Introduction: Why Deduplication? Labeled Records and Classification Manual record identification for 824 inventors² Which correspond to the **same unique person**? \implies 98,762 labeled USPTO inventor records Year First Middle Assignee Last St City (1) **Train classification models** 2004 Stanford University Millar David Stanford CA Fair Haven NJ 1995 Miller David UNC Miller David A.B. CA Stanfrod University 2002 Stanford (2) Calculate pairwise matching probabilities CA Lucent Technologies Inc. 1996 Miller David Andrew Stanford Fair Haven NJ Lucent Technologes, Inc. 1996 Miller David Andrew for all record-pairs ($p_{ij} > 0.5 \implies$ Match) Los Angeles CA Agilent Technologies, Inc. 2000 Miller David 2 Akinsanmi et al (2013) MA Lucent Technologies, Inc. 2002 Miller David D. Billerca Source: United States Patent & Trademark Office (USPTO) **Conditioning and Classifier Aggregation** Which **addresses** are the same? **Different types of record-pairs** different matching characteristics Address different classifiers \Longrightarrow 123 East Main Street, Pittsburgh, Pennsylvania (1) Condition on a feature of the record-pairs 0123 E. Main St., Pgh., Pa. 123 Main St., East Pittsburgh, PA (2) Train a classifier on each conditional subset 123 East Street, Portland, MA (3) Use appropriate classifier for prediction 123 East Street, Portland, WA **Approach to Deduplication Forest of Random Forests** "Too much" training data Goal: **Determine unique entities** in set of records $\implies \binom{98,782}{2} > 1$ billion labeled record-pairs (1) Compare pairs of text records (2) Determine pairwise matching probability (1) Randomly partition the training data (3) Identify clusters of duplicated entities (2) Train classifiers on each subset (4) Consolidate records within clusters Notation: Record, Field, Similarity, Distance Combine this approach with conditioning when any conditional subset is "too large" x_i : the i^{th} record, i = 1, ..., n x_{im} : the m^{th} field of x_i , m = 1, ..., M**Clustering to Resolve Transitivity Violations** γ_{ijm} : the similarity of $\mathbf{x}_{im}, \mathbf{x}_{jm}$ **Comparing Pairs of Text Records**

Long Strings (e.g. last name): Jaro-Winkler¹

 $\gamma_{\textit{ijm}} \in [\mathbf{0}, \mathbf{1}]$ $\gamma_{\textit{ijm}} = \mathbf{1}$ if $\mathbf{x}_{\textit{im}} = \mathbf{x}_{\textit{jm}}$

Short Strings (e.g. state code): Exact matching

$$\gamma_{ijm} = \mathbf{1}$$
 if $\mathbf{x}_{im} = \mathbf{x}_{jm}$
 $\gamma_{ijm} = \mathbf{0}$ if $\mathbf{x}_{im} \neq \mathbf{x}_{jm}$

Lists (e.g. co-authors): Intersection / Union

$$\gamma_{ijm} = rac{|\mathbf{x}_{im} \cap \mathbf{x}_{jm}|}{|\mathbf{x}_{im} \cup \mathbf{x}_{jm}|}$$

1 Winkler, W.E. (1990)

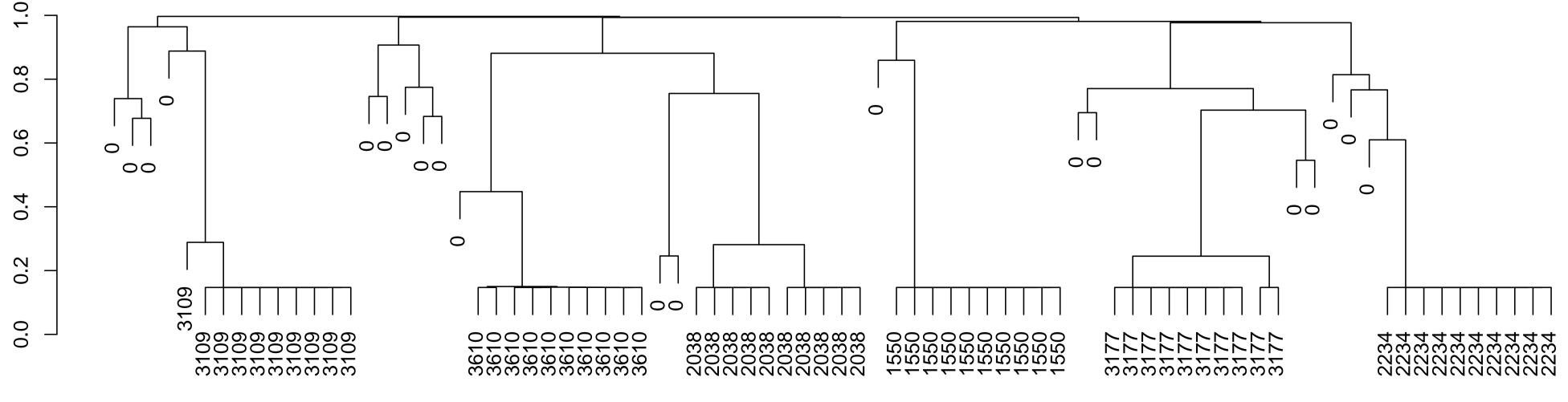
A Supervised Learning Approach for Identifying Duplicated Text Records

- on pairwise comparisons of labeled records

- (3) Aggregate (average) predictions of classifiers

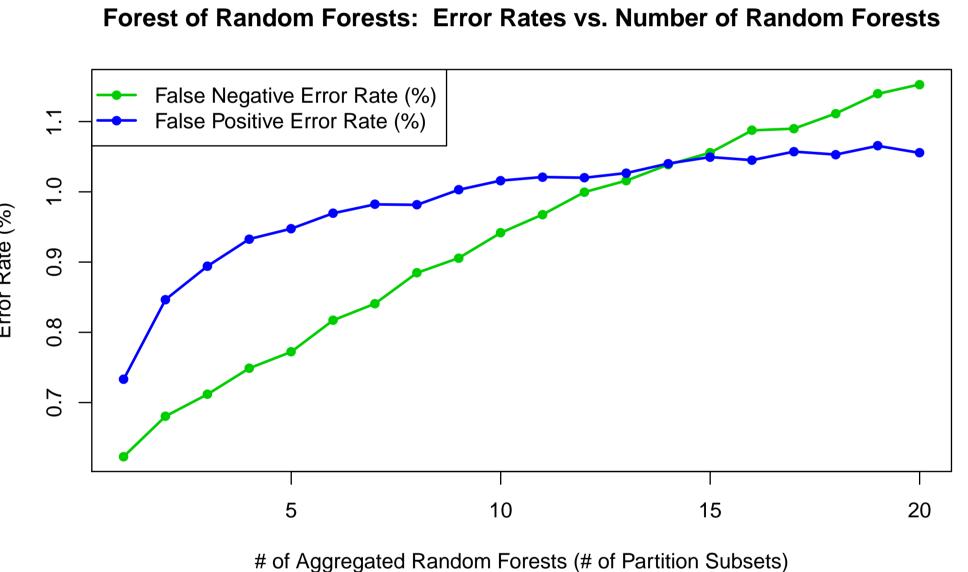
Cluster with **D**, the distance matrix, where $D[i, j] = d_{ij} = h(p_{ij}) =$ the distance between x_i, x_j

Deduplication Dendrogram: 6 Labeled Individuals, Distance = 1-p



Lai et Rando

Cond



6 Unique Individuals with IDs (0 = not identified) hclust (*, "complete")

Dedu Lai et Classific

Logistic Rando

Random Forests: low, balanced error rates

3 Lai et al (2009): weighted sums of similarity scores

Carnegie Mellon University Department of Statistics

uplication	False	False
1ethod	Negatives (%)	Positives (%)
al (2009) ³	8.39	4.13
cation Trees	2.23	2.49
Regression	1.68	1.64
om Forests	0.62	0.74

Deduplication	False	False
Method	Negatives (%)	Positives (%)
Lai et al (2009)	8.39	4.13
Random Forests	0.62	0.74
Conditional RF	0.61	0.64

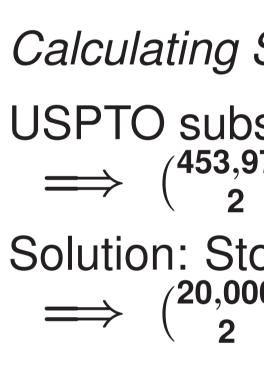
Conditioning further reduces error rates

Issues & Future Work

Reducing Deduplication Run Time



All $\binom{n}{2}$ comparisons: $O(n^2)$ run time (left) Solution: Blocking to limit # of comparisons (right) Current/Future Work: Use labeled training data to find blocking schemes which reduce run time without increasing deduplication error



Current/Future Work: Automatically determine the best conditioning scheme given input data

Extending to US Census Contexts

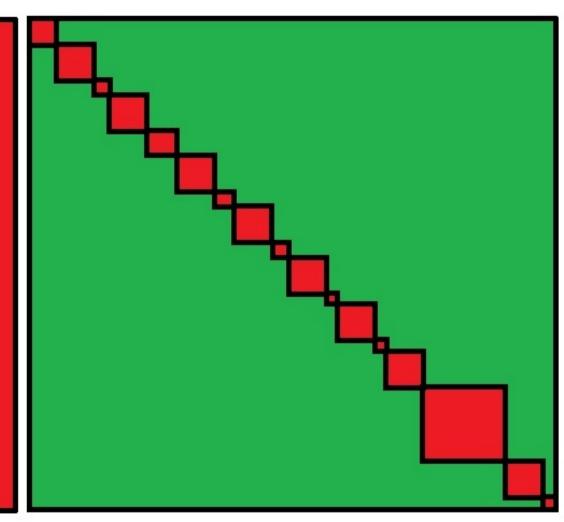
Acknowledgments

NSF SES-1130706 – NCRN-MN: "Data Integration, Online Data Collection and Privacy Protection for Census 2020" NSF RTG Grant 25951-1-1121631 NSF Science of Science and Innovation Policy Grant, "Quantifying the Resilience of Industry Innovation Ecosystems" (Award 0830354) NSF Science of Science and Innovation Policy Grant, "CAREER: Rethinking National Innovation Systems – Economic Downturns, Offshoring, and the Global Evolution of Technology" (Award 1056955)

Corresponding Author: Samuel Ventura. Mail: sventura@stat.cmu.edu

Samuel L. Ventura and Rebecca Nugent





Calculating Similarity Scores

USPTO subset: 453,972 inventor records $\implies \binom{453,972}{2} \approx 100 \text{ billion calculations / field}$

Solution: Store / re-use repeated comparisons $\implies \binom{20,000}{2} \approx 200 \text{ million calculations / field}$

Automate Choice of Conditioning Scheme

(1) Labeled Census records SSNs are susceptible to error SSN "labels" for classification?

(2) **Computational feasibility for Census data** 2010 Census: 300 million records Scale classification/clustering approach?