUSBUREAU

# Jaro-Winkler Variations

- Similar characters

U S C E N S U S B U R E A U

# Jaro-Winkler Variations

- Similar characters

- Prefix adjustment

# Jaro-Winkler Variations

- Similar characters

- Prefix adjustment

- Long suffix adjustment

# Similar Characters

- Attempt to compensate for common misspellings or typos

# Similar Characters

- Attempt to compensate for common misspellings or typos

- List of 36 pairs of characters deemed similar (*e.g.* most vowell pairs)

# Similar Characters

- Attempt to compensate for common misspellings or typos

- List of 36 pairs of characters deemed similar (*e.g.* most vowell pairs)

- After common characters designated, remaining characters checked for similar pairs

USCENSUSBUREAU

# Similar Characters

- Attempt to compensate for common misspellings or typos

- List of 36 pairs of characters deemed similar (*e.g.* most vowell pairs)

- After common characters designated, remaining characters checked for similar pairs

- Each similar pair is scored as 0.3 of a common pair

USCENSUSBUREAU

# Similar Characters, Cont.

- Revised character count

$$c_s = c + 0.3s$$

# Similar Characters, Cont.

- Revised character count

$$c_s = c + 0.3s$$

- Adjusted comparator score

$$x_s = \frac{1}{3}\left(\frac{c_s}{m} + \frac{c_s}{n} + \frac{c - t}{c}\right)$$

USCENSUSBUREAU

# Similar Characters, Cont.

- For example. *abc* and *ebc* have 2 common characters and the remaining pair (*a,e*) are similar, so

$$
\begin{aligned}
x_s & = \frac{1}{3}\left(\frac{2}{3} + \frac{2}{3} + 1\right) + \frac{1}{3}\left(\frac{0.3}{3} + \frac{0.3}{3}\right) \\
& = \frac{7}{9} + \frac{1}{15} \\
& = \frac{38}{45}
\end{aligned}
$$

USCENSUSBUREAU

# Common Prefix

- Spelling mistakes tend to occur later in the string (Winkler)

# Common Prefix

- Spelling mistakes tend to occur later in the string (Winkler)

- Check for common prefix of up to 4 characters

# Common Prefix

- Spelling mistakes tend to occur later in the string (Winkler)

- Check for common prefix of up to 4 characters

- If length of common prefix is $p$, adjust score $x$ by

$$x_p = x + \frac{p\,(1 - x)}{10}$$

USCENSUSBUREAU

# Long String Adjustment

- Adjust score for longer strings with several common characters beyond common prefix

# Long String Adjustment

- Adjust score for longer strings with several common characters beyond common prefix

- Conditions for using the adjustment

USCENSUSBUREAU

# Long String Adjustment

- Adjust score for longer strings with several common characters beyond common prefix

- Conditions for using the adjustment
  1. $m \geq 5$

USCENSUSBUREAU

# Long String Adjustment

- Adjust score for longer strings with several common characters beyond common prefix

- Conditions for using the adjustment
  1. $m \geq 5$
  2. $c - p \geq 2$

# Long String Adjustment

- Adjust score for longer strings with several common characters beyond common prefix

- Conditions for using the adjustment
  1. $m \geq 5$
  2. $c - p \geq 2$
  3. $c - p \geq \frac{m-p}{2}$

# Long String Adjustment, Cont.

- That is,

# Long String Adjustment, Cont.

- That is,

  1. Both strings are at least 5 characters long

# Long String Adjustment, Cont.

- That is,

  1. Both strings are at least 5 characters long
  2. There are at least two common characters besides the agreeing prefix characters

# Long String Adjustment, Cont.

- That is,

  1. Both strings are at least 5 characters long
  2. There are at least two common characters besides the agreeing prefix characters
  3. We want the strings outside the common prefixes to be fairly rich in common characters, so that the remaining common characters are at least half of the remaining common characters of the shorter string

USCENSUSBUREAU

# Long String Adjustment, Cont.

- If conditions met, then adjust score by

$$x_l = x + (1 - x) \, \frac{c - (p + 1)}{m + n - 2 \, (p - 1)}$$

US CENSUS BUREAU

# Long String Adjustment, Cont.

- In *barnes, anderson* example, conditions are met, so the adjusted score is

$$
\begin{aligned}
x_l &= \frac{271}{360} + \left(1 - \frac{271}{360}\right) \frac{5-1}{6+8+2} \\
&= \frac{391}{480} \\
&\doteq 0.8146
\end{aligned}
$$

# Jaro-Winkler Comparator

- Slower algorithm (quadratic)

U S C E N S U S B U R E A U

# Jaro-Winkler Comparator

- Slower algorithm (quadratic)
- Performs very well in tests

USCENSUSBUREAU

# Edit Distance String Comparators

- The minimum number of edits required to convert sting $\alpha$ to string $\beta$, lengths $m \leq n$

# Edit Distance String Comparators

- 🔴 The minimum number of edits required to convert sting $\alpha$ to string $\beta$, lengths $m \leq n$
  - 🟢 Insert

USCENSUSBUREAU

# Edit Distance String Comparators

- The minimum number of edits required to convert sting $\alpha$ to string $\beta$, lengths $m \leq n$
  - Insert
  - Delete

USCENSUSBUREAU

# Edit Distance String Comparators

- The minimum number of edits required to convert sting $\alpha$ to string $\beta$, lengths $m \leq n$
  - Insert
  - Delete
  - Substitute

# Edit Distance String Comparators

- The minimum number of edits required to convert sting $\alpha$ to string $\beta$, lengths $m \leq n$
  - Insert
  - Delete
  - Substitute
- Dynamic programming algorithm, quadratic complexity $O(mn)$

USCENSUSBUREAU

# Edit Distance Algorithm

- For $\alpha_i$ prefix of $\alpha$ of length $i$, $\beta_j$ prefix of $\beta$ of length $j$

# Edit Distance Algorithm

- For $\alpha_i$ prefix of $\alpha$ of length $i$, $\beta_j$ prefix of $\beta$ of length $j$

- Initialize

$$
\begin{aligned}
e\left(\alpha_i, \varepsilon\right) &= i \\
e\left(\varepsilon, \beta_j\right) &= j \\
e\left(\varepsilon, \varepsilon\right) &= 0
\end{aligned}
$$

USCENSUSBUREAU

# Edit Distance Algorithm, Cont.

- Compute

$$e\left(\alpha_i, \beta_j\right) = \min \begin{cases} e\left(\alpha_{i-1}, \beta_j\right) + 1 \\ e\left(\alpha_i, \beta_{j-1}\right) + 1 \\ \begin{cases} e\left(\alpha_{i-1}, \beta_{j-1}\right) & \text{if } a_i = b_j \\ e\left(\alpha_{i-1}, \beta_{j-1}\right) + 1 & \text{if } a_i \neq b_j \end{cases} \end{cases}$$

USCENSUSBUREAU

# Edit Distance Algorithm, Cont.

- Compute

$$e\left(\alpha_i, \beta_j\right) = \min \begin{cases} e\left(\alpha_{i-1}, \beta_j\right) + 1 \\ e\left(\alpha_i, \beta_{j-1}\right) + 1 \\ \begin{cases} e\left(\alpha_{i-1}, \beta_{j-1}\right) & \text{if } a_i = b_j \\ e\left(\alpha_{i-1}, \beta_{j-1}\right) + 1 & \text{if } a_i \neq b_j \end{cases} \end{cases}$$

- Distance

$$e = e\left(\alpha, \beta\right) = e\left(\alpha_m, \beta_n\right)$$

USCENSUSBUREAU

# Edit Distance Similarity Function

- Edit distance is a metric

# Edit Distance Similarity Function

- Edit distance is a metric

- Similarity function

$$x_e = 1 - \frac{e}{n}$$

# Edit Distance Example

- For example, for *barnes, anderson*, have possible minimal edit path

$$(b, \varepsilon)\, a\, (r, n)\, (n, d)\, e\, (\varepsilon, r)\, s\, (\varepsilon, o)\, (\varepsilon, n)$$

# Edit Distance Example

- For example, for *barnes, anderson*, have possible minimal edit path

$$(b, \varepsilon) \, a \, (r, n) \, (n, d) \, e \, (\varepsilon, r) \, s \, (\varepsilon, o) \, (\varepsilon, n)$$

- So

$$x_e = 1 - \frac{6}{8} = \frac{1}{4}$$

USCENSUSBUREAU

# Edit Distance Example

- For example, for *barnes, anderson*, have possible minimal edit path

$$(b, \varepsilon)\, a\, (r, n)\, (n, d)\, e\, (\varepsilon, r)\, s\, (\varepsilon, o)\, (\varepsilon, n)$$

- So

$$x_e = 1 - \frac{6}{8} = \frac{1}{4}$$

- Note order of characters very important

## U S C E N S U S B U R E A U

# Longest Common Subsequence

- Length of longest common subsequence (*lcs*)

# Longest Common Subsequence

- Length of longest common subsequence (*lcs*)

- Similar dynamic programming algorithm, without substitutions

$$
e\left(\alpha_i, \beta_j\right) = \min \begin{cases} e\left(\alpha_{i-1}, \beta_j\right) + 1 \\ e\left(\alpha_i, \beta_{j-1}\right) + 1 \\ e\left(\alpha_{i-1}, \beta_{j-1}\right) \quad \text{if } a_i = b_j \end{cases}
$$

# LCS Similarity Function

- Similarity function

$$x_c = \frac{l}{m}$$

# LCS Similarity Function

- Similarity function

$$x_c = \frac{l}{m}$$

- Example *lcs*$=(a, n, e, s)$, similarity score

$$x_c = \frac{4}{6} = \frac{2}{3}$$

USCENSUSBUREAU

# Combination Similarity Function

- Compute both edit distance and *lcs*

U S C E N S U S B U R E A U

# Combination Similarity Function

- Compute both edit distance and *lcs*

- Combined score

$$x_{ec} = \frac{1}{2}\left(\left(1 - \frac{e}{n}\right) + \frac{l}{m}\right)$$

USCENSUSBUREAU

# Combination Similarity Function

- Compute both edit distance and *lcs*

- Combined score

$$x_{ec} = \frac{1}{2}\left(\left(1 - \frac{e}{n}\right) + \frac{l}{m}\right)$$

- Example

$$x_{ec} = \frac{1}{2}\left(\frac{1}{4} + \frac{2}{3}\right) = \frac{11}{24} \doteq 0.4583$$

USCENSUSBUREAU

# Evaluating String Comparators

- Yancey, "Evaluating String Comparator Performance for Record Linkage," 2005, http://www.census.gov/srd/www/byname.html

U S C E N S U S B U R E A U

# Evaluating String Comparators

- Yancey, "Evaluating String Comparator Performance for Record Linkage," 2005, http://www.census.gov/srd/www/byname.html

- Compare performance

U S C E N S U S B U R E A U

# Evaluating String Comparators

- Yancey, "Evaluating String Comparator Performance for Record Linkage," 2005, http://www.census.gov/srd/www/byname.html

- Compare performance
  - Jaro-Winkler

USCENSUSBUREAU

# Evaluating String Comparators

- Yancey, "Evaluating String Comparator Performance for Record Linkage," 2005, http://www.census.gov/srd/www/byname.html

- Compare performance
  - Jaro-Winkler
  - Edit distance

U S C E N S U S B U R E A U

# Evaluating String Comparators, Cont.

- Jaro-Winkler, with and without modifications

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Jaro-Winkler, with and without modifications
  - Prefix adjustment

# Evaluating String Comparators, Cont.

- Jaro-Winkler, with and without modifications
  - Prefix adjustment
  - Similar characters

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Jaro-Winkler, with and without modifications
    - Prefix adjustment
    - Similar characters
    - Long suffix adjustment

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Edit distance

U S C E N S U S B U R E A U

# Evaluating String Comparators, Cont.

- Edit distance
  - Edit distance similarity

# Evaluating String Comparators, Cont.

- Edit distance
  - Edit distance similarity
  - Markov edit distance (J. Wei. "Markov Edit Distance". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol 26, No. 3, pp. 311–321, 2004)

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Edit distance
  - Edit distance similarity
  - Markov edit distance (J. Wei. "Markov Edit Distance". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol 26, No. 3, pp. 311–321, 2004)
  - With and without *lcs*

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Lots of data

# Evaluating String Comparators, Cont.

- Lots of data

- Truth decks from 1990 and 2000 U.S. Census

# Evaluating String Comparators, Cont.

- Lots of data

- Truth decks from 1990 and 2000 U.S. Census

- $M$: All non-identical, non-blank names from matched record pairs

USCENSUSBUREAU

# Evaluating String Comparators, Cont.

- Lots of data

- Truth decks from 1990 and 2000 U.S. Census

- $M$: All non-identical, non-blank names from matched record pairs

- $U$: All cross pairs of these names

USCENSUSBUREAU

# Results of String Comparator Evaluatio

- Jaro-Winkler did well

# Results of String Comparator Evaluatio

- Jaro-Winkler did well
  - Prefix adjustment always helps

# Results of String Comparator Evaluation

- Jaro-Winkler did well
  - Prefix adjustment always helps
  - Similar character adjustment generally helps a bit

USCENSUSBUREAU

# Results of String Comparator Evaluatio

- Jaro-Winkler did well
  - Prefix adjustment always helps
  - Similar character adjustment generally helps a bit
  - Long suffix adjustment sometime helps a little

USCENSUSBUREAU

# Results of String Comparator Evaluati[on]

- Adding *lcs* significantly improves edit distance and Markov edit distance

# Results of String Comparator Evaluati...

- Adding *lcs* significantly improves edit distance and Markov edit distance

- Edit distance always better than Markov edit distance

U S C E N S U S B U R E A U

# Results of String Comparator Evaluati

- Adding *lcs* significantly improves edit distance and Markov edit distance

- Edit distance always better than Markov edit distance

- Jaro-Winkler (full) comparable to edit distance/*lcs*

USCENSUSBUREAU

# Results of String Comparator Evaluatio

- Adding *lcs* significantly improves edit distance and Markov edit distance

- Edit distance always better than Markov edit distance

- Jaro-Winkler (full) comparable to edit distance/*lcs*
  - Usually

USCENSUSBUREAU

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

- For Jaro-Winkler

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

- For Jaro-Winkler
  - $s\left(\alpha, \alpha\beta\right) = \frac{5}{6}$

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

- For Jaro-Winkler
  - $s\left(\alpha, \alpha\beta\right) = \frac{5}{6}$
  - In $n \geq 4$, with prefix adjustment, $s\left(\alpha, \alpha\beta\right) = \frac{9}{10}$

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

- For Jaro-Winkler
  - $s\left(\alpha, \alpha\beta\right) = \frac{5}{6}$
  - In $n \geq 4$, with prefix adjustment, $s\left(\alpha, \alpha\beta\right) = \frac{9}{10}$
  - $s\left(\beta, \alpha\beta\right) = 0$

# Jaro-Winkler Anomaly

- Let $\alpha, \beta$ be strings of length $n$ with no common characters

- For Jaro-Winkler
  - $s\left(\alpha, \alpha\beta\right) = \frac{5}{6}$
  - In $n \geq 4$, with prefix adjustment, $s\left(\alpha, \alpha\beta\right) = \frac{9}{10}$
  - $s\left(\beta, \alpha\beta\right) = 0$

- For edit-distance/*lcs*, $s\left(\alpha, \alpha\beta\right) = s\left(\beta, \alpha\beta\right) = \frac{3}{4}$

USCENSUSBUREAU

# Hybrid Comparator

- Compute both Jaro-Winkler and edit distance/*lcs*

# Hybrid Comparator

- Compute both Jaro-Winkler and edit distance/*lcs*

- Use larger of Jaro-Winkler and (scaled) edit distance/*lcs*

# Hybrid Comparator

- Compute both Jaro-Winkler and edit distance/*lcs*

- Use larger of Jaro-Winkler and (scaled) edit distance/*lcs*

- Where J-W does well, hybrid does a little better than either

# Hybrid Comparator

- Compute both Jaro-Winkler and edit distance/*lcs*

- Use larger of Jaro-Winkler and (scaled) edit distance/*lcs*

- Where J-W does well, hybrid does a little better than either

- Where J-W does significantly worse, hybrid does nearly as well as edit distance/*lcs*

USCENSUSBUREAU

# Hybrid Comparator, Cont.

- Can see some improvement in actual record linkage results

# Hybrid Comparator, Cont.

- Can see some improvement in actual record linkage results

- Calculation takes a long time

USCENSUSBUREAU

# String Comparator Summary

- String comparator improves record linkage

# String Comparator Summary

- String comparator improves record linkage

- String comparator takes significant amount of record linkage computation time

# String Comparator Summary

- String comparator improves record linkage

- String comparator takes significant amount of record linkage computation time

  - For J-W, about $30\%$

U S C E N S U S B U R E A U

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

- Does $U$ have any natural partitions?

U S C E N S U S B U R E A U

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

- Does $U$ have any natural partitions?

- For Census data

USCENSUSBUREAU

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

- Does $U$ have any natural partitions?

- For Census data
  - $U_1$, different people, same household

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

- Does $U$ have any natural partitions?

- For Census data
  - $U_1$, different people, same household
  - $U_2$, different people, different household

# More Than Two Latent Classes

- EM algorithm generalizes to more than 2 classes, $M, U$

- Does $U$ have any natural partitions?

- For Census data
  - $U_1$, different people, same household
  - $U_2$, different people, different household

- Classes have to be implicit in the matching data

USCENSUSBUREAU

# EM for Three Classes

- Use EM to estimate $\Pr(U_1)$, $\Pr(U_2)$, and marginal probabilities $\Pr(\gamma_i | U_1)$, $\Pr(\gamma_i | U_2)$

# EM for Three Classes

- Use EM to estimate $\Pr\left(U_1\right), \Pr\left(U_2\right)$, and marginal probabilities $\Pr\left(\gamma_i|U_1\right), \Pr\left(\gamma_i|U_2\right)$

- Recombine

$$\Pr\left(\gamma_i|U\right) = \frac{\Pr\left(\gamma_i|U_1\right)\Pr\left(U_1\right) + \Pr\left(\gamma_i|U_2\right)\Pr\left(U_2\right)}{\Pr\left(U_1\right) + \Pr\left(U_2\right)}$$

U S C E N S U S B U R E A U

# More Than Two Comparison Values

- Can have more than {agree, disagree}

# More Than Two Comparison Values

- Can have more than {agree, disagree}

- For $m$ comparison values, EM algorithm must estimate $2\,(m-1)$ parameters

US CENSUS BUREAU

# More Than Two Comparison Values

- Can have more than {agree, disagree}

- For $m$ comparison values, EM algorithm must estimate $2(m-1)$ parameters

- We have used {agree, disagree, missing} when data is often missing but has distinguishing power when present

USCENSUSBUREAU

# More Than Two Comparison Values

- Can have more than {agree, disagree}

- For $m$ comparison values, EM algorithm must estimate $2\,(m-1)$ parameters

- We have used {agree, disagree, missing} when data is often missing but has distinguishing power when present
  - For example, middle initial

# More Than Two Compr. Values, Cont.

- Reasonability check for parameter estimation

$$\log \frac{\Pr\left(\mathsf{blank}|M\right)}{\Pr\left(\mathsf{blank}|U\right)} \approx 0$$

USCENSUSBUREAU

# One-to-one Matching

- If both files have no duplication within them, then it is preferable to have output with each record linked to no more than one record in the other file

# One-to-one Matching

- If both files have no duplication within them, then it is preferable to have output with each record linked to no more than one record in the other file

- All records that are compared with each other are within a block

USCENSUSBUREAU

# One-to-one Matching

- If both files have no duplication within them, then it is preferable to have output with each record linked to no more than one record in the other file

- All records that are compared with each other are within a block

- Linear assignment algorithm used to find optimal one-to-one matches within a block

U S C E N S U S B U R E A U

# Linear Assignment Algorithm

- For agreement weights in block

|        | $B_1$    | $B_2$    | $B_3$    | $\cdots$ | $B_n$    |
|--------|----------|----------|----------|----------|----------|
| $A_1$  | $w_{11}$ | $w_{12}$ | $w_{13}$ |          | $w_{1n}$ |
| $A_2$  | $w_{21}$ | $w_{22}$ | $w_{23}$ |          | $w_{2n}$ |
| $A_3$  | $w_{31}$ | $w_{32}$ | $w_{33}$ |          | $w_{3n}$ |
| $\vdots$ |        |          |          |          |          |
| $A_n$  | $w_{n1}$ | $w_{n2}$ | $w_{n3}$ |          | $w_{nn}$ |

USCENSUSBUREAU

# Linear Assignment Algorithm

■ Find permutation $\bar{\sigma}$ that maximizes

$$\sum_{i=1}^{n} w_{i,\sigma(i)}$$

# Linear Assignment Algorithm

- Find permutation $\bar{\sigma}$ that maximizes

$$\sum_{i=1}^{n} w_{i,\sigma(i)}$$

- Not a greedy algorithm

# Linear Assignment Algorithm

- Find permutation $\bar{\sigma}$ that maximizes

$$\sum_{i=1}^{n} w_{i,\sigma(i)}$$

- Not a greedy algorithm

| | | | | |
|---|---|---|---|---|
| Father | 40 | $\leftrightarrow$ | Mother | 39 |
| Mother | 39 | $\leftrightarrow$ | Daughter | 16 |
| Daughter | 16 | $\leftrightarrow$ | Son | 13 |
| Son | 13 | | | |

US CENSUS BUREAU

# Error Rates

- **False Match Rate**

$$\mu = \Pr\left(L \mid U\right) = \Pr\left(w\left(\gamma\right) < T_\mu \mid U\right)$$

USCENSUSBUREAU

# Error Rates

- False Match Rate

$$\mu = \Pr\left(L \mid U\right) = \Pr\left(w\left(\gamma\right) < T_\mu \mid U\right)$$

- False Non-match Rate

$$\lambda = \Pr\left(N \mid M\right) = \Pr\left(w\left(\gamma\right) > T_\lambda \mid U\right)$$

USCENSUSBUREAU

# Practical Considerations

- Question:  Relative to what sample space?

# Practical Considerations

- Question: Relative to what sample space?
  - $A \times B$

# Practical Considerations

- Question: Relative to what sample space?
  - $A \times B$
  - Pairs in blocking scheme

# Practical Considerations

- Question: Relative to what sample space?
  - $A \times B$
  - Pairs in blocking scheme
  - After 1-1 matching

# Practical Considerations

- Question: Relative to what sample space?
  - $A \times B$
  - Pairs in blocking scheme
  - After 1-1 matching
- Each step presumably filters out a lot of low-weight pairs

# False Non-Match Rate

- Difficult to determine as well as define

# False Non-Match Rate

- Difficult to determine as well as define

- May as well try to estimate number of undiscovered matches in $A \times B$

# False Non-Match Rate

- Difficult to determine as well as define

- May as well try to estimate number of undiscovered matches in $A \times B$

- Can try capture-recapture using *independent* blocking schemes

# False Match Rate

- Bellin-Rubin

# False Match Rate

- Bellin-Rubin
- Larsen

# False Match Rate

- Bellin-Rubin
- Larsen
- Larsen, Rubin, Winkler

# Bellin-Rubin

- Bellin, T.R. and Rubin, D.B. (1995) "A Method for Calibrating False-Match Rates in Record Linkage," *Journal of the American Statistical Association*, 90,pp.694–707.

# Bellin-Rubin

- Bellin, T.R. and Rubin, D.B. (1995) "A Method for Calibrating False-Match Rates in Record Linkage," *Journal of the American Statistical Association*, 90,pp.694–707.

- Consider sample space without 1-1 matching

# Bellin-Rubin

- Bellin, T.R. and Rubin, D.B. (1995) "A Method for Calibrating False-Match Rates in Record Linkage," *Journal of the American Statistical Association*, 90,pp.694–707.

- Consider sample space without 1-1 matching

- Model as a mixture of 2 normal distributions (Box-Cox)

USCENSUSBUREAU

# Bellin-Rubin

- Bellin, T.R. and Rubin, D.B. (1995) "A Method for Calibrating False-Match Rates in Record Linkage," *Journal of the American Statistical Association*, 90,pp.694–707.

- Consider sample space without 1-1 matching

- Model as a mixture of 2 normal distributions (Box-Cox)

- $M$ and $U$ must be well-separated and unimodal

US CENSUS BUREAU

# Larsen

- Larsen, M.D. "Hierarchical Bayesian Record Linkage Theory," Iowa State University, Statiistics Department Technical Report

# Larsen

- Larsen, M.D. "Hierarchical Bayesian Record Linkage Theory," Iowa State University, Statiistics Department Technical Report

- Estimate error rates with 1-1 matching

# Larsen

- Larsen, M.D. "Hierarchical Bayesian Record Linkage Theory," Iowa State University, Statiistics Department Technical Report

- Estimate error rates with 1-1 matching

- Complicated restrained optimization

USCENSUSBUREAU

# Larsen

- Larsen, M.D. "Hierarchical Bayesian Record Linkage Theory," Iowa State University, Statiistics Department Technical Report

- Estimate error rates with 1-1 matching

- Complicated restrained optimization

- Metropolis-Hastings procedure

USCENSUSBUREAU

# Improved Parameter Estimates

- Recall, if we had correct parameter values (and model), under Fellegi-Sunter, error rates are known

USCENSUSBUREAU

# Improved Parameter Estimates

- Recall, if we had correct parameter values (and model), under Fellegi-Sunter, error rates are known

- Improve parameter estimates using training data

# Extended Likelihood Function

- For unlabled sample space $S$ and labeled training data set $T$, extended likelihood function

$$L = \left( \prod_{(a,b) \in S} \Pr\left(\gamma\left(a,b\right)\right) \right)^{1-\lambda} \left( \prod_{(a,b) \in T} \Pr\left(\gamma\left(a,b\right)\right) \right)^{\lambda}$$

for $0 \leq \lambda \leq 1$

# Extended Likelihood Function

- For unlabled sample space $S$ and labeled training data set $T$, extended likelihood function

$$L = \left( \prod_{(a,b) \in S} \Pr\left(\gamma\left(a,b\right)\right) \right)^{1-\lambda} \left( \prod_{(a,b) \in T} \Pr\left(\gamma\left(a,b\right)\right) \right)^{\lambda}$$

for $0 \leq \lambda \leq 1$

- Estimate using EM

USCENSUSBUREAU

# Larsen, Rubin

- Larsen, M.D. and Rubin, D.B. (2001) "Iterative Automated Record Linkage Using Mixture Models," *Journal of the American Statistical Association* 79, pp.32–41

# Larsen, Rubin

- Larsen, M.D. and Rubin, D.B. (2001) "Iterative Automated Record Linkage Using Mixture Models," *Journal of the American Statistical Association* 79, pp.32–41

- $T$ is sample of pairs from the clerical review region that have been clerically reviewed

# Winkler

- Winkler, W.E. "Automatically Estimating Record Linkage False Match Rates," (2007) http://www.census.gov/srd/www/byname.html

# Winkler

- Winkler, W.E. "Automatically Estimating Record Linkage False Match Rates," (2007) http://www.census.gov/srd/www/byname.html

- $T$ is sample of "pseudo-truth" data: pairs with sufficiently high or sufficiently low weight

USCENSUSBUREAU

# Data Preparation

- Files must have matching fields of fixed length and location

# Data Preparation

- Files must have matching fields of fixed length and location

- Matching fields are compared on a character by character basis

# Data Preparation

- Files must have matching fields of fixed length and location

- Matching fields are compared on a character by character basis

- Unnecessary inconsistencies must be removed before matching is done

U S C E N S U S B U R E A U

# Basic Preparation

- Consistently encode categorical variables

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race
  - Date, age

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race
  - Date, age
- Spelling standardization

USCENSUSBUREAU

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race
  - Date, age
- Spelling standardization
  - Titles: Dr, Dr., Doctor

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race
  - Date, age
- Spelling standardization
  - Titles: Dr, Dr., Doctor
  - Nicknames: Bill, William

USCENSUSBUREAU

# Basic Preparation

- Consistently encode categorical variables
  - Sex, race
  - Date, age
- Spelling standardization
  - Titles: Dr, Dr., Doctor
  - Nicknames: Bill, William
  - Standard words: Co, Co., Cmpny, Company

# Basic Preparation, Cont.

- Identify and parse components

# Basic Preparation, Cont.

- Identify and parse components
  - Names: last, first

# Basic Preparation, Cont.

- Identify and parse components
  - Names: last, first
  - Addresses: number, street, unit

USCENSUSBUREAU

# Address Parsing

- 16 W Main ST APT 16

  RR 2 BX 215

  Fuller BLDG SUITE 405

  14588 HWY 16 W

US CENSUS BUREAU

# Address Parsing

16 W Main ST APT 16

RR 2 BX 215

Fuller BLDG SUITE 405

14588 HWY 16 W

| Pre2 | Hsnm | Stnm | RR | Box | Post1 | Post2 | Unit1 | Unit2 | Bldg |
|------|-------|--------|----|-----|-------|-------|-------|-------|-------|
| W | 16 | Main | | | | | | 16 | |
| | | | 2 | 215 | | | | | |
| | | | | | | | | 405 | Fuller |
| | 14588 | HWY 16 | | | | W | | | |

US CENSUS BUREAU

# Business Lists

- Much harder

# Business Lists

- Much harder

- May have fewer comparison fields

# Business Lists

- Much harder

- May have fewer comparison fields
  - Name

USCENSUSBUREAU

# Business Lists

- Much harder

- May have fewer comparison fields
  - Name
  - Address

# Business Lists

- Much harder

- May have fewer comparison fields
  - Name
  - Address
  - Phone

# Business Lists

- Much harder

- May have fewer comparison fields
  - Name
  - Address
  - Phone

- These may not be unique

# Business Lists

- Much harder

- May have fewer comparison fields
  - Name
  - Address
  - Phone

- These may not be unique

- May be difficult to parse

# Example of Business Name Parsing

DR John J Smith MD

Smith DRY FRM

Smith & Son ENTP

USCENSUSBUREAU

# Example of Business Name Parsing

- DR John J Smith MD

  Smith DRY FRM

  Smith & Son ENTP

-

| Pre | First | Mid | Last | Post1 | Post2 | Bus1 | Bus2 |
|-----|-------|-----|------|-------|-------|------|------|
| DR | John | J | Smith | MD | | | |
| | | | Smith | | | DRY | FRM |
| | | | Smith | | Son | ENTP | |

USCENSUSBUREAU

# Two Kinds of Standardizer

- Deterministic

# Two Kinds of Standardizer

- Deterministic
  - Rule based

# Two Kinds of Standardizer

- Deterministic
  - Rule based
- Probabilistic

USCENSUSBUREAU

# Two Kinds of Standardizer

- Deterministic
  - Rule based

- Probabilistic
  - Hidden Markov model

# Rule-Based Standardizer

- U.S. Census Bureau software

# Rule-Based Standardizer

- U.S. Census Bureau software

- Based on extensive expert experience

USCENSUSBUREAU

# Rule-Based Standardizer

- U.S. Census Bureau software

- Based on extensive expert experience

- Created for a specific sample space

USCENSUSBUREAU

# Hidden Markov Standardizer

- Adaptable to different sample spaces

US CENSUS BUREAU

# Hidden Markov Standardizer

- Adaptable to different sample spaces

- Based on training data

# Hidden Markov Standardizer Referenc

- P. Christen, T. Churches, J.X. Jhu. (2002) "Probabilistic Name and Address Cleaning and Standardization." *The Australasian Data Mining Workshop*. http://datamining.anu.eedu.au/projects/linkage.html

# Hidden Markov Standardizer Referenc

- P. Christen, T. Churches, J.X. Jhu. (2002) "Probabilistic Name and Address Cleaning and Standardization." *The Australasian Data Mining Workshop*. http://datamining.anu.eedu.au/projects/linkage.html

- T. Churches, P. Christen, J. Lu, J.X. Zhu. (2002) "Preparation of Name and Address Data for Record Linkage Using Hidden Markov Models." *BioMed Central Medical Informatics and Decision Making,* 2(9), http://www.biomedcentral.com/1472-6947/2/9.

U S C E N S U S B U R E A U

# Hidden Markov Standardizer Referenc

- P. Christen, T. Churches, J.X. Jhu. (2002) "Probabilistic Name and Address Cleaning and Standardization." *The Australasian Data Mining Workshop*. http://datamining.anu.eedu.au/projects/linkage.html

- T. Churches, P. Christen, J. Lu, J.X. Zhu. (2002) "Preparation of Name and Address Data for Record Linkage Using Hidden Markov Models." *BioMed Central Medical Informatics and Decision Making,* 2(9), http://www.biomedcentral.com/1472-6947/2/9.

- FEBRL Project (Freely Extensible Biomedical Record Linkage)

# Hidden Markov Model

- Identify a finite number of hidden Markov states

U S C E N S U S B U R E A U

# Hidden Markov Model

- Identify a finite number of hidden Markov states
  - first, last1, last2, mi, prefix, suffix

USCENSUSBUREAU

# Hidden Markov Model

- Identify a finite number of hidden Markov states

  - first, last1, last2, mi, prefix, suffix

- Use training data to assign transition probabilities from one state to the next

# Hidden Markov Model

- Identify a finite number of hidden Markov states

  - first, last1, last2, mi, prefix, suffix

- Use training data to assign transition probabilities from one state to the next

- Use training data to assign probabilities for observations having given hidden state

# Hidden Markov Model

- Identify a finite number of hidden Markov states
  - first, last1, last2, mi, prefix, suffix
- Use training data to assign transition probabilities from one state to the next
- Use training data to assign probabilities for observations having given hidden state
  - Look-up lists

# Hidden Markov Model

- Identify a finite number of hidden Markov states
  - first, last1, last2, mi, prefix, suffix
- Use training data to assign transition probabilities from one state to the next
- Use training data to assign probabilities for observations having given hidden state
  - Look-up lists
  - Coded rules

US CENSUS BUREAU

# Hidden Markov Model, Cont.

- Break object into component observations, assign them initial Markov states

# Hidden Markov Model, Cont.

- Break object into component observations, assign them initial Markov states
  - "sir", "mick", "jagger", "mbe"

# Hidden Markov Model, Cont.

- Break object into component observations, assign them initial Markov states
  - "sir", "mick", "jagger", "mbe"
- Compute the highest probability sequence of hidden states for the given observations

# Viterbi Algorithm

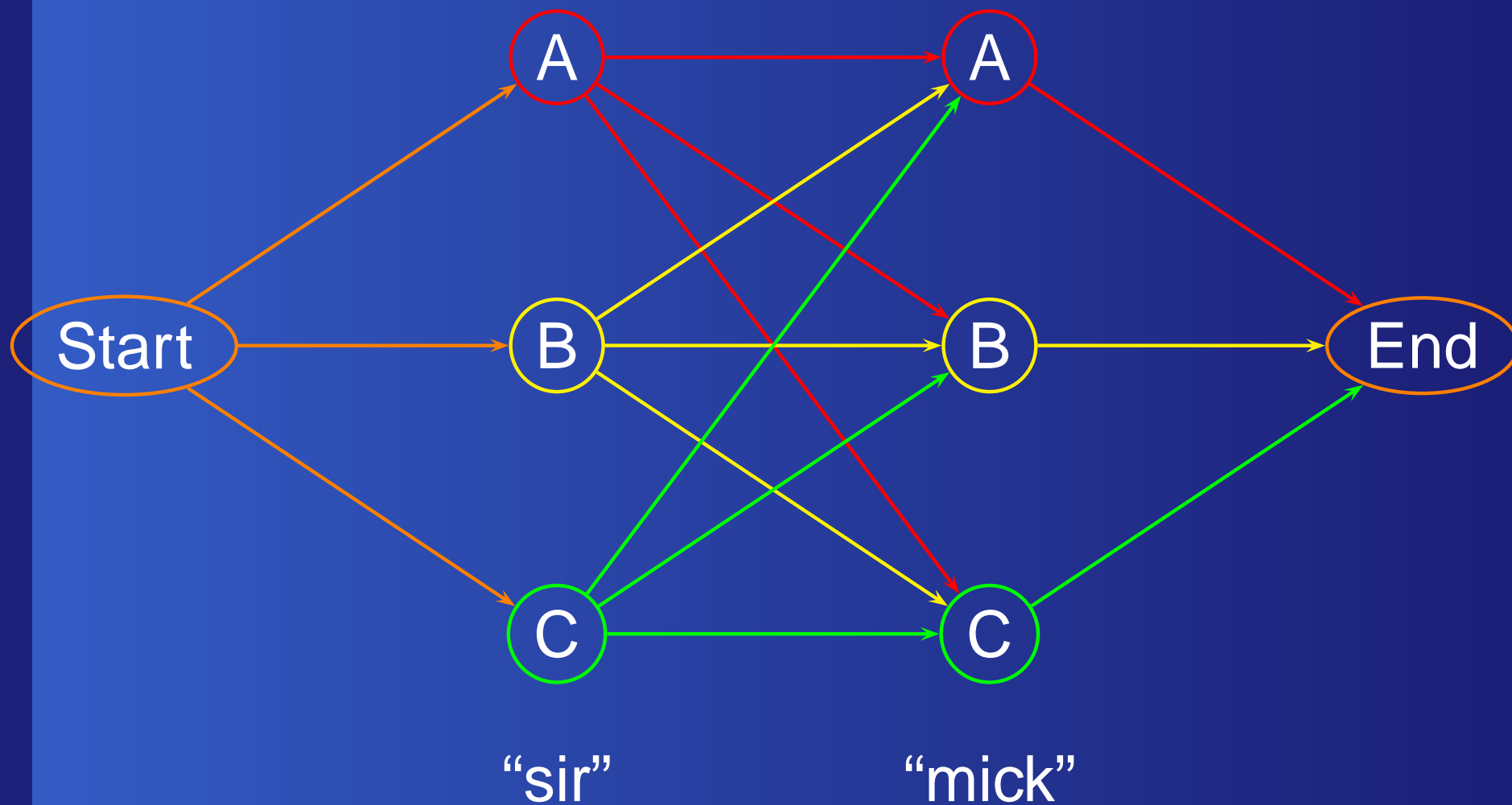- Not feasible to compute probabilities for all possible paths $O\left(n^{l}\right)$

# Viterbi Algorithm

- Not feasible to compute probabilities for all possible paths $O\left(n^l\right)$

- Dynamic programming algorithm $O\left(nl\right)$

# Viterbi Algorithm

- Not feasible to compute probabilities for all possible paths $O\left(n^l\right)$

- Dynamic programming algorithm $O\left(nl\right)$

- Each state is arrived at by the most probable subpath (Markov property)

USCENSUSBUREAU

# HMM Diagram



"sir"   "mick"

# Standardization Summary

- Much more time is likely to be spent preparing the data than performing the record linkage

# Standardization Summary

- Much more time is likely to be spent preparing the data than performing the record linkage

- Records that fail to be standardized will probably fail to be matched

US CENSUS BUREAU

# U.S. Census Bureau Software

- 🔴 Matching programs

# U.S. Census Bureau Software

- Matching programs
  - Matcher

U S C E N S U S B U R E A U

# U.S. Census Bureau Software

- Matching programs
  - Matcher
  - Bigmatch

USCENSUSBUREAU

# U.S. Census Bureau Software

- 🔴 Matching programs
  - 🟢 Matcher
  - 🟢 Bigmatch
- 🔴 Auxiliary programs

U S C E N S U S B U R E A U

# U.S. Census Bureau Software

- Matching programs
  - Matcher
  - Bigmatch
- Auxiliary programs
  - Counter

USCENSUSBUREAU

# U.S. Census Bureau Software

- Matching programs
  - Matcher
  - Bigmatch
- Auxiliary programs
  - Counter
  - EM

USCENSUSBUREAU

# U.S. Census Bureau Software

- Matching programs
  - Matcher
  - Bigmatch
- Auxiliary programs
  - Counter
  - EM
  - Standardizer

USCENSUSBUREAU

# Matching Programs: Matcher

- Matcher

# Matching Programs: Matcher

- Matcher
    - One-to-one matching

# Matching Programs: Matcher

- Matcher
  - One-to-one matching
    - Files should not have duplicates

# Matching Programs: Matcher

- Matcher
  - One-to-one matching
    - Files should not have duplicates
  - Pre-sort files according to blocking scheme

# Matching Programs: Matcher

- Matcher
  - One-to-one matching
    - Files should not have duplicates
  - Pre-sort files according to blocking scheme
  - Can re-run program on residual files

# Matching Programs: Matcher

- Matcher
  - One-to-one matching
    - Files should not have duplicates
  - Pre-sort files according to blocking scheme
  - Can re-run program on residual files
    - Resort files according to new blocking scheme

# Matching Programs: Bigmatch

- Bigmatch

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching
    - Can be used for deduplicating file

US CENSUS BUREAU

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching
    - Can be used for deduplicating file
  - Do not pre-sort files

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching
    - Can be used for deduplicating file
  - Do not pre-sort files
  - Can run several blocking schemes

## U S C E N S U S B U R E A U

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching
    - Can be used for deduplicating file
  - Do not pre-sort files
  - Can run several blocking schemes
  - Can match several files to one file

U S C E N S U S B U R E A U

# Matching Programs: Bigmatch

- Bigmatch
  - No one-to-one matching
    - Can be used for deduplicating file
  - Do not pre-sort files
  - Can run several blocking schemes
  - Can match several files to one file
  - One file must fit into memory

USCENSUSBUREAU

# Auxiliary Programs: Counter

- Counter program

US CENSUS BUREAU

# Auxiliary Programs: Counter

- Counter program
  - Simplified matching program

# Auxiliary Programs: Counter

- Counter program
  - Simplified matching program
  - Counts number of times each matching pattern occurs

# Auxiliary Programs: Counter

- Counter program
  - Simplified matching program
  - Counts number of times each matching pattern occurs
  - String comparator has (high) cutoff

USCENSUSBUREAU

# Auxiliary Programs: Counter

- Counter program
  - Simplified matching program
  - Counts number of times each matching pattern occurs
  - String comparator has (high) cutoff
  - Provides input for EM

USCENSUSBUREAU

# Auxiliary Programs: EM

- EM algorithm program

# Auxiliary Programs: EM

- EM algorithm program
  - Estimates probability parameters for given file and blocking scheme

# Auxiliary Programs: EM

- EM algorithm program
  - Estimates probability parameters for given file and blocking scheme
  - Has 2-class and 3-class versions

U S C E N S U S B U R E A U

# Auxiliary Programs, Standardizer

- Standardizer

# Auxiliary Programs, Standardizer

- Standardizer
  - Standardizes names and addresses

USCENSUSBUREAU

# Auxiliary Programs, Standardizer

- Standardizer
  - Standardizes names and addresses
  - Rule-based parsing

USCENSUSBUREAU