



Statistical Methods for the Forensic Analysis of Geolocated Event Data

By: **Christopher Galbraith** (Department of Statistics, University of California, Irvine), Padhraic Smyth (Department of Computer Science, University of California, Irvine), and Hal S. Stern (Department of Statistics, University of California, Irvine)

From the proceedings of
The Digital Forensic Research Conference
DFRWS 2020 USA
Virtual -- July 20-24

DFRWS is dedicated to the sharing of knowledge and ideas about digital forensics research. Ever since it organized the first open workshop devoted to digital forensics in 2001, DFRWS continues to bring academics and practitioners together in an informal environment.

As a non-profit, volunteer organization, DFRWS sponsors technical working groups, annual conferences and challenges to help drive the direction of research and development.

<https://dfrws.org>

Statistical Methods for the Forensic Analysis of Geolocated Event Data

Christopher Galbraith

Padhraic Smyth

Hal S. Stern

DFRWS USA

7.22.20



UCIRVINE
UNIVERSITY of CALIFORNIA • IRVINE



The material presented here is based upon work supported by the National Institute of Science and Technology under Award No. 70NANB15H176. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Institute of Science and Technology, nor of the Center for Statistics and Applications in Forensic Evidence.

Outline

1 Motivation

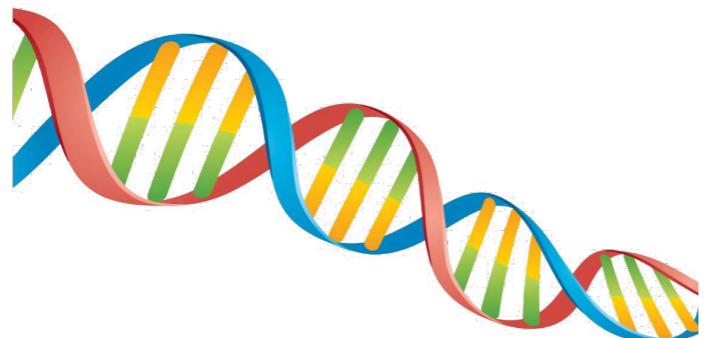
2 Quantifying Strength of Evidence

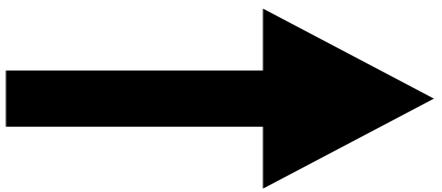
3 Application to Geolocated Event Data

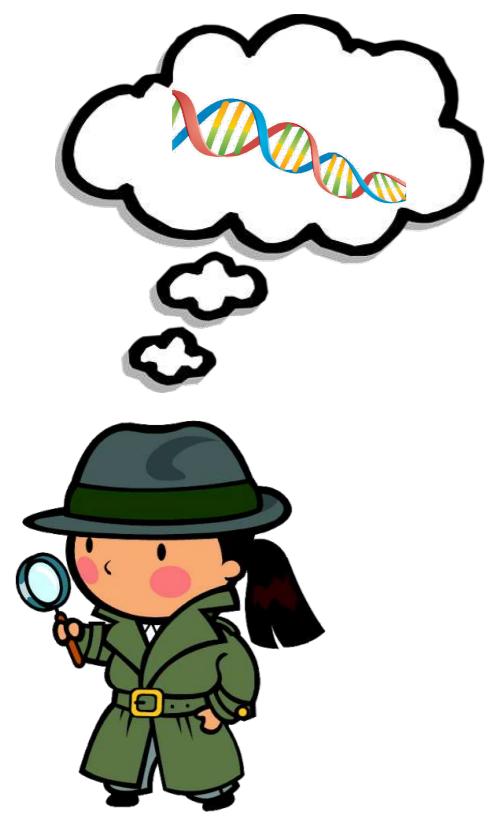
4 Future Directions and Conclusions

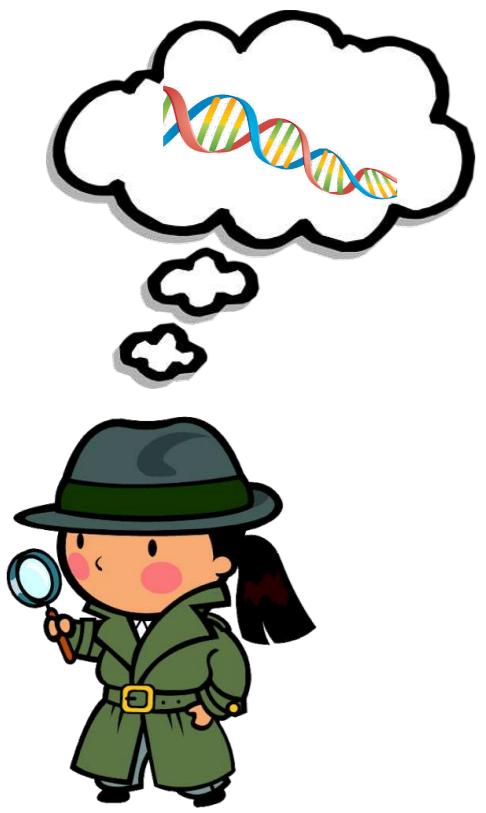
Motivation



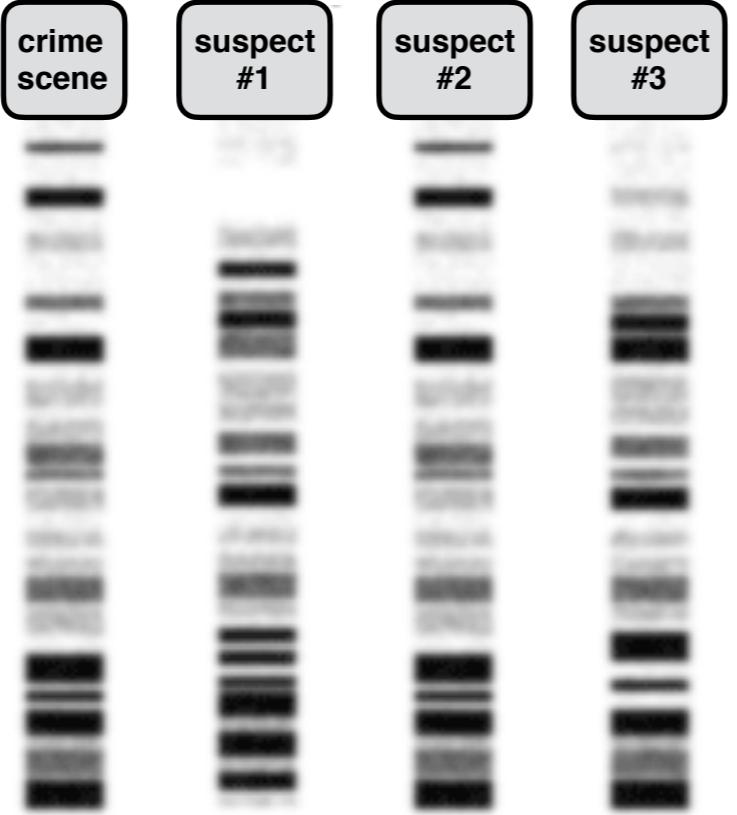


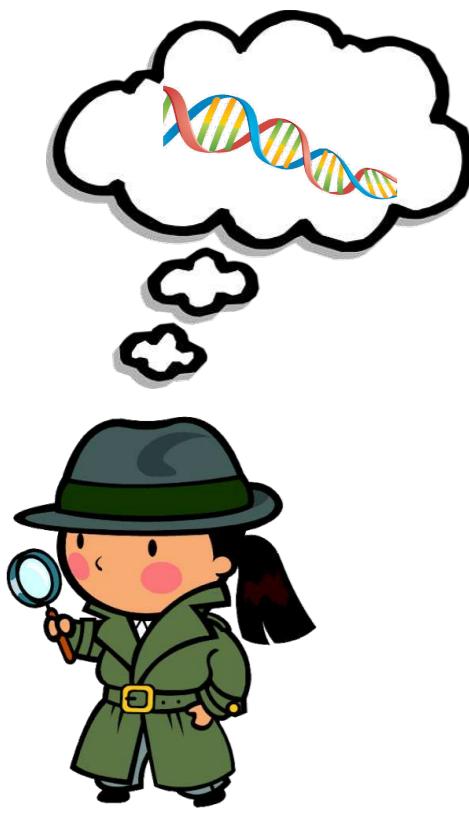




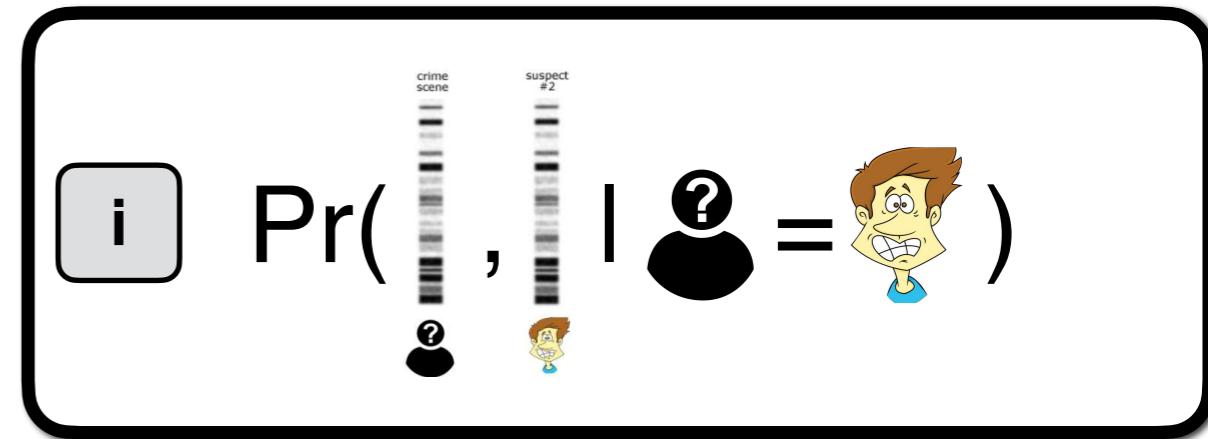
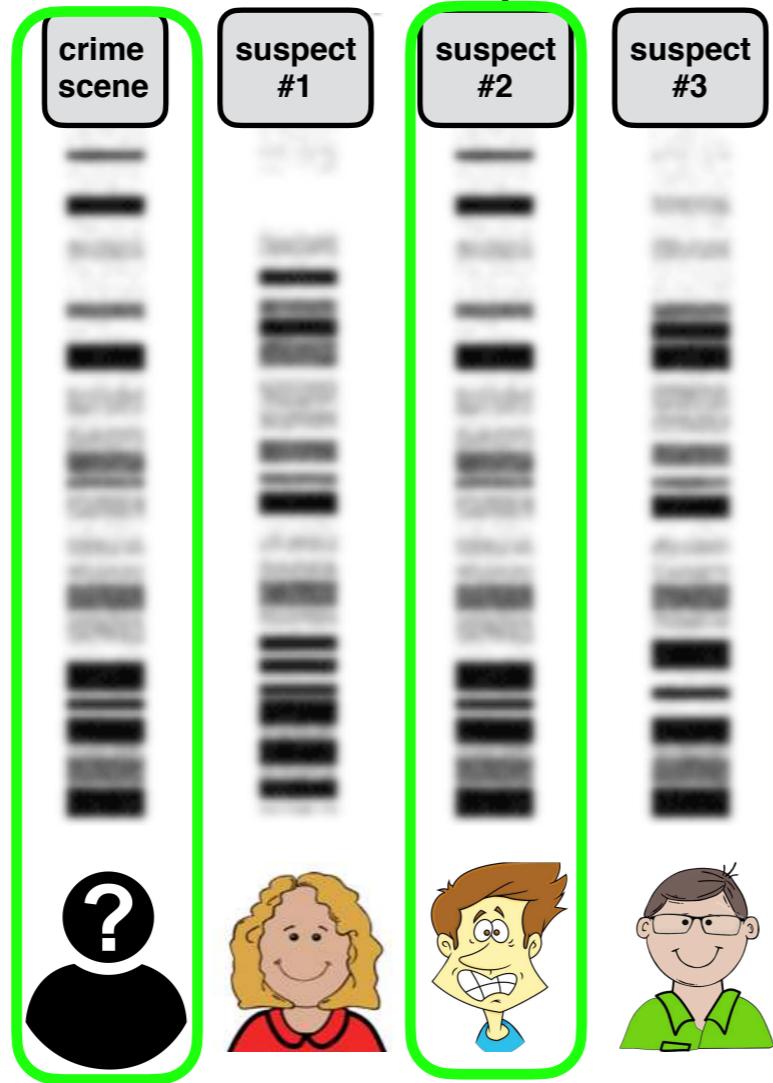


DNA Samples



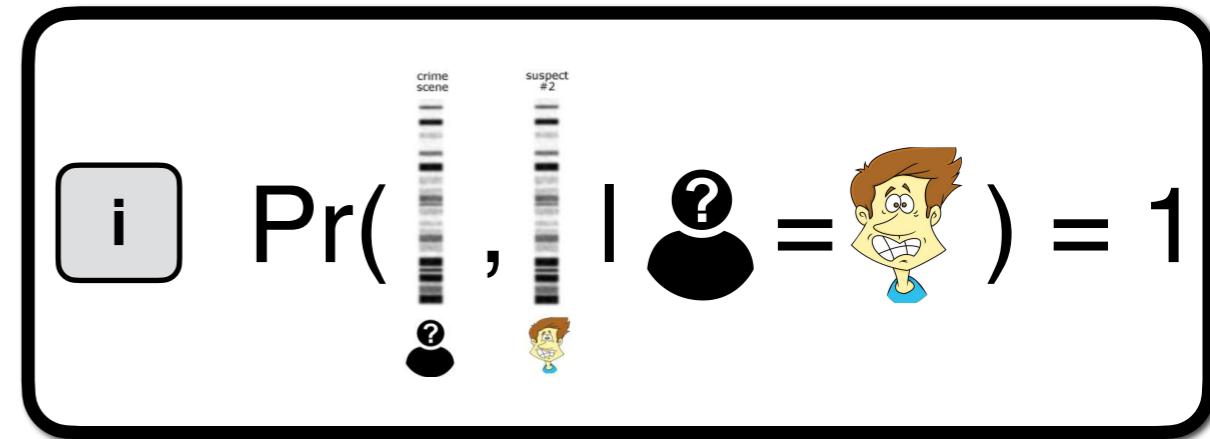
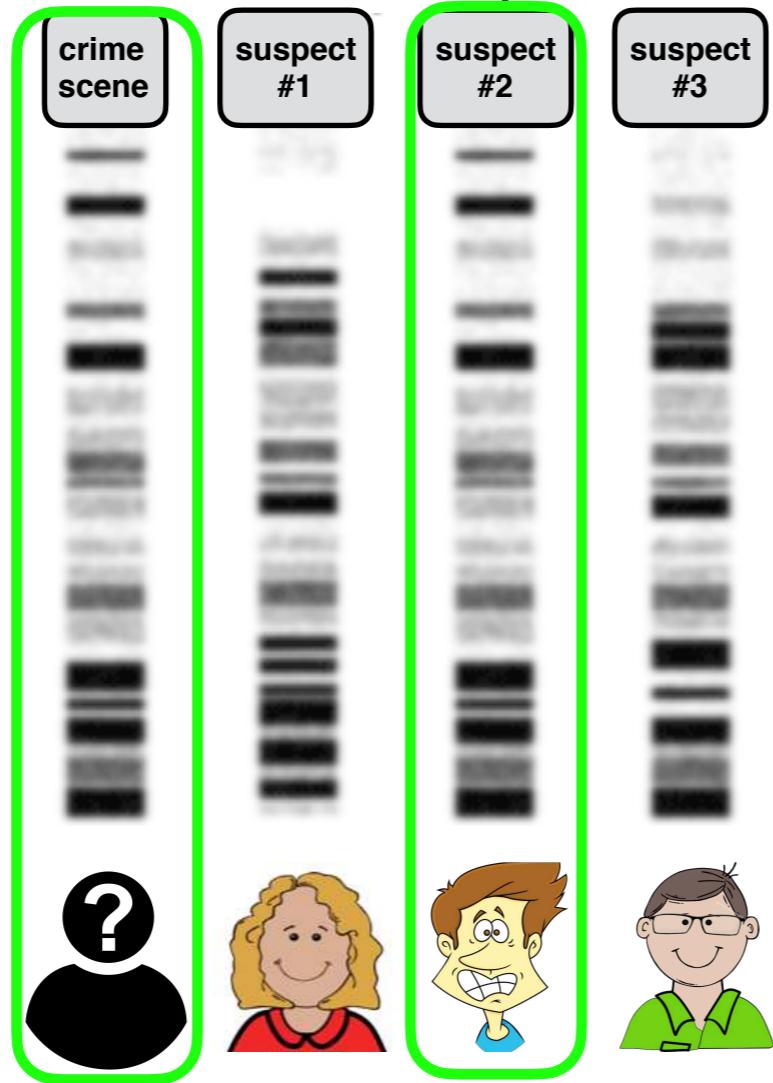


DNA Samples



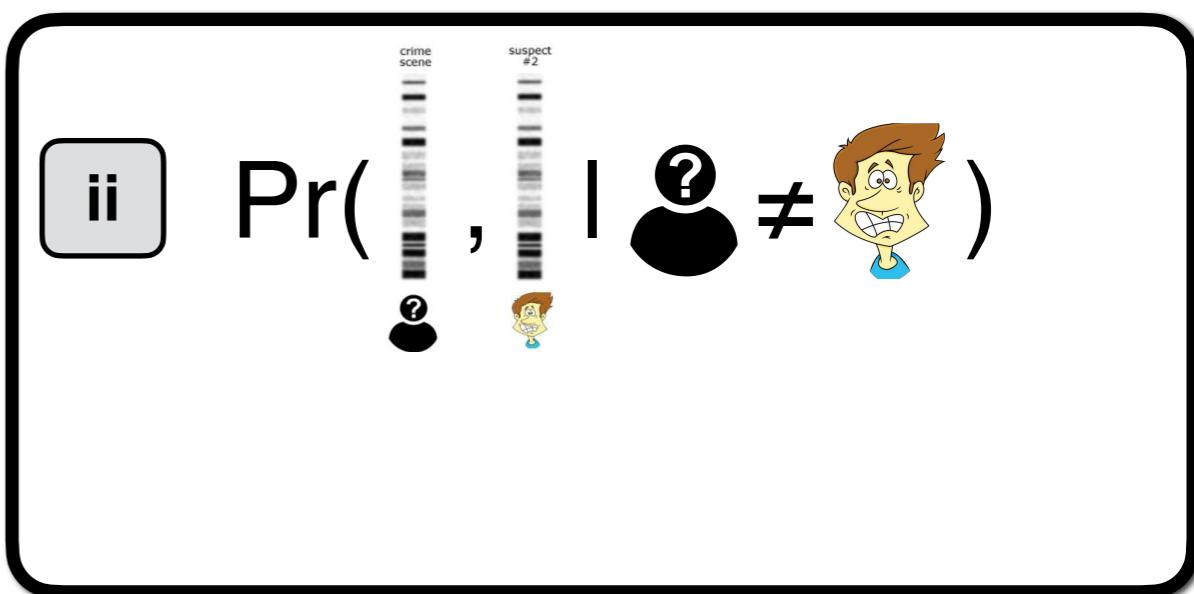
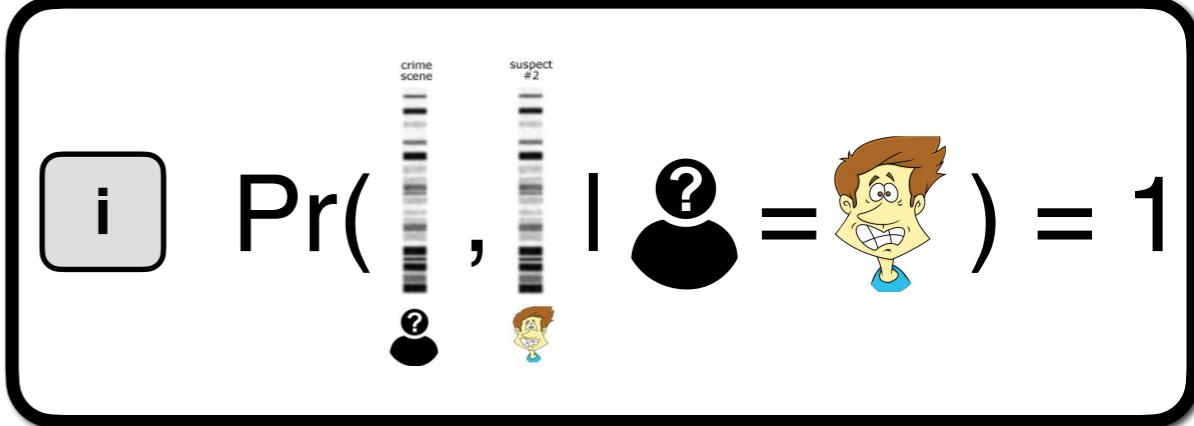
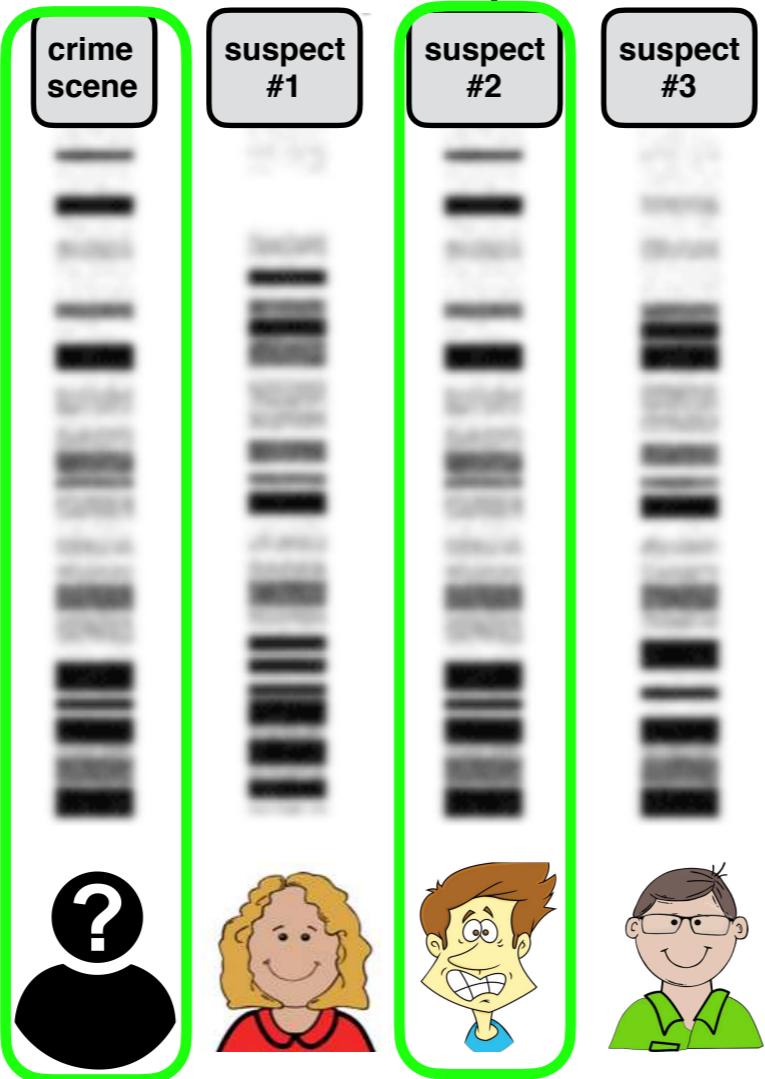


DNA Samples



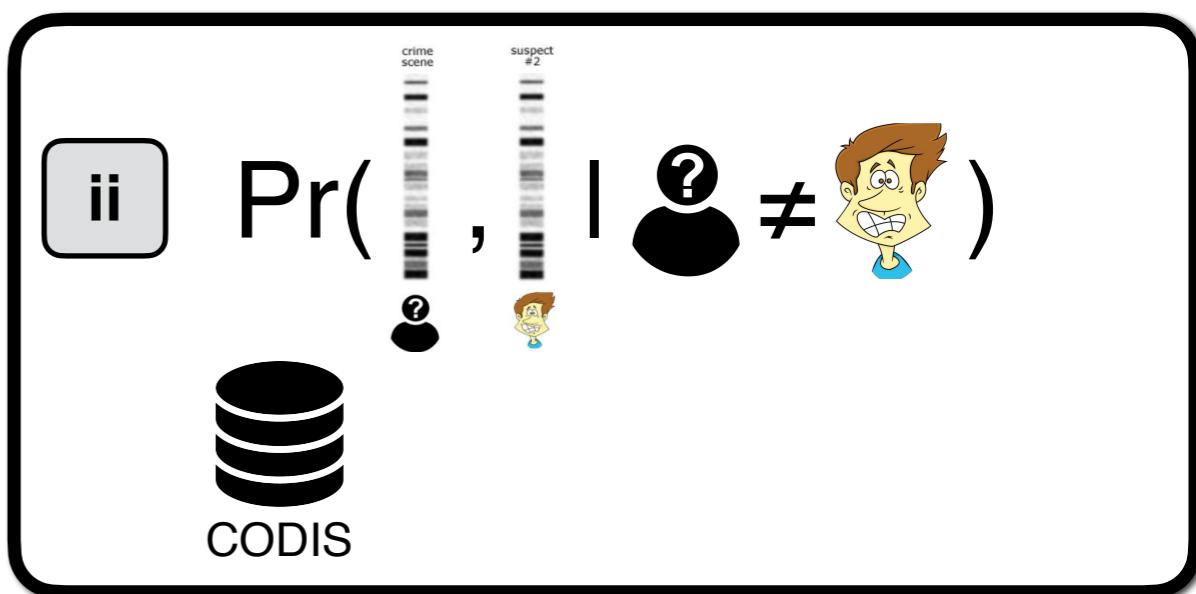
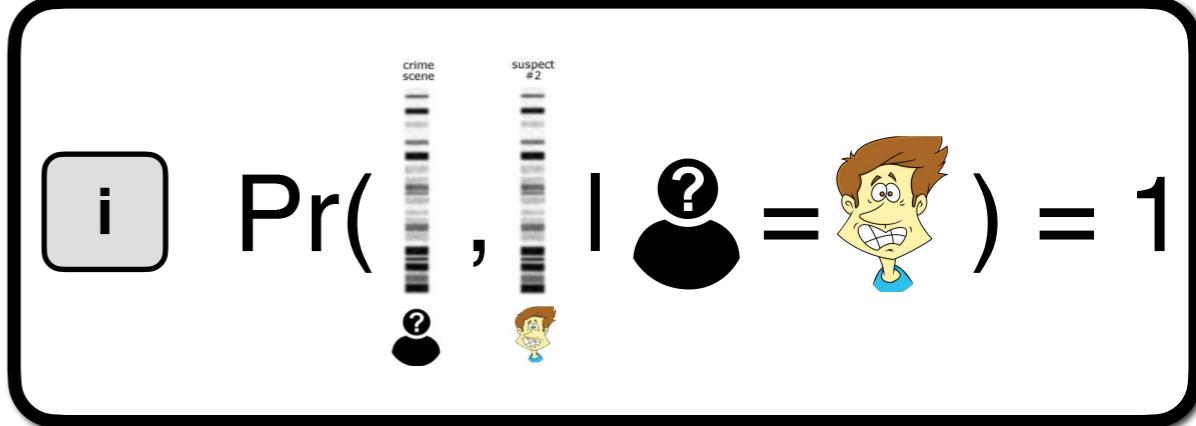
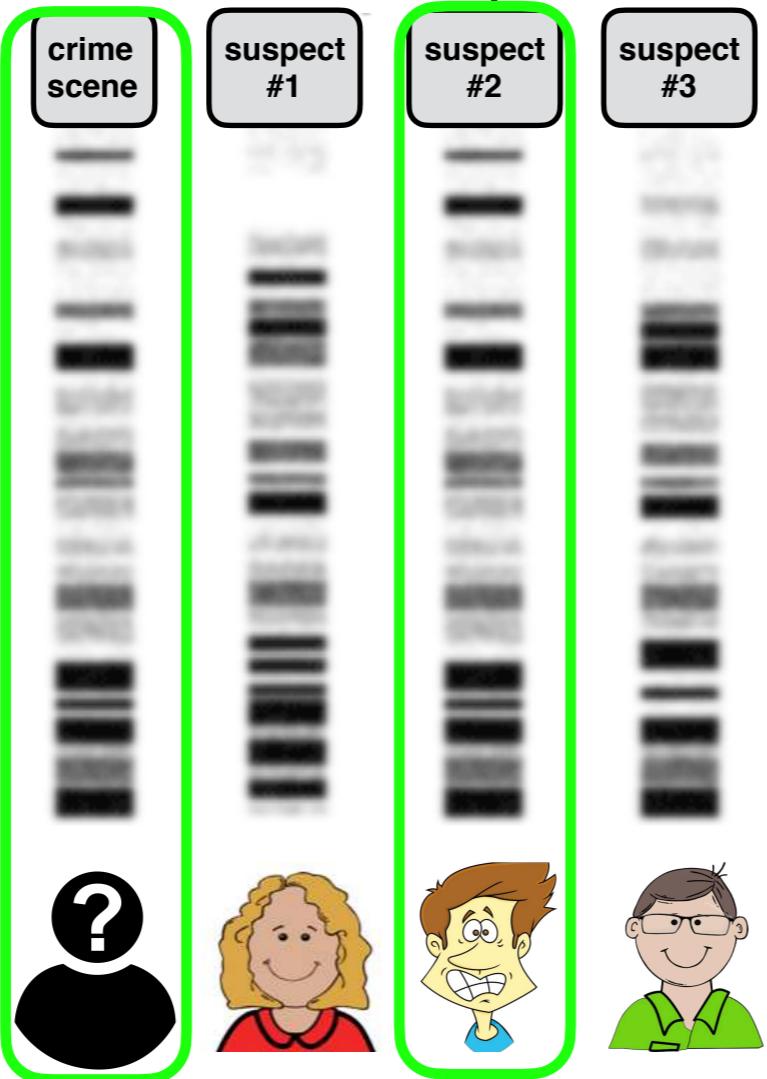


DNA Samples



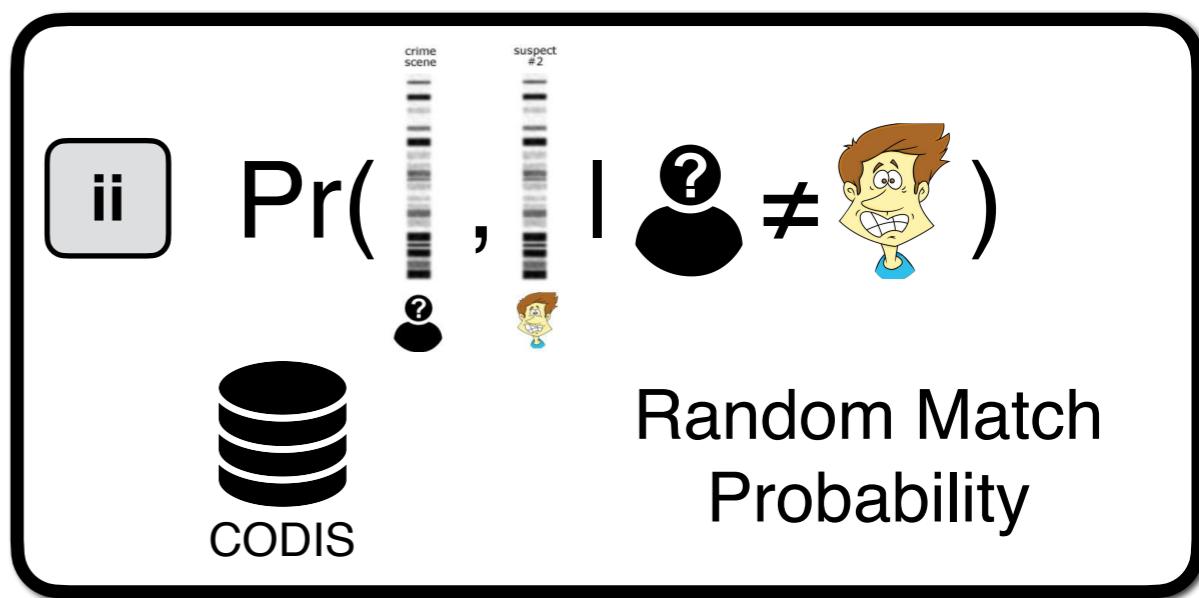
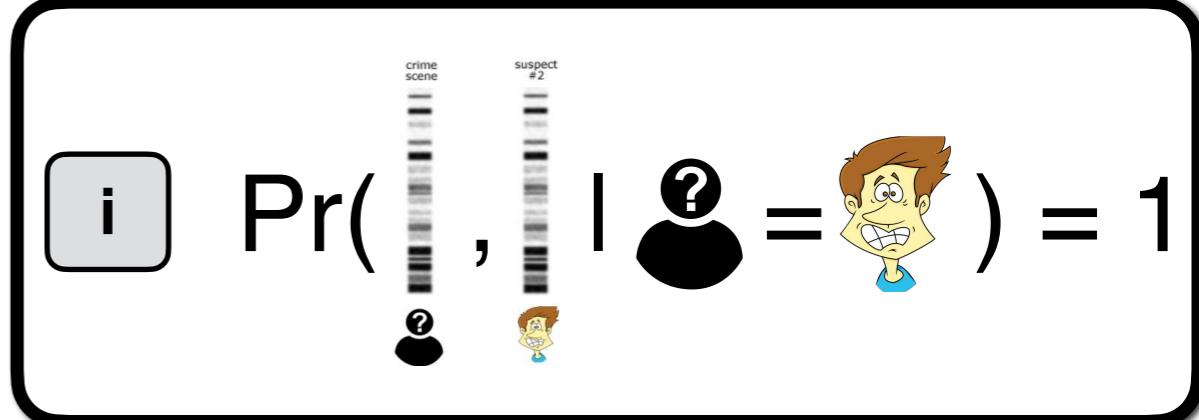
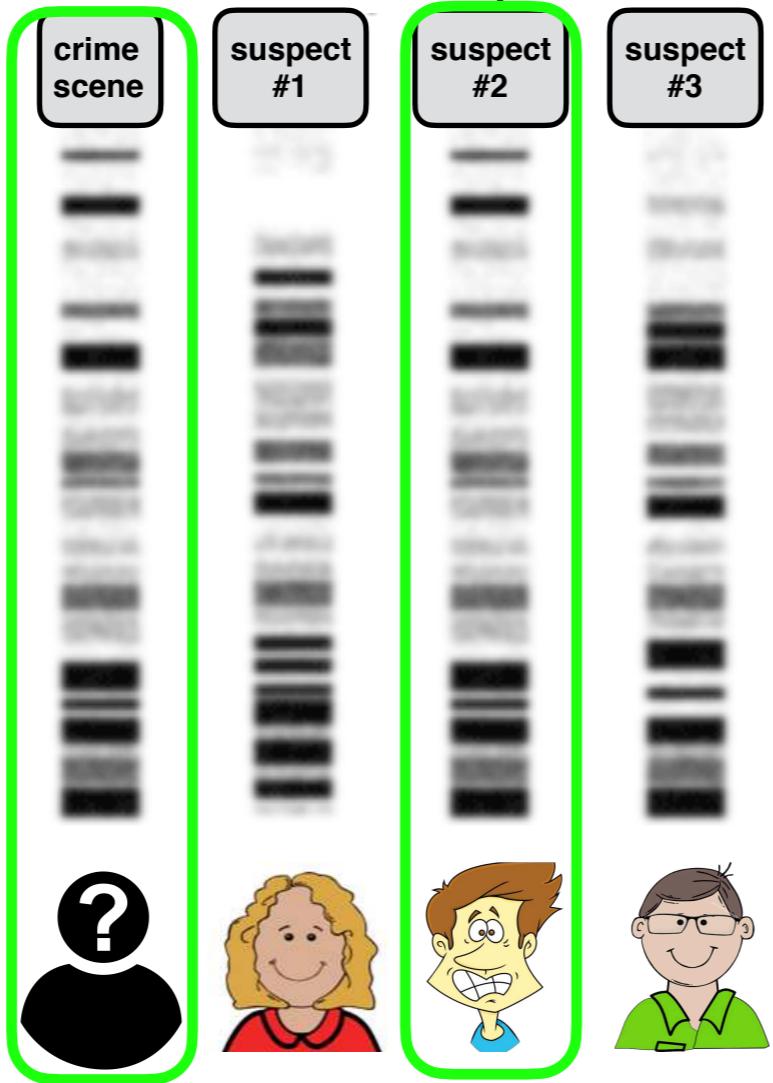


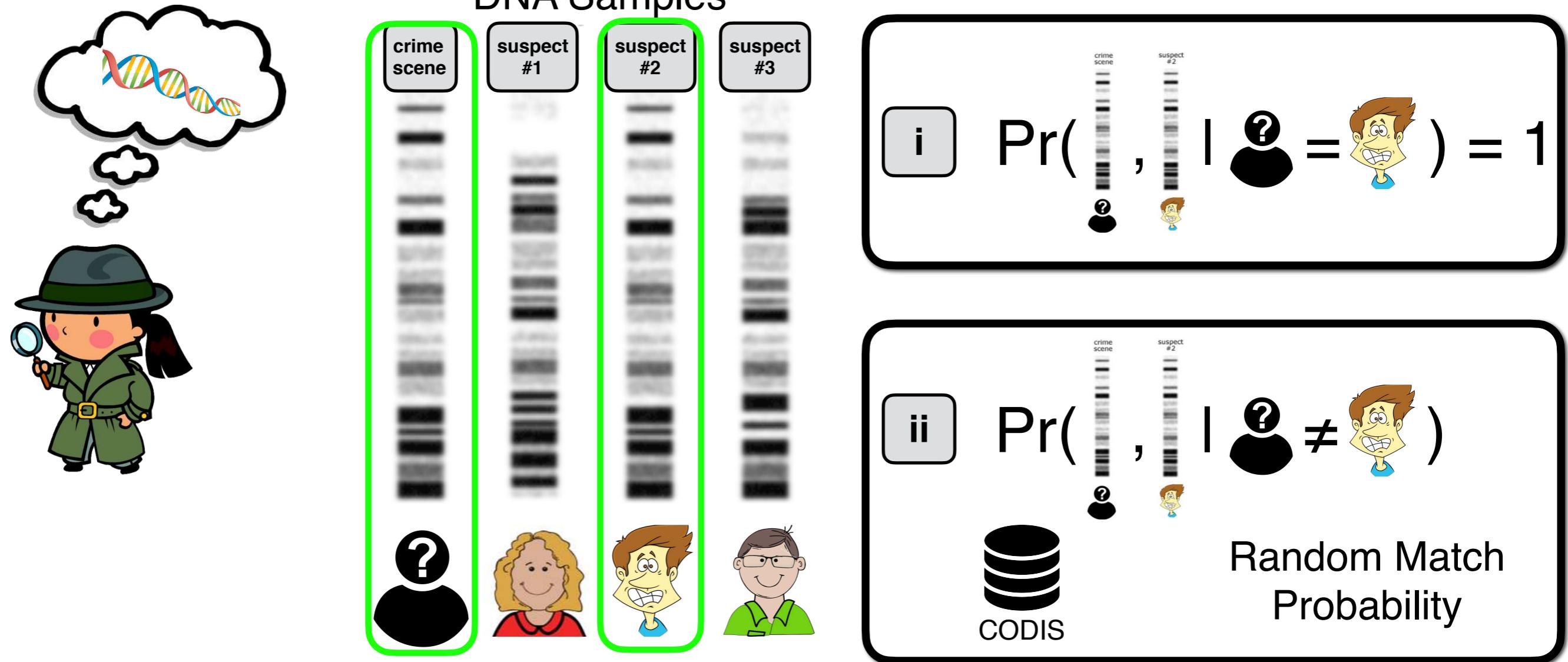
DNA Samples





DNA Samples





Likelihood Ratio

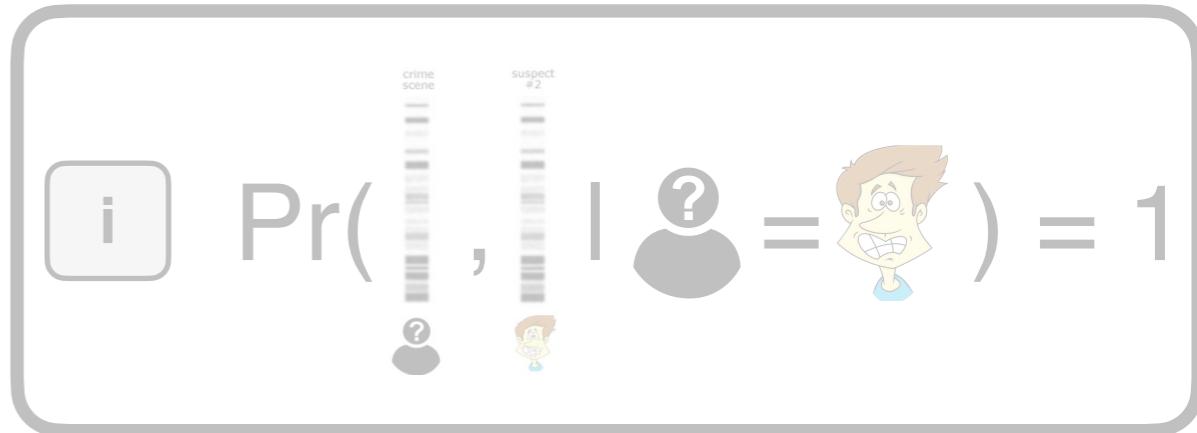
$$\frac{i}{ii}$$

< 1 Samples from different sources

$= 1$ Inconclusive

> 1 Samples from same source

DNA Samples



ch

sources

source







Extraction

[SWDGE, 2019;
Roussev, 2016;
Casey, 2011]

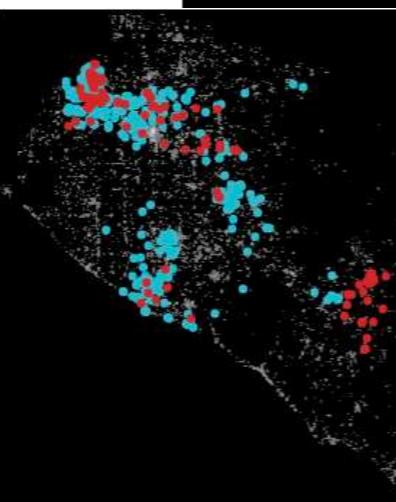
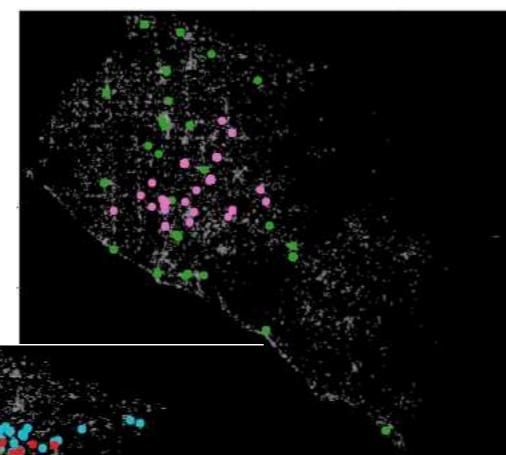
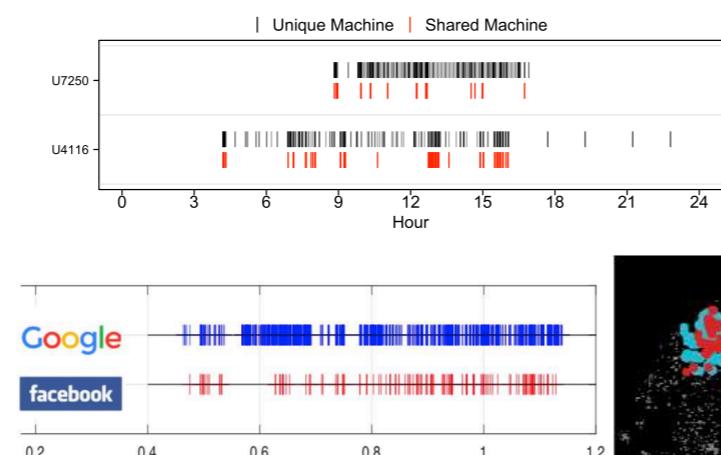
Browser requests
Web searches
Email activity
Phone/SMS
Social media activity
GPS locations
File access
Network activity
Exercise/movement
...



Extraction

[SWDGE, 2019;
Roussev, 2016;
Casey, 2011]

Browser requests
Web searches
Email activity
Phone/SMS
Social media activity
GPS locations
File access
Network activity
Exercise/movement
...



Analysis & Visualization

[Buchholz and Falk, 2005;
Grier, 2011;
Koven et al., 2016;
Gresty et al., 2016;
Kirchler et al., 2016]

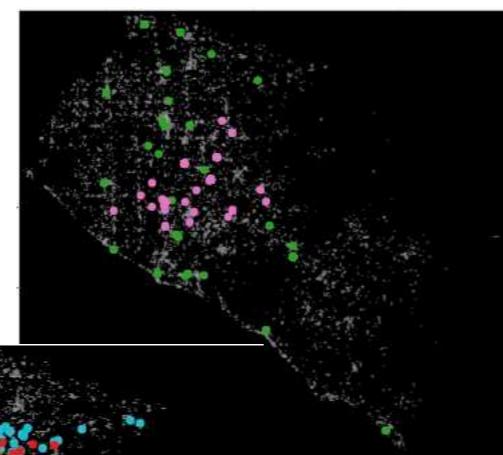
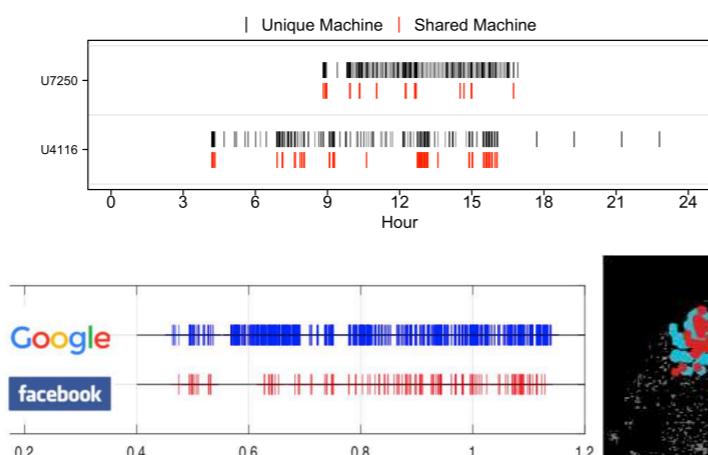


Extraction

[SWDGE, 2019;
Roussev, 2016;
Casey, 2011]

Browser requests
Web searches
Email activity
Phone/SMS
Social media activity
GPS locations
File access
Network activity
Exercise/movement
...

Probabilistic conclusions regarding source,
e.g., Likelihood Ratio



Analysis & Visualization

[Buchholz and Falk, 2005;
Grier, 2011;
Koven et al., 2016;
Gresty et al., 2016;
Kirchler et al., 2016]



Extraction

[SWDGE, 2019;
Roussev, 2016;
Casey, 2011]

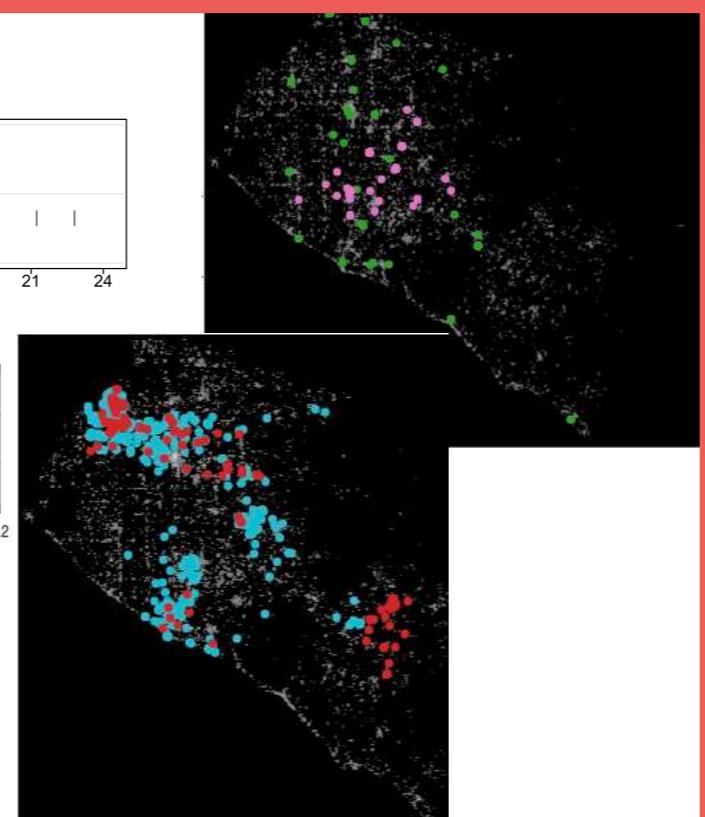
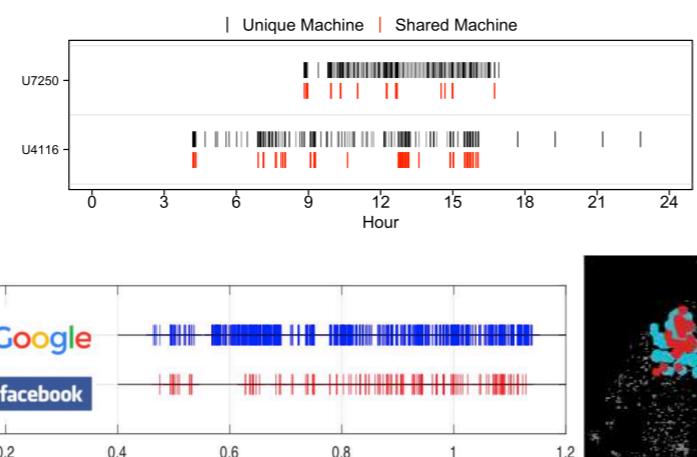
Browser requests
Web searches
Email activity
Phone/SMS
Social media activity
GPS locations
File access
Network activity
Exercise/movement
...

Analysis & Visualization

[Buchholz and Falk, 2005;
Grier, 2011;
Koven et al., 2016;
Gresty et al., 2016;
Kirchler et al., 2016]

Probabilistic conclusions regarding source, e.g., Likelihood Ratio

TOPIC OF DFRWS PAPER

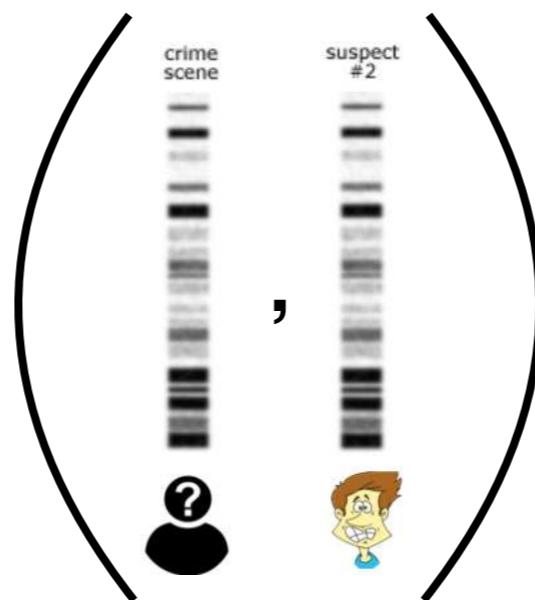


BACKGROUND

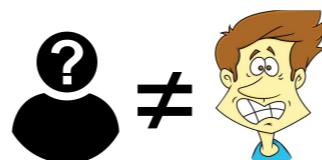
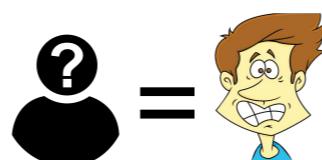
Statistical Approaches for Evaluating Forensic Evidence

Goal

Assess the likelihood of observing

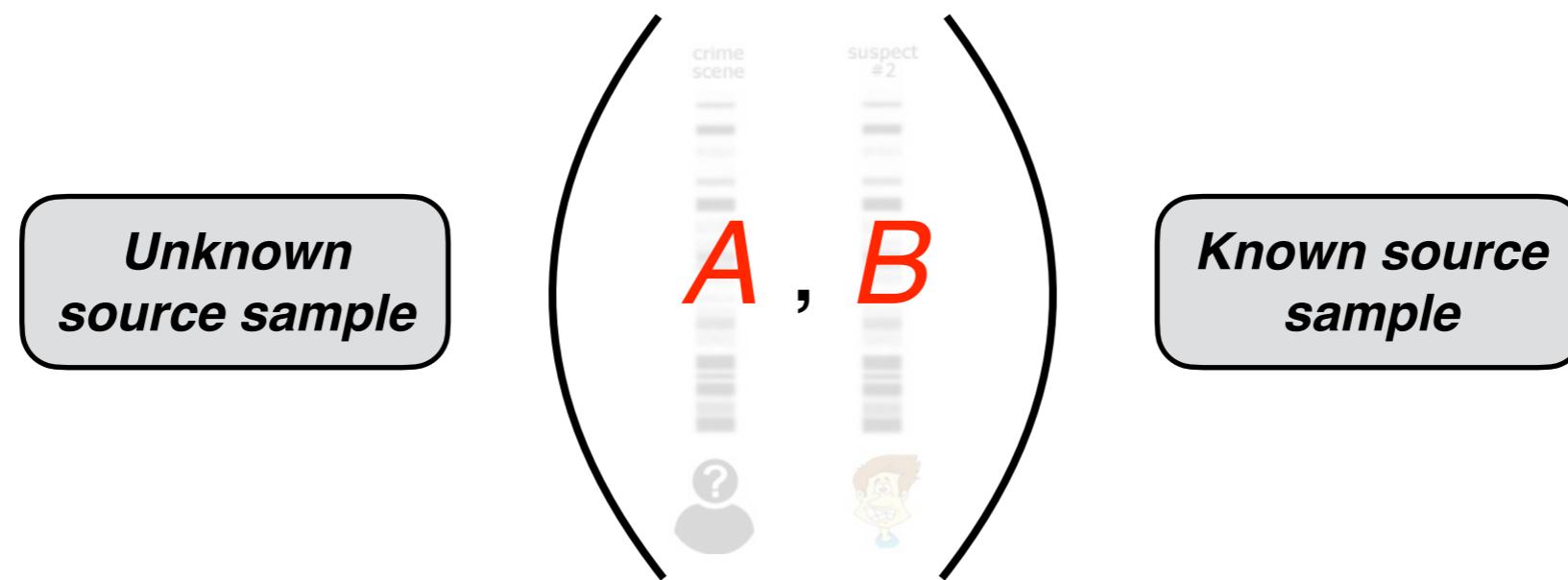


Under two competing hypotheses



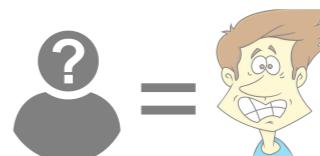
Goal

Assess the likelihood of observing



Under two competing hypotheses

H_s : (A, B) came from the same source



H_d : (A, B) came from the different sources



**Wait...why aren't we interested in the probability
of the source hypothesis *given the evidence*?**

**Wait...why aren't we interested in the probability
of the source hypothesis given the evidence?**

$$\underbrace{\frac{Pr(H_s | A, B)}{Pr(H_d | A, B)}}_{\text{posterior odds}}$$

Wait...why aren't we interested in the probability of the source hypothesis given the evidence?

$$\underbrace{\frac{Pr(H_s | A, B)}{Pr(H_d | A, B)}}_{\text{posterior odds}} = \underbrace{\frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{Pr(H_s)}{Pr(H_d)}}_{\text{prior odds}}$$



$$\underbrace{\frac{Pr(H_s | A, B)}{Pr(H_d | A, B)}}_{\text{posterior odds}} = \underbrace{\frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{Pr(H_s)}{Pr(H_d)}}_{\text{prior odds}}$$





$$\underbrace{\frac{Pr(H_s | A, B)}{Pr(H_d | A, B)}}_{\text{posterior odds}} = \underbrace{\frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{Pr(H_s)}{Pr(H_d)}}_{\text{prior odds}}$$



“Strength of Evidence”

$$\frac{\Pr(H_s | A, B)}{\Pr(H_d | A, B)} = \underbrace{\frac{Pr(H_s)}{Pr(H_d)}}_{\text{posterior odds}} \cdot \underbrace{\frac{\Pr(A, B | H_s)}{\Pr(A, B | H_d)}}_{\text{likelihood ratio}}.$$



“Weight of Evidence”

[Pierce, 1878]

$$\underbrace{\frac{\Pr(H_s)}{\Pr(H_d)}}_{\text{prior odds}}$$



The Likelihood Ratio

- ☐ Widely accepted as a “logically defensible way” to asses the strength of evidence [Willis et al., 2016]

The Likelihood Ratio

- Widely accepted as a “logically defensible way” to assess the strength of evidence [Willis et al., 2016]

- Has been applied in a variety of forensic disciplines

The Likelihood Ratio

- Widely accepted as a “logically defensible way” to assess the strength of evidence [Willis et al., 2016]
- Has been applied in a variety of forensic disciplines
 - DNA [Aitken & Stoney, 1991; Evett & Weir, 1998; Steele & Balding, 2014]
 - Fingerprints [Champod & Evett, 2001]
 - Handwriting [Bozza et al., 2008]
 - Speaker Recognition [Champod & Meuwly, 2000]

The Likelihood Ratio

- Widely accepted as a “logically defensible way” to assess the strength of evidence [Willis et al., 2016]
- Has been applied in a variety of forensic disciplines
 - DNA [Aitken & Stoney, 1991; Evett & Weir, 1998; Steele & Balding, 2014]
 - Fingerprints [Champod & Evett, 2001]
 - Handwriting [Bozza et al., 2008]
 - Speaker Recognition [Champod & Meuwly, 2000]
- Studies demonstrating its understanding

The Likelihood Ratio

- Widely accepted as a “logically defensible way” to assess the strength of evidence [Willis et al., 2016]
- Has been applied in a variety of forensic disciplines
 - DNA [Aitken & Stoney, 1991; Evett & Weir, 1998; Steele & Balding, 2014]
 - Fingerprints [Champod & Evett, 2001]
 - Handwriting [Bozza et al., 2008]
 - Speaker Recognition [Champod & Meuwly, 2000]
- Studies demonstrating its understanding
 - Misconceptions [Martire et al., 2013, Thompson and Newman, 2015, Thompson et al., 2018]

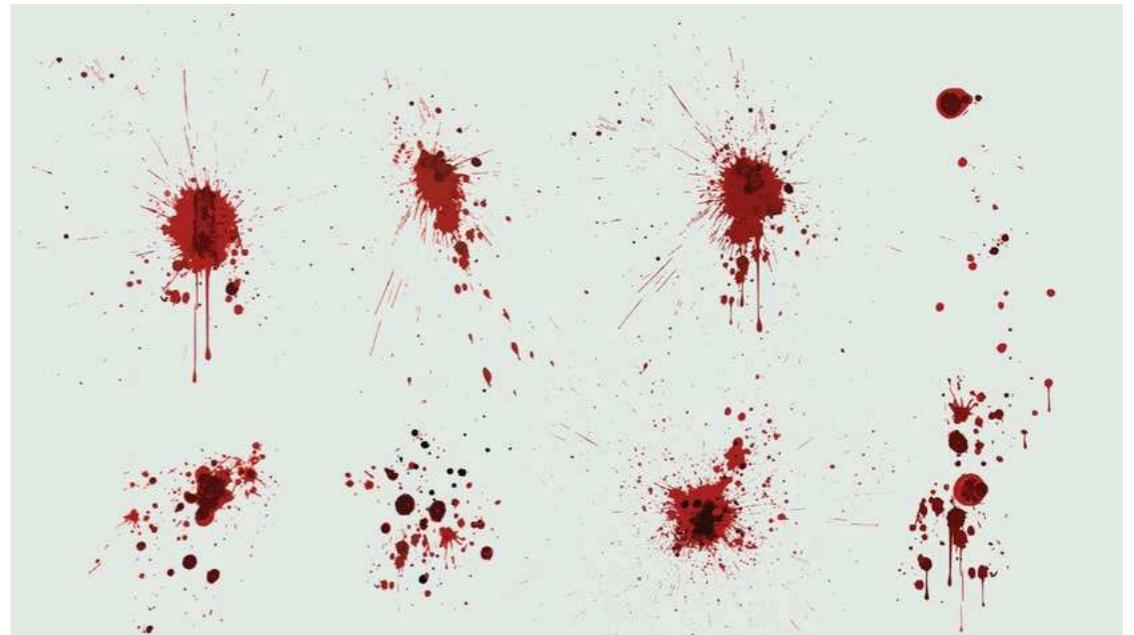
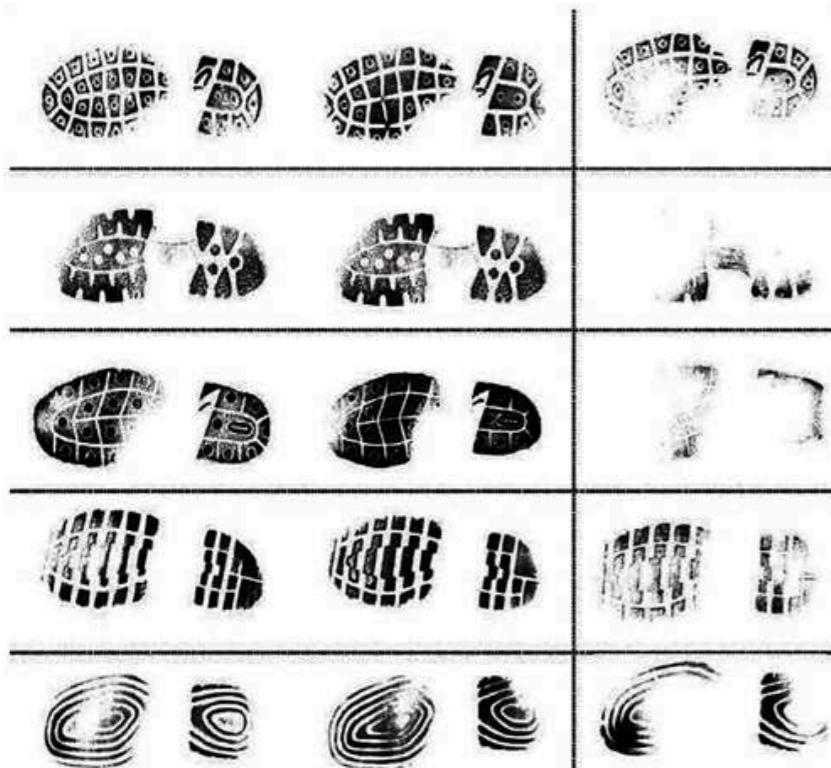
The Likelihood Ratio

- Widely accepted as a “logically defensible way” to assess the strength of evidence [Willis et al., 2016]
- Has been applied in a variety of forensic disciplines
 - DNA [Aitken & Stoney, 1991; Evett & Weir, 1998; Steele & Balding, 2014]
 - Fingerprints [Champod & Evett, 2001]
 - Handwriting [Bozza et al., 2008]
 - Speaker Recognition [Champod & Meuwly, 2000]
- Studies demonstrating its understanding
 - Misconceptions [Martire et al., 2013, Thompson and Newman, 2015, Thompson et al., 2018]
 - Verbal Equivalents [e.g., AFSP, 2009]

Why not always use the LR?

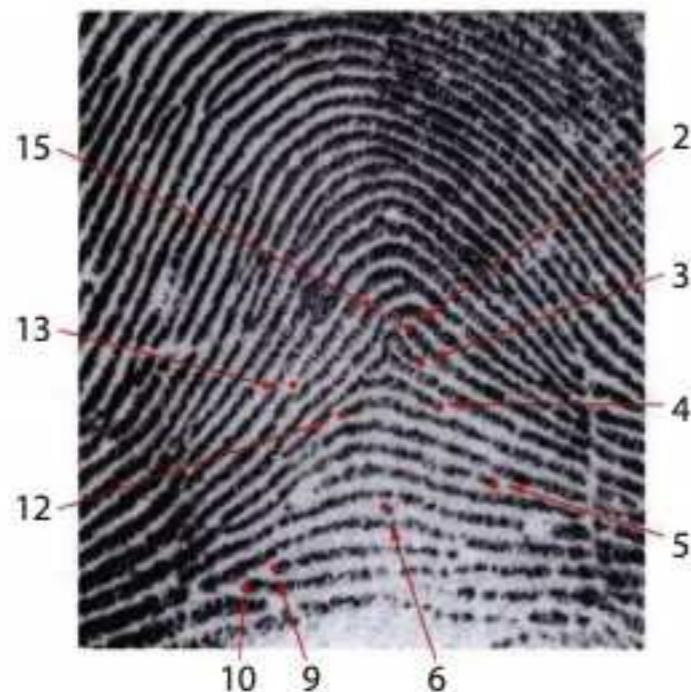
Why not always use the LR?

Complexity: Evidence can be high-dimensional



Why not always use the LR?

- **Complexity:** Evidence can be high-dimensional
- **Feature Selection:** Wide variety of features to consider



[Fine, 2016]



[Stern, 2017]

Why not always use the LR?

- **Complexity:** Evidence can be high-dimensional
- **Feature Selection:** Wide variety of features to consider
- **Appropriate Probability Models:** Must describe variation within a given source and between different sources

$$H_s: \text{?} = \text{!}$$
$$H_d: \text{?} \neq \text{!}$$

Why not always use the LR?

- Complexity:** Evidence can be high-dimensional
- Feature Selection:** Wide variety of features to consider
- Appropriate Probability Models:** Must describe variation within a given source and between different sources
- Reference Population:** Difficult to identify a *relevant* reference population to estimate model parameters & perform validation studies

Score-based Approaches

- Measure similarity between A and B via a *score function*

$$\Delta(A, B)$$

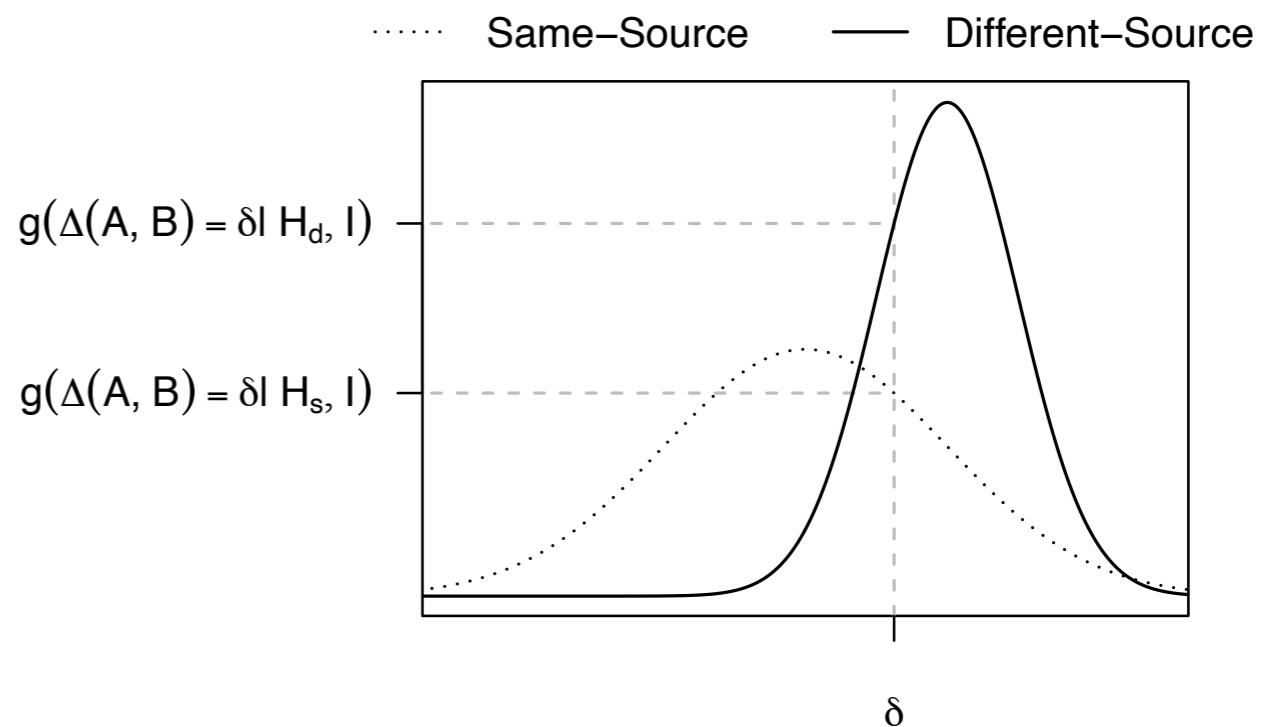
Score-based Approaches

- Measure similarity between A and B via a *score function*

$$\Delta(A, B)$$

- **Score-based Likelihood Ratio:** Compute a LR for the observed score

$$SLR_{\Delta} = \frac{g(\Delta(A, B) = \delta | H_s)}{g(\Delta(A, B) = \delta | H_d)}$$



Score-based Approaches

- Measure similarity between A and B via a *score function*

$$\Delta(A, B)$$

- **Score-based Likelihood Ratio:** Compute a LR for the observed score

$$SLR_{\Delta} = \frac{g(\Delta(A, B) = \delta | H_s)}{g(\Delta(A, B) = \delta | H_d)}$$

Gaining popularity in a variety of forensic disciplines

- Chemical Concentrations [Bolck et al., 2015]
- Speaker Recognition [Gonzalez-Rodriguez et al., 2007]
- Fingerprints [Alberink et al., 2013; Neumann et al., 2015]
- Handwriting [Hepler et al., 2012]

Evidence Evaluation Approaches

- **Likelihood Ratio:** Models evidence directly

$$LR = \frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}$$

- **Score-based Likelihood Ratio:** Models low-dimensional summary of the evidence, $\Delta(A, B)$

$$SLR_{\Delta} = \frac{g(\Delta(A, B) = \delta | H_s)}{g(\Delta(A, B) = \delta | H_d)}$$

Evidence Evaluation Approaches

- **Likelihood Ratio:** Models evidence directly

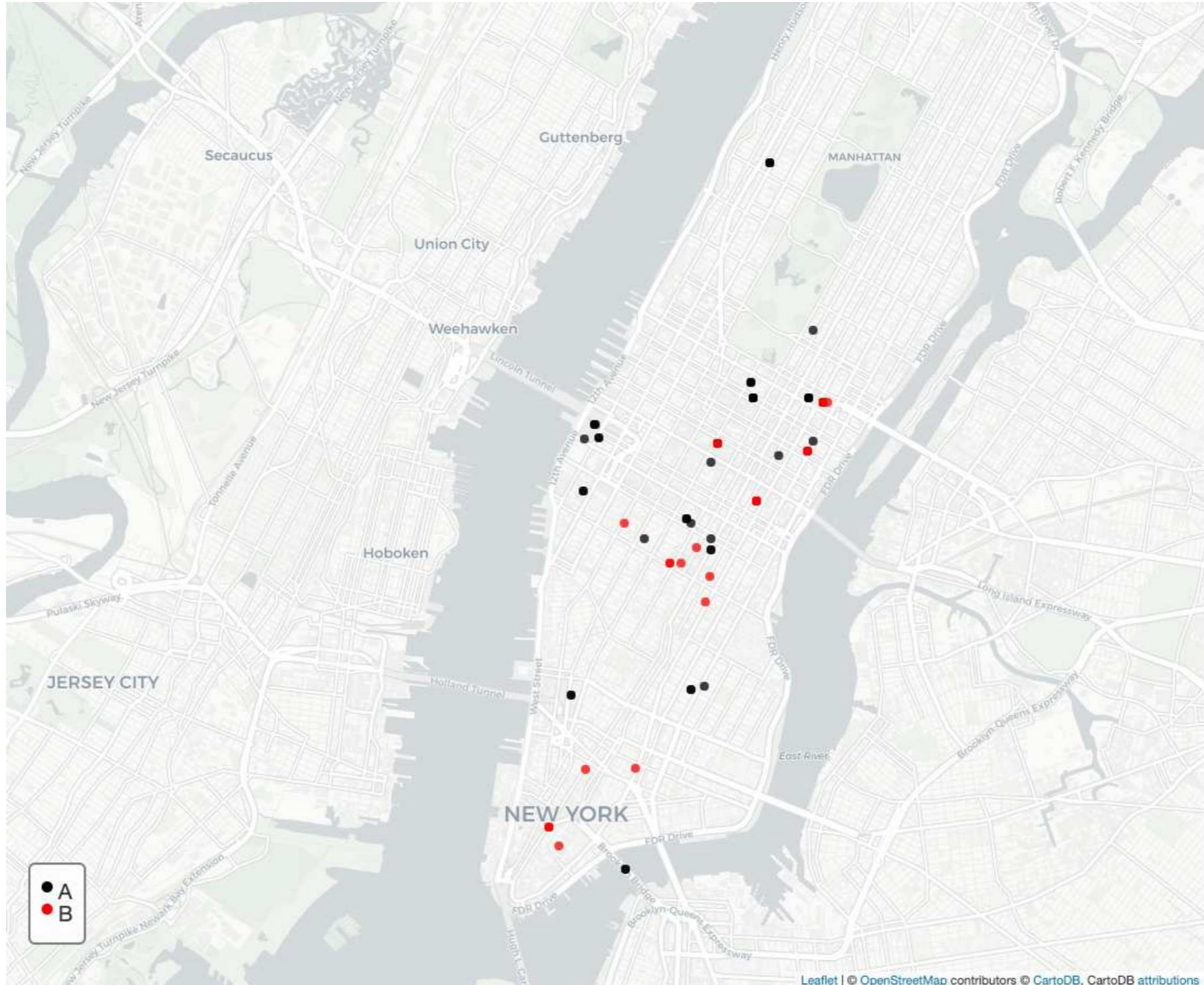
$$LR = \frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}$$

- **Score-based Likelihood Ratio:** Models low-dimensional summary of the evidence, $\Delta(A, B)$

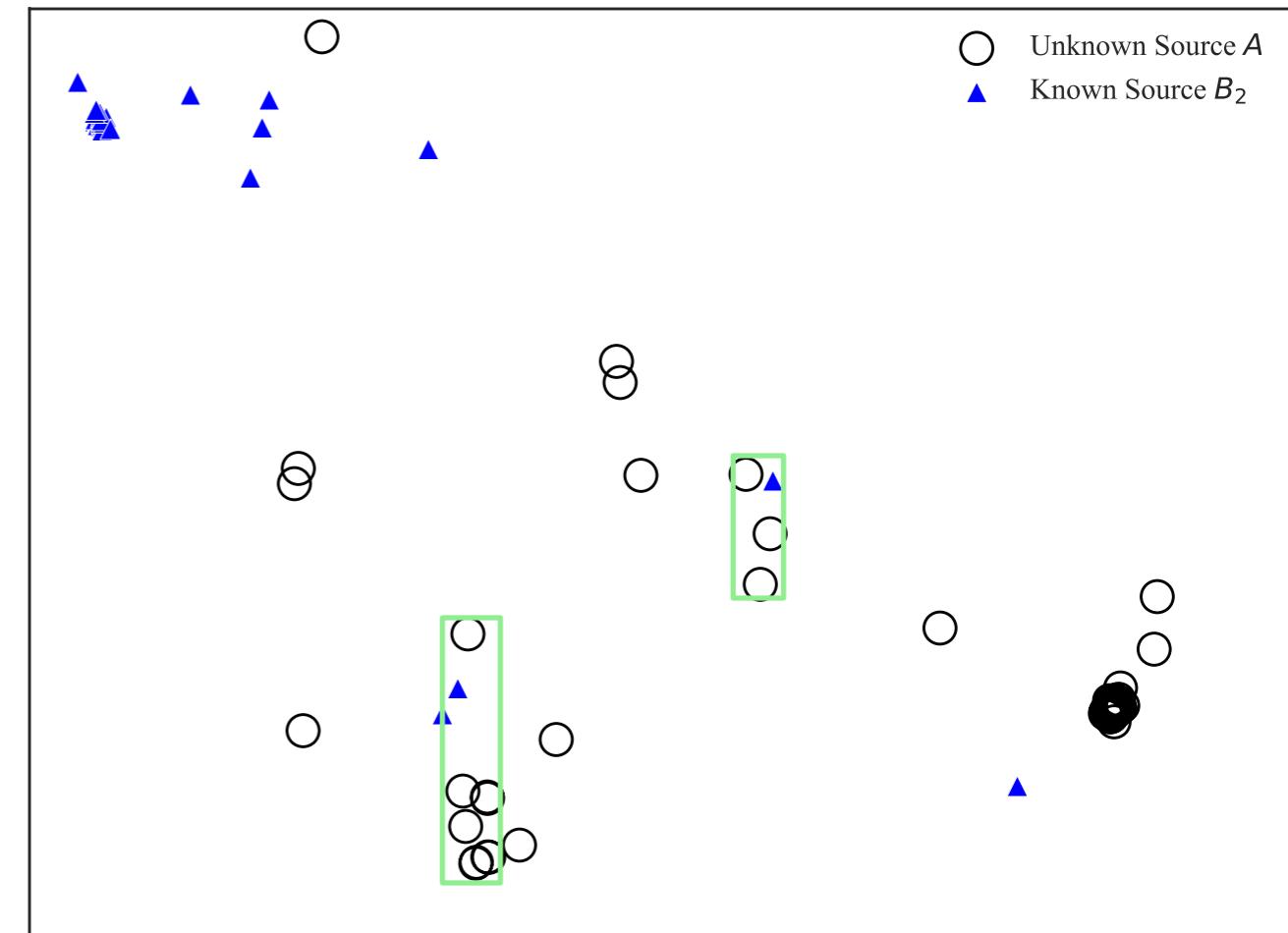
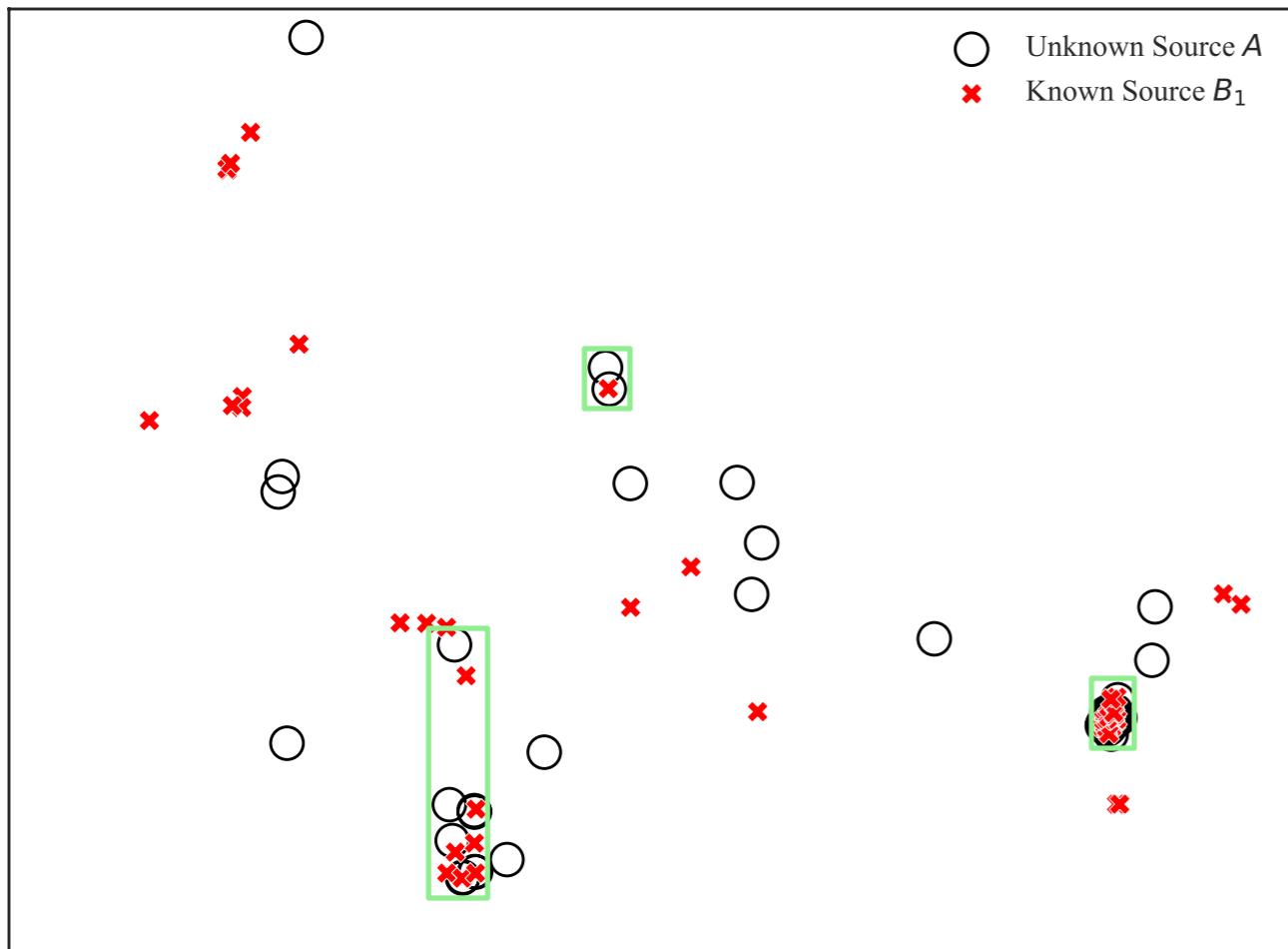
$$SLR_{\Delta} = \frac{g(\Delta(A, B) = \delta | H_s)}{g(\Delta(A, B) = \delta | H_d)}$$

CONTRIBUTION

Quantifying the Strength of Geolocated Event Evidence



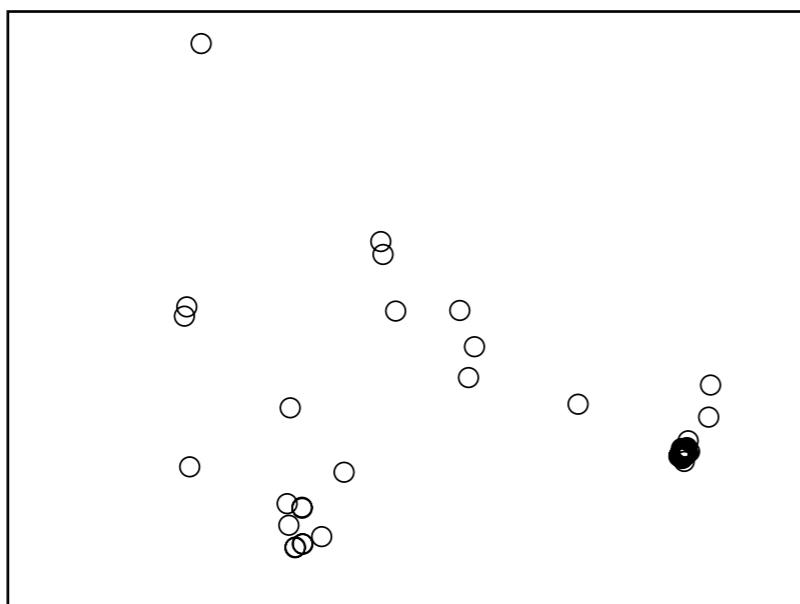
Geofence Warrants



Revisiting the LR

$$LR = \frac{Pr(A, B | H_s)}{Pr(A, B | H_d)}$$

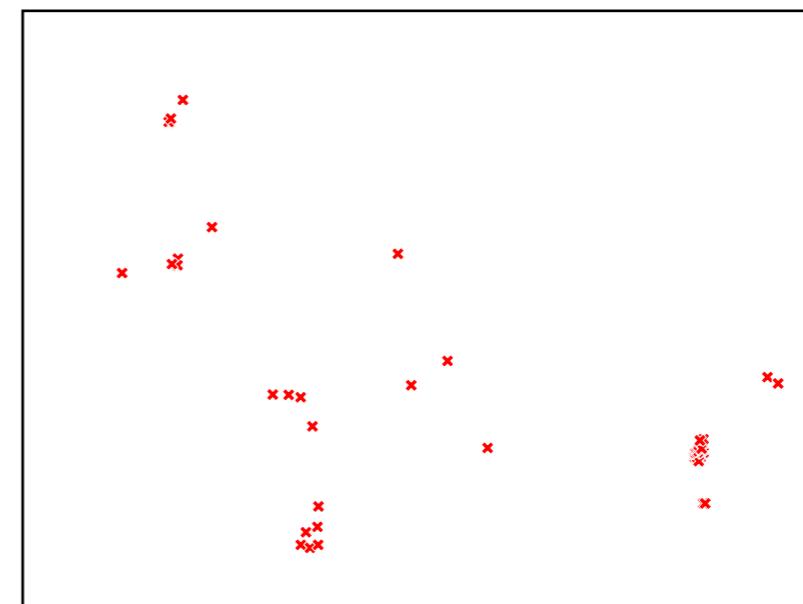
A



Unknown source events



B



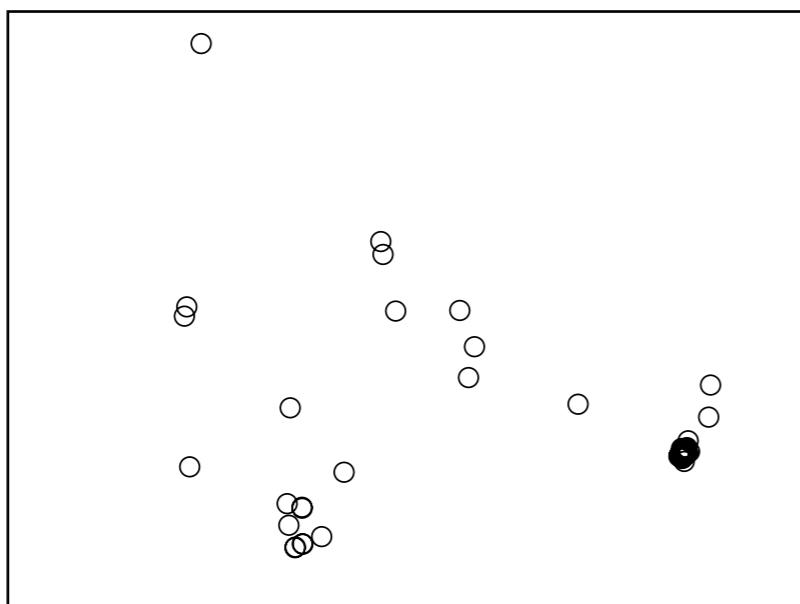
Known source events



Revisiting the LR

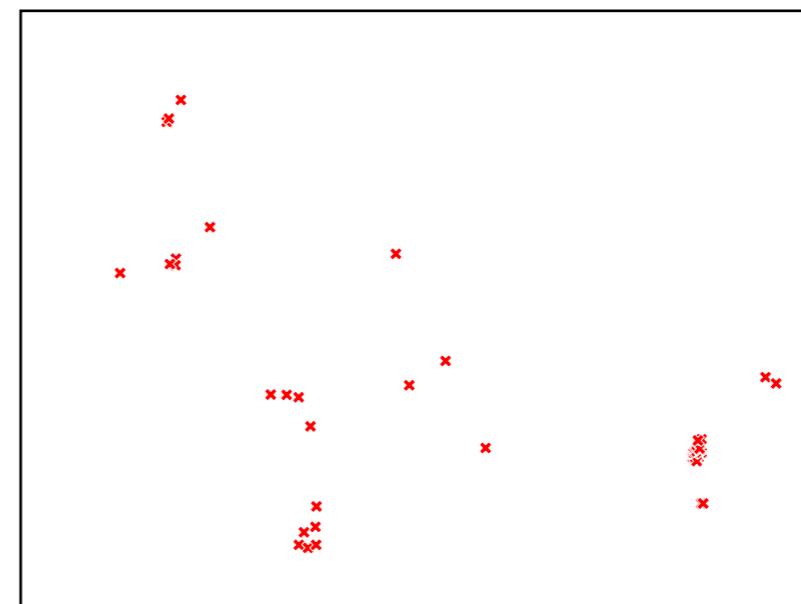
$$LR = \frac{Pr(A, B | H_s)}{Pr(A, B | H_d)} \Rightarrow \dots \Rightarrow LR = \frac{f(B | A, H_s)}{f(B | H_d)}$$

A



Unknown source events

B

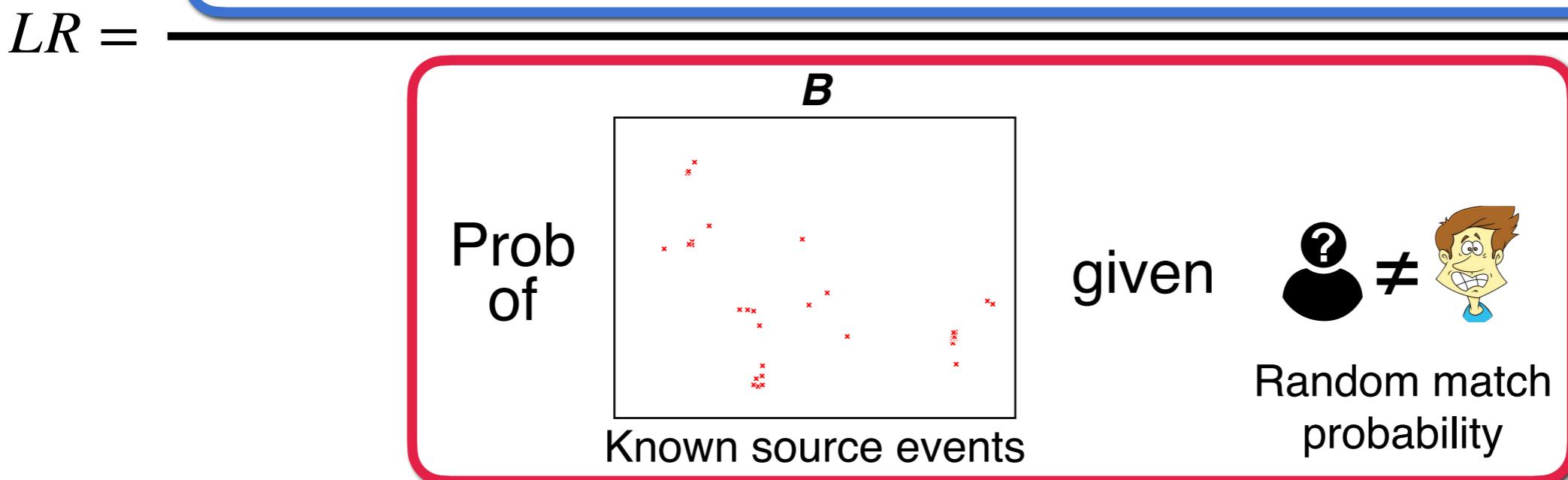
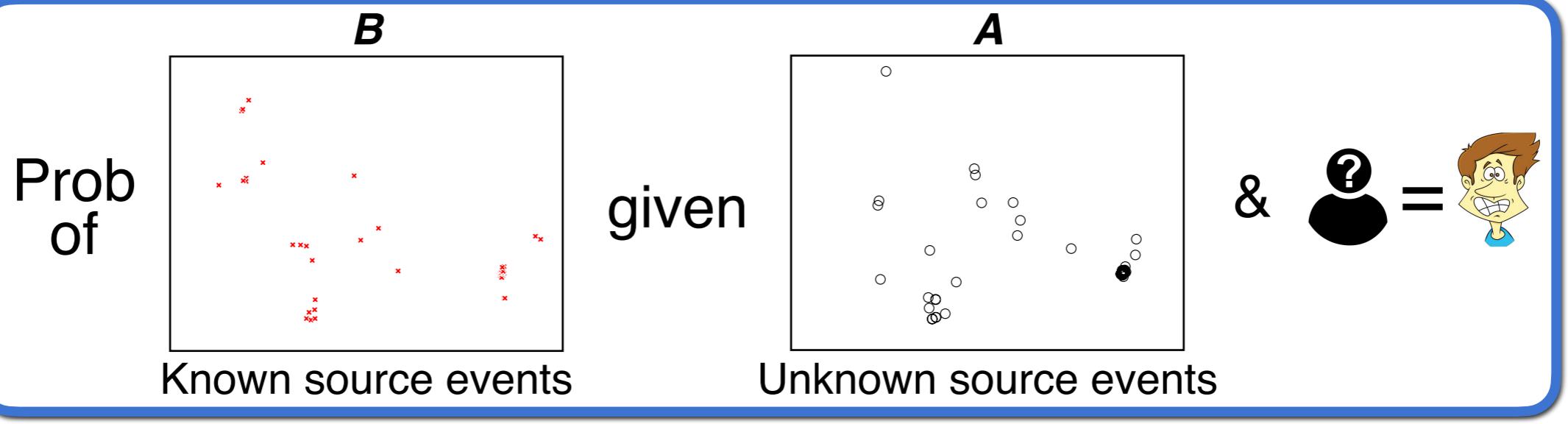


Known source events

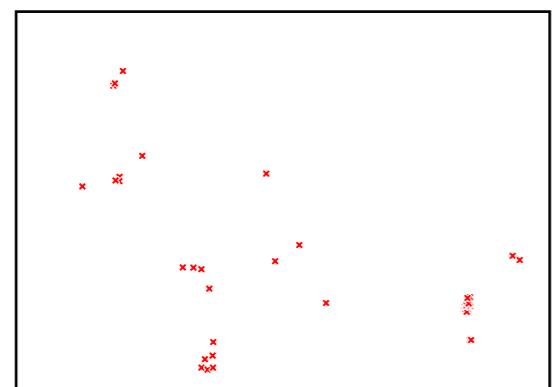


Revisiting the LR

$$LR = \frac{Pr(A, B | H_s)}{Pr(A, B | H_d)} \Rightarrow \dots \quad \text{Maths} \quad \Rightarrow LR = \frac{f(B | A, H_s)}{f(B | H_d)}$$



Prob
of

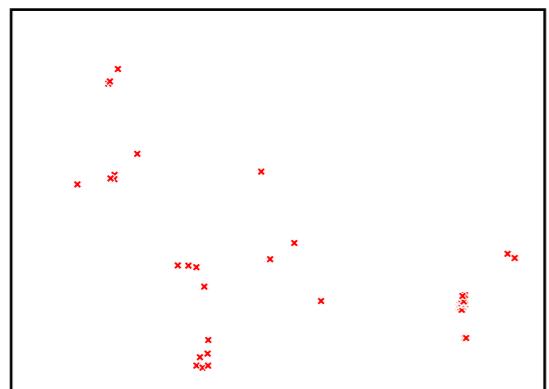


Known source events

given
Random match probability

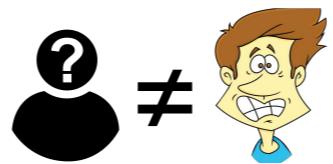
Prob
of

B

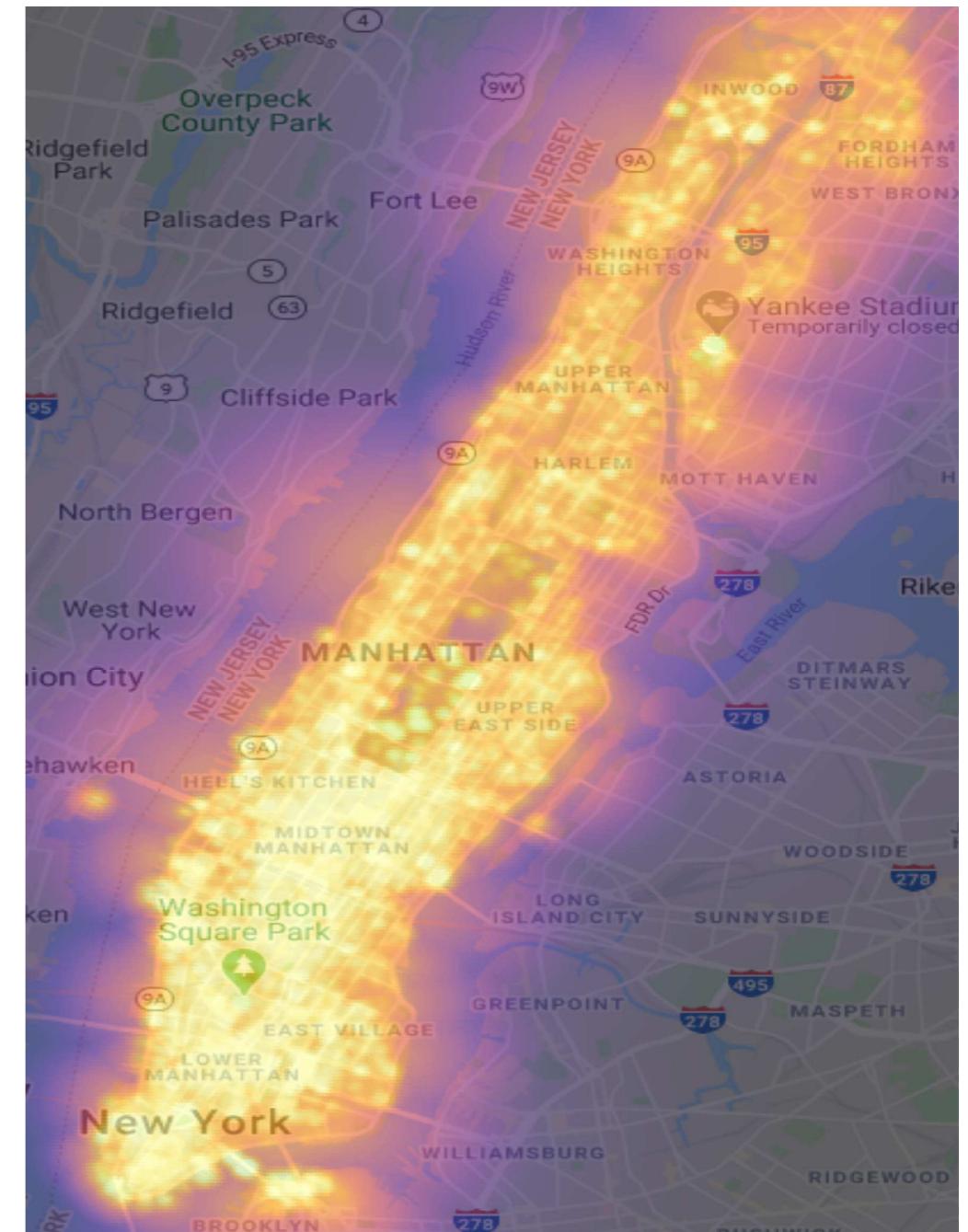
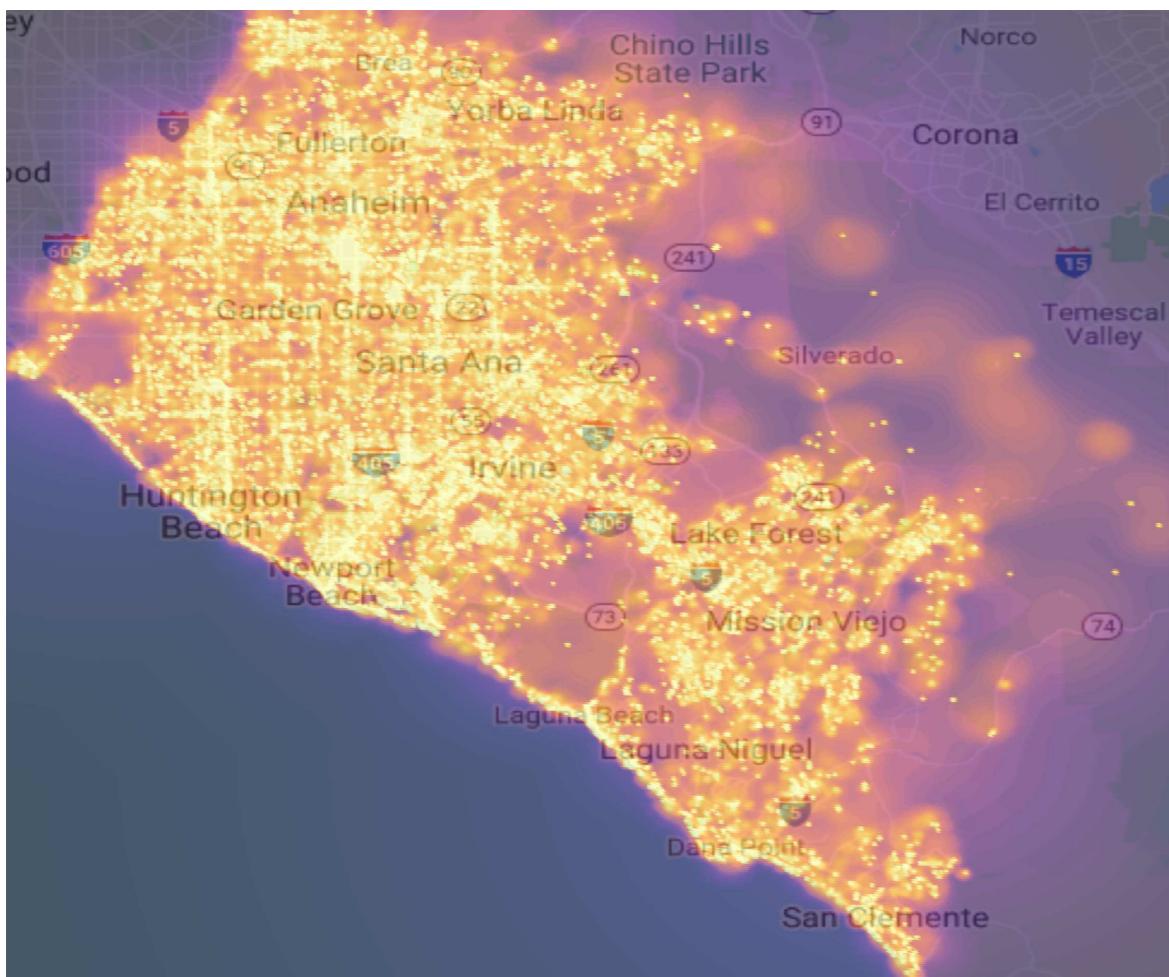


Known source events

given



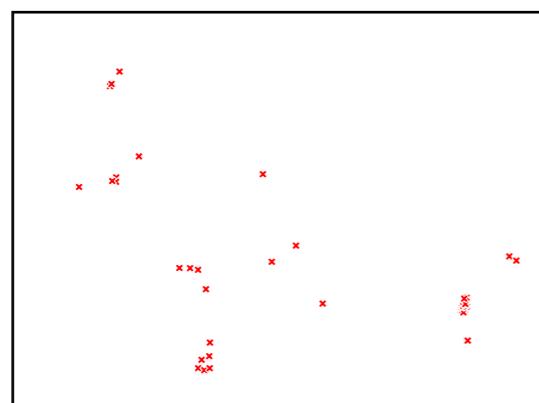
Random match
probability



Adaptive Bandwidth
Kernel Density Estimators
[Breiman et al., 1977]

Prob
of

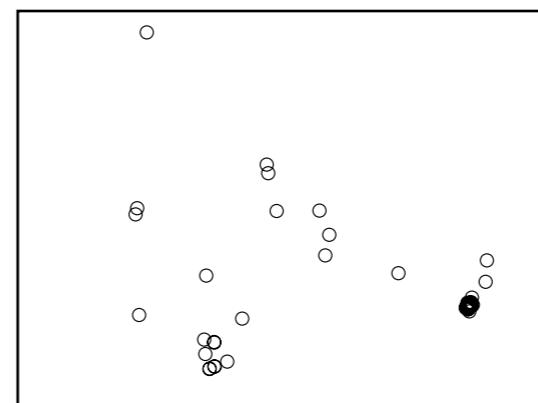
B



Known source events

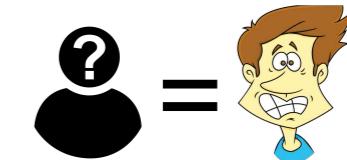
given

A



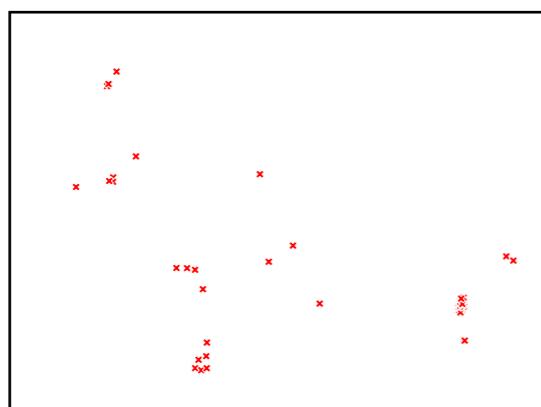
Unknown source events

&



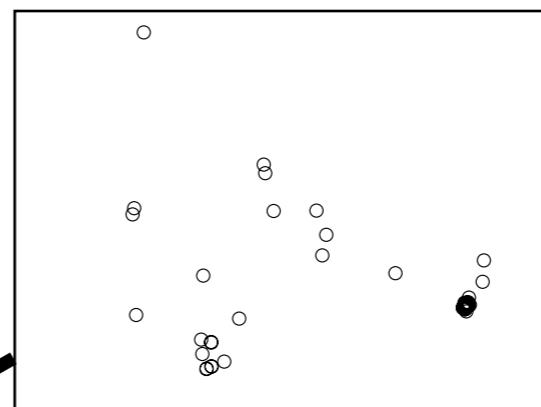
Prob
of

B



Known source events

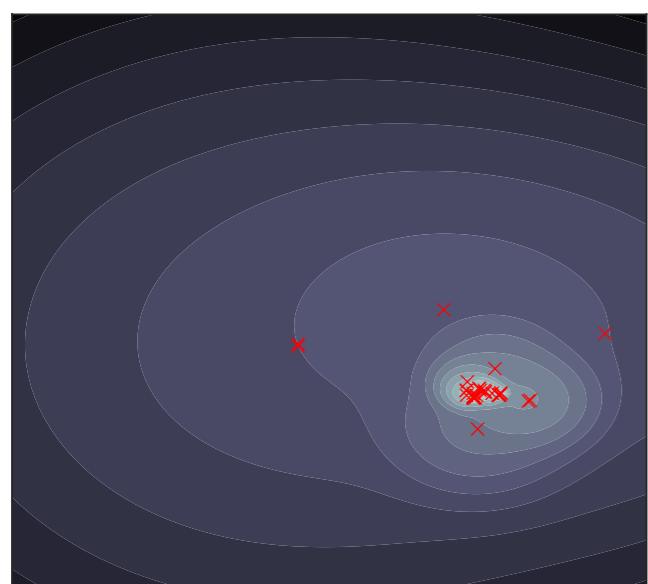
A



given

Unknown source events

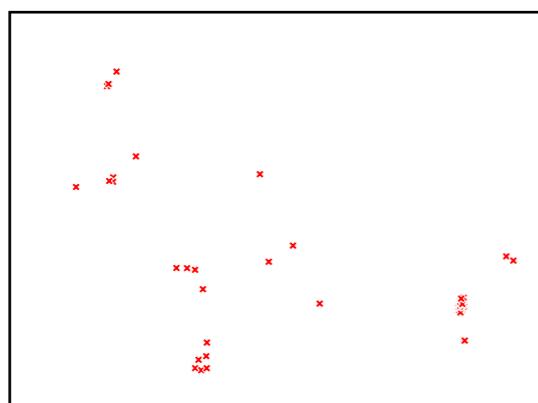
&  = 



Individual Component

Prob
of

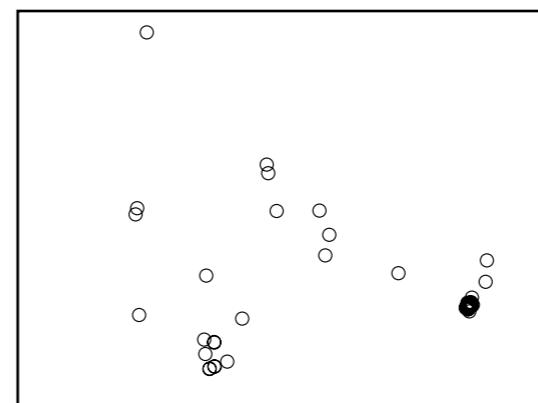
B



Known source events

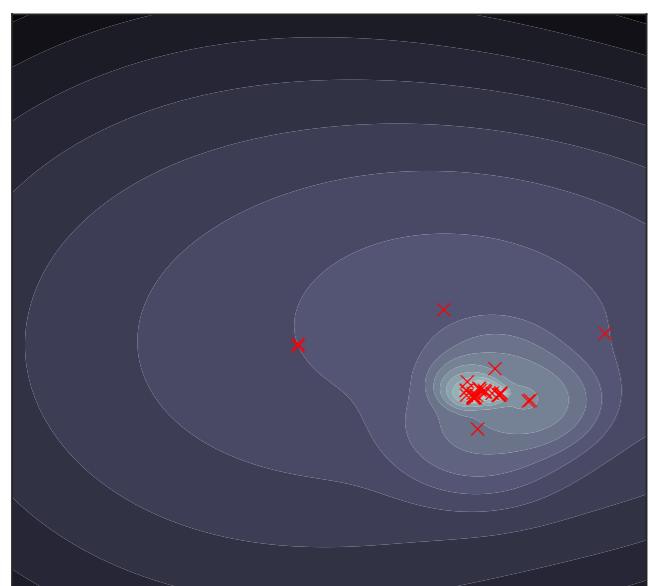
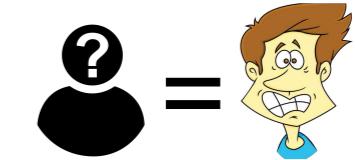
given

A

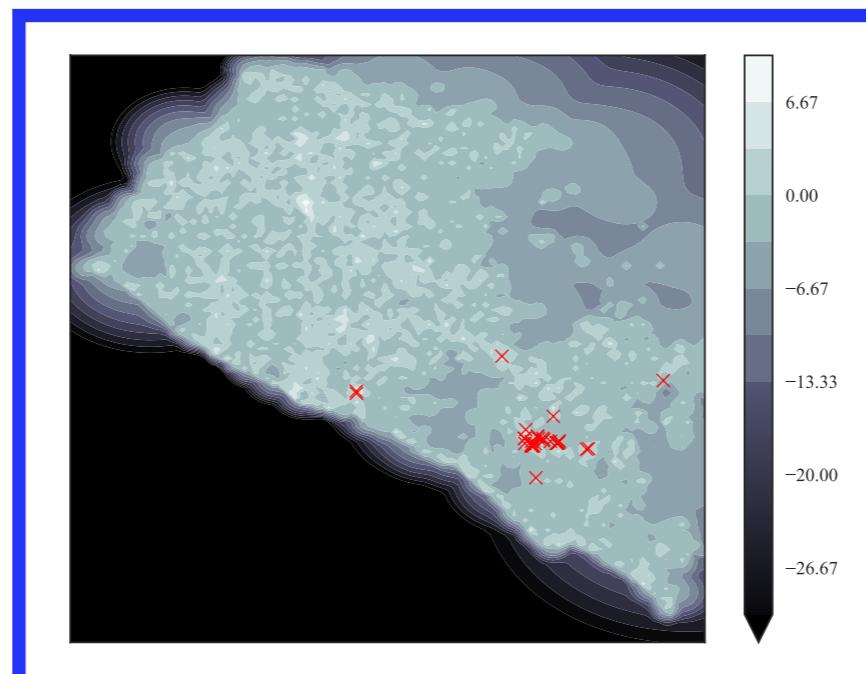


Unknown source events

&



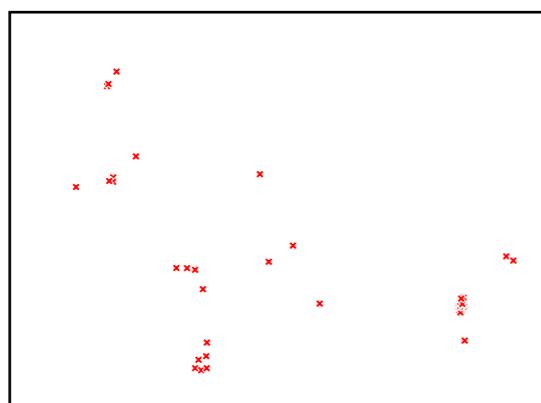
Individual Component



Population Component

Prob
of

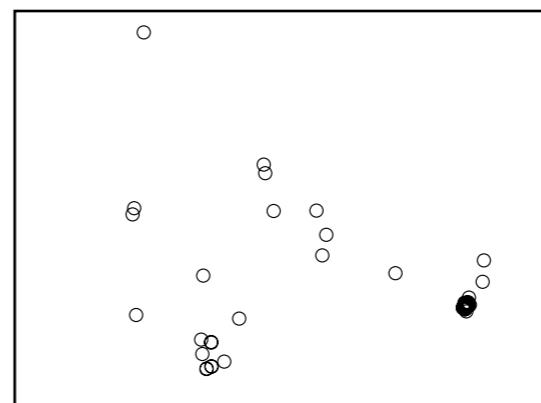
B



Known source events

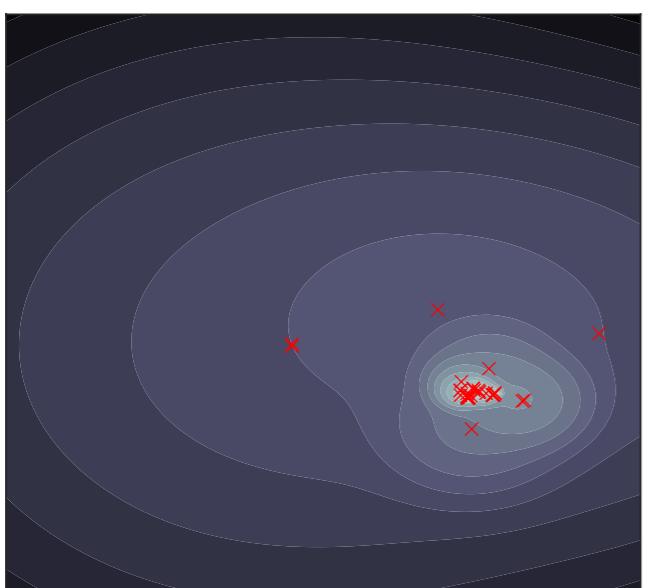
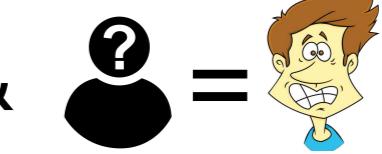
given

A

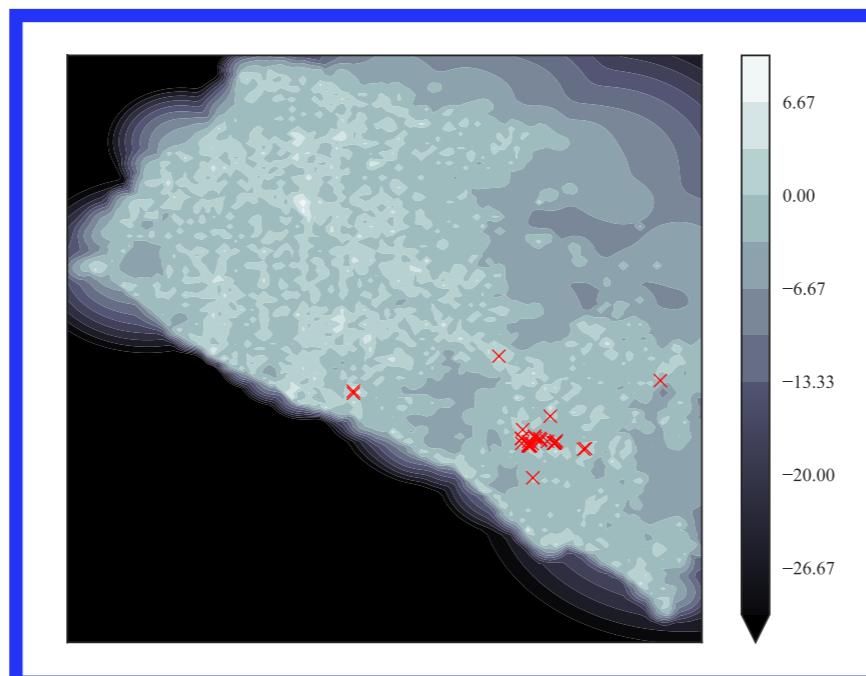


Unknown source events

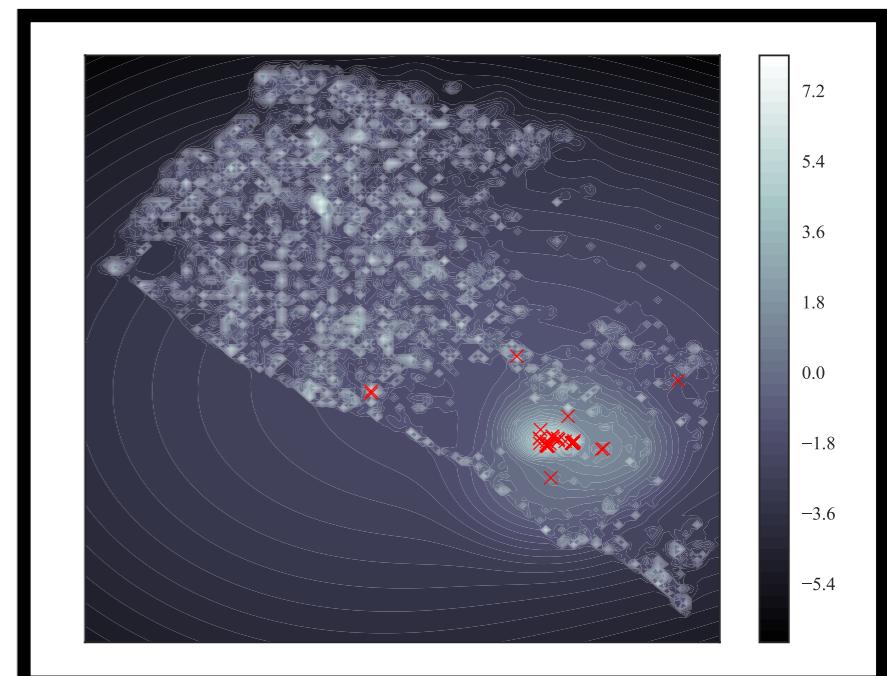
&



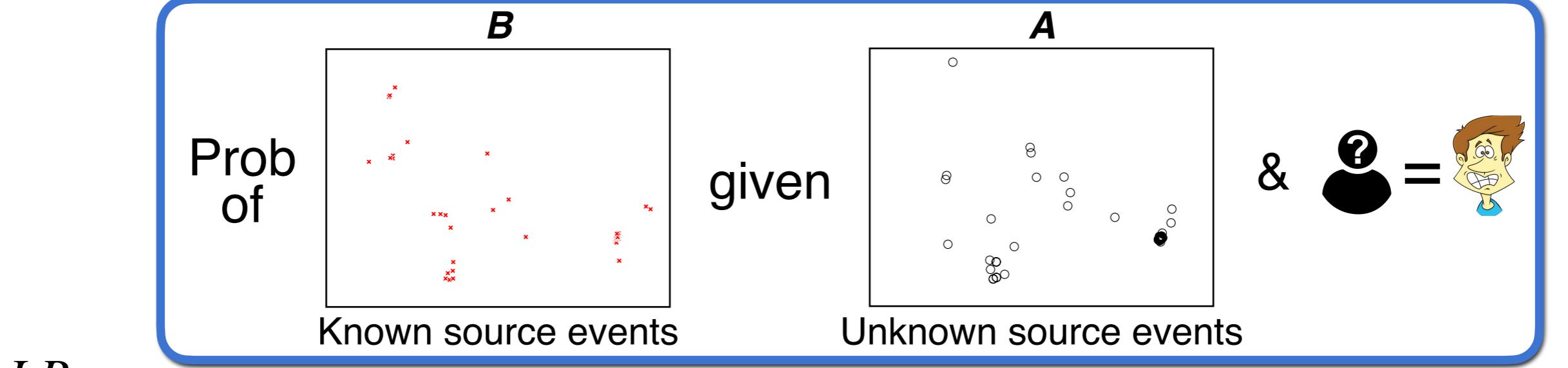
Individual Component



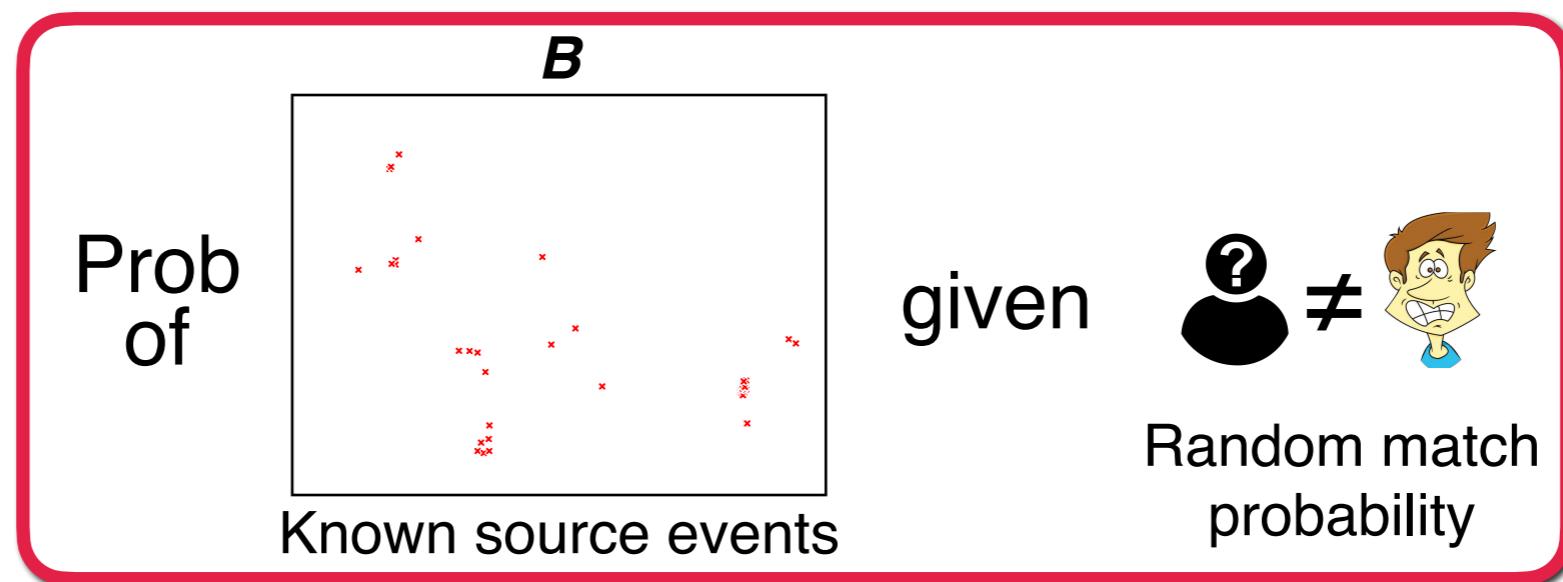
Population Component



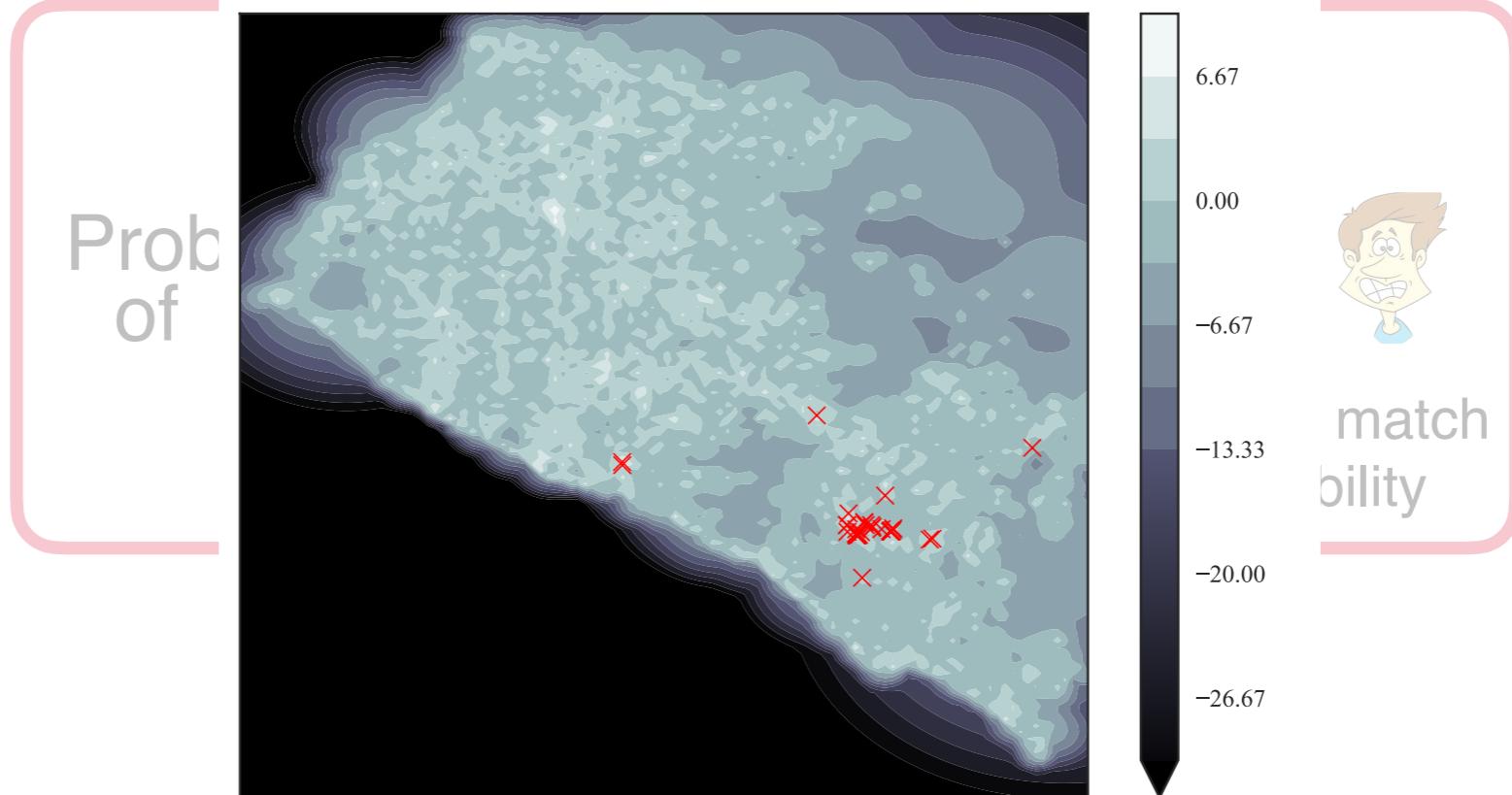
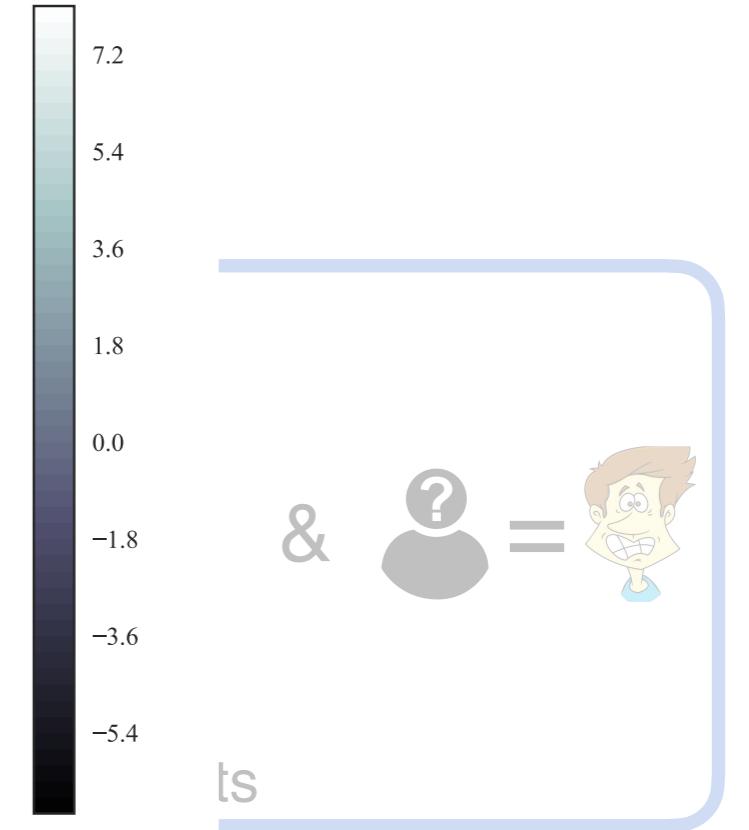
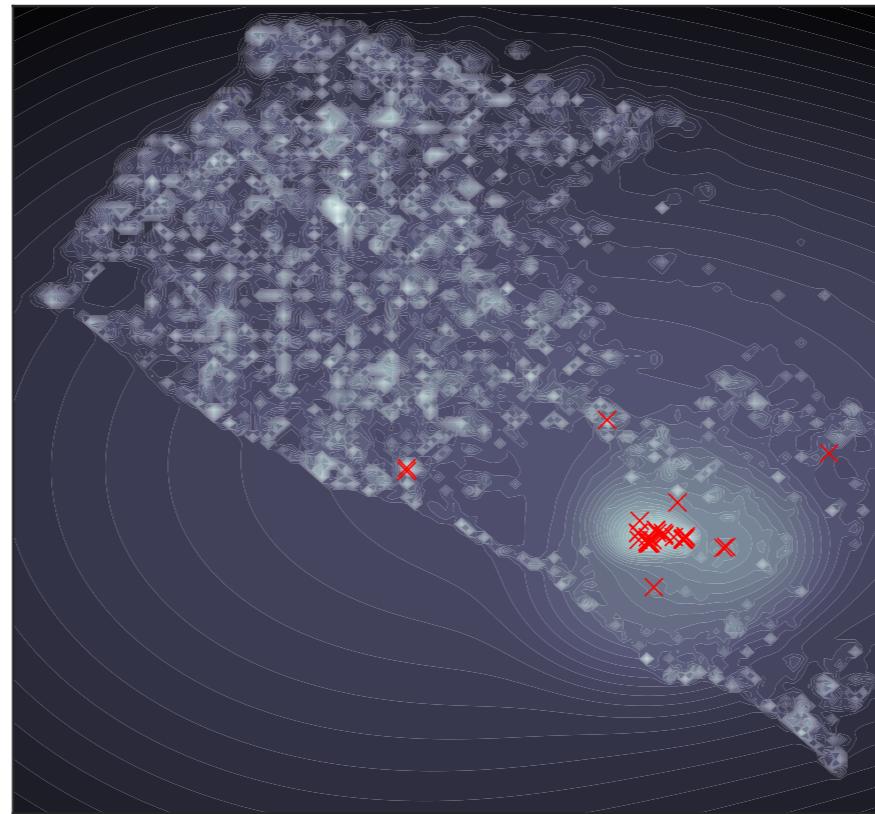
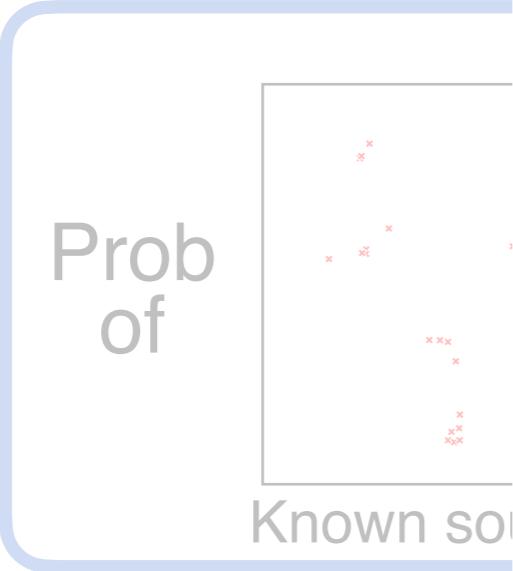
Mixture
*80% individual,
20% population*



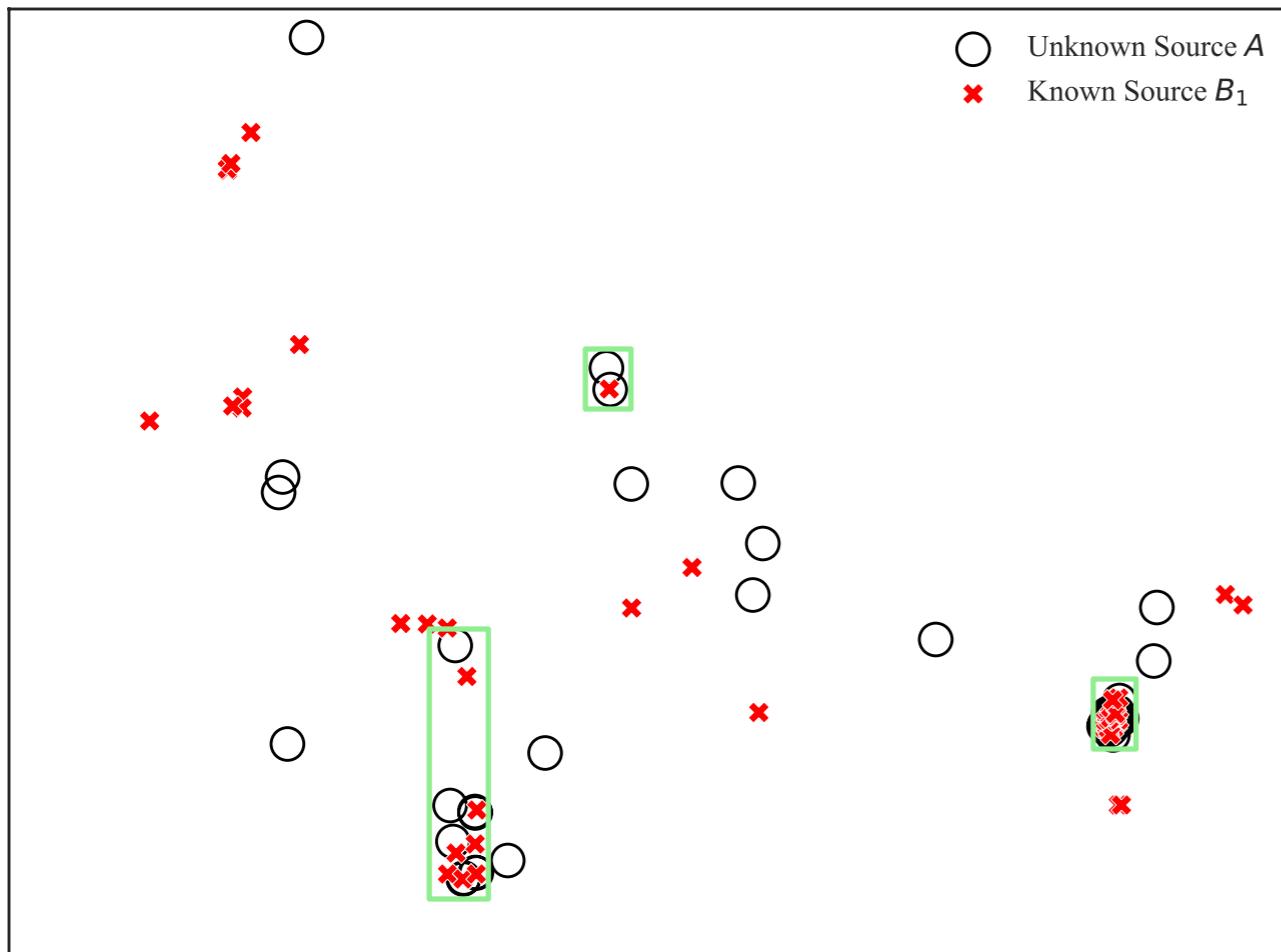
$$LR = \frac{\text{Prob of } B \text{ given } A}{\text{Prob of } B \text{ given } \text{Random match probability}}$$



$$LR = \frac{\text{Prob of Known sol}}{\text{Prob of matchability}}$$

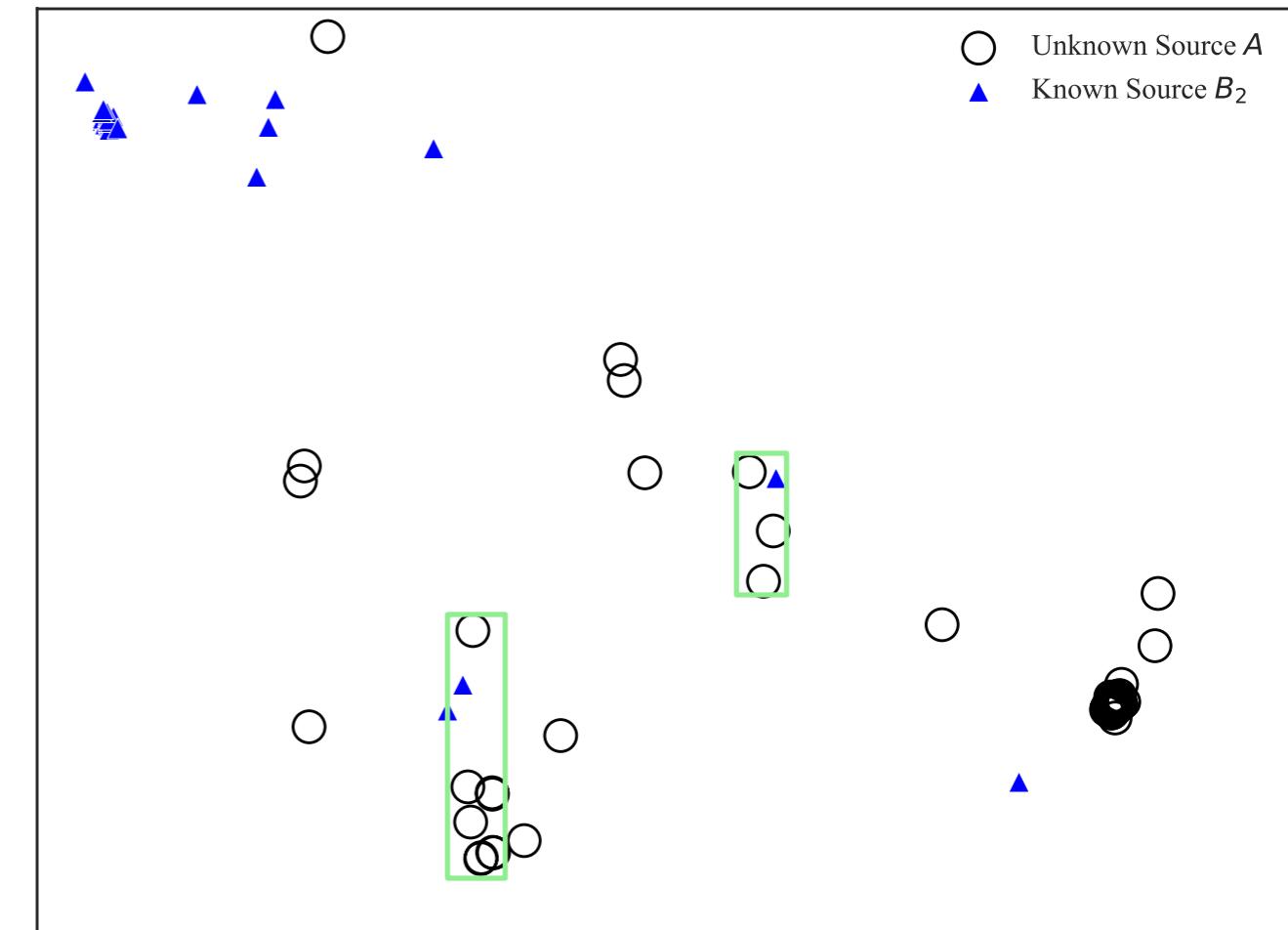


Geofence Warrants Revisited



$$LR \approx 1137$$

Strong support for same-source hypothesis



$$LR \approx 2.8 \times 10^{-28}$$

Exclusion

Case Study

- Collected Twitter data from May 2015 to Feb 2016
 - Orange County, CA
 - Manhattan, New York, NY
- *A* and *B* are consecutive months from the same account

Region	Accounts	Visits in <i>A</i>	Visits in <i>B</i>
OC	6,714	44,310 (6.6)	38,697 (5.8)
NY	13,523	72,799 (5.4)	65,852 (4.9)

Case Study

- Collected Twitter data from May 2015 to Feb 2016
 - Orange County, CA
 - Manhattan, New York, NY
- A and B are consecutive months from the same account

Region	Accounts	Visits in A	Visits in B
OC	6,714	44,310 (6.6)	38,697 (5.8)
NY	13,523	72,799 (5.4)	65,852 (4.9)

- Results based on stratified sample based on n_a and n_b for different-source evidence

Results

Region	Method ¹	TP Rate ²	FP Rate ²	AUC
OC	LR	0.380	0.038	0.845
	SLR	0.614	0.162	0.783
NY	LR	0.285	0.089	0.768
	SLR	0.511	0.235	0.685

(1) LR with $\alpha(n_a)$ weights; SLR using earth mover's distance with account weights

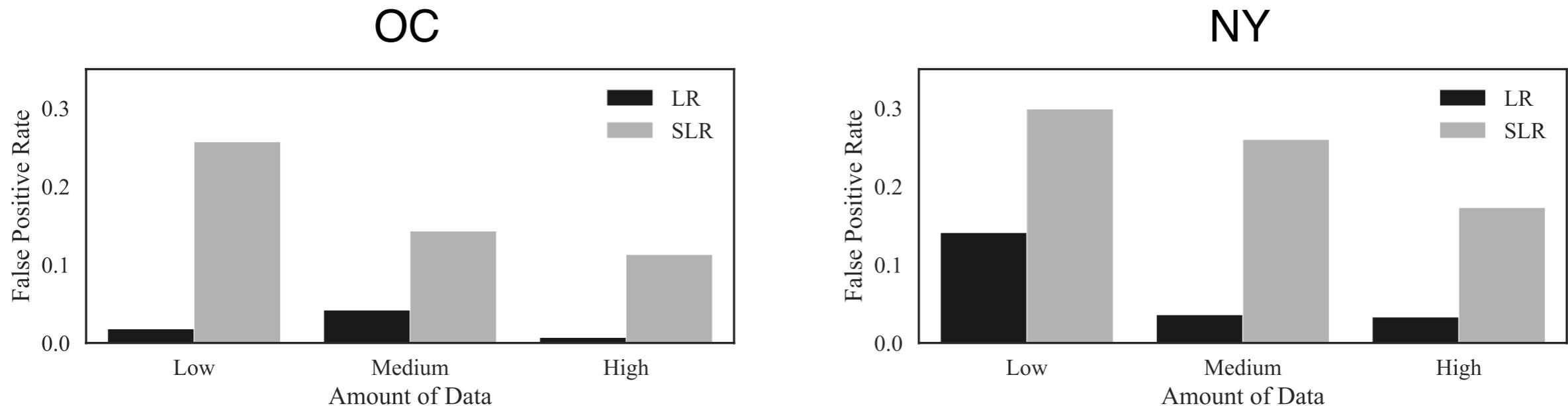
(2) LR & SLR threshold is 1

Results

Region	Method ¹	TP Rate ²	FP Rate ²	AUC
OC	LR	0.380	0.038	0.845
	SLR _{EMD}	0.614	0.162	0.783
NY	LR	0.285	0.089	0.768
	SLR _{EMD}	0.511	0.235	0.685

(1) LR with $\alpha(n_a)$ weights; SLR using earth mover's distance with account weights

(2) LR & SLR threshold is 1



Future Directions and Summary

Future Directions

- Reference Data:** Collect & share relevant digital data amongst law enforcement & researchers, e.g., start to build CODIS-like databases.

- Discovery:** Finding the most likely known source in a database given an unknown source sample...quickly.

Summary

- **Statistical approaches** play a key role in the **forensic analysis** of a wide variety of evidence.
- **Digital evidence** is lagging behind other forensic disciplines.
- Presented a framework for estimating **likelihood ratios** for forensic applications with **geolocated event data**.

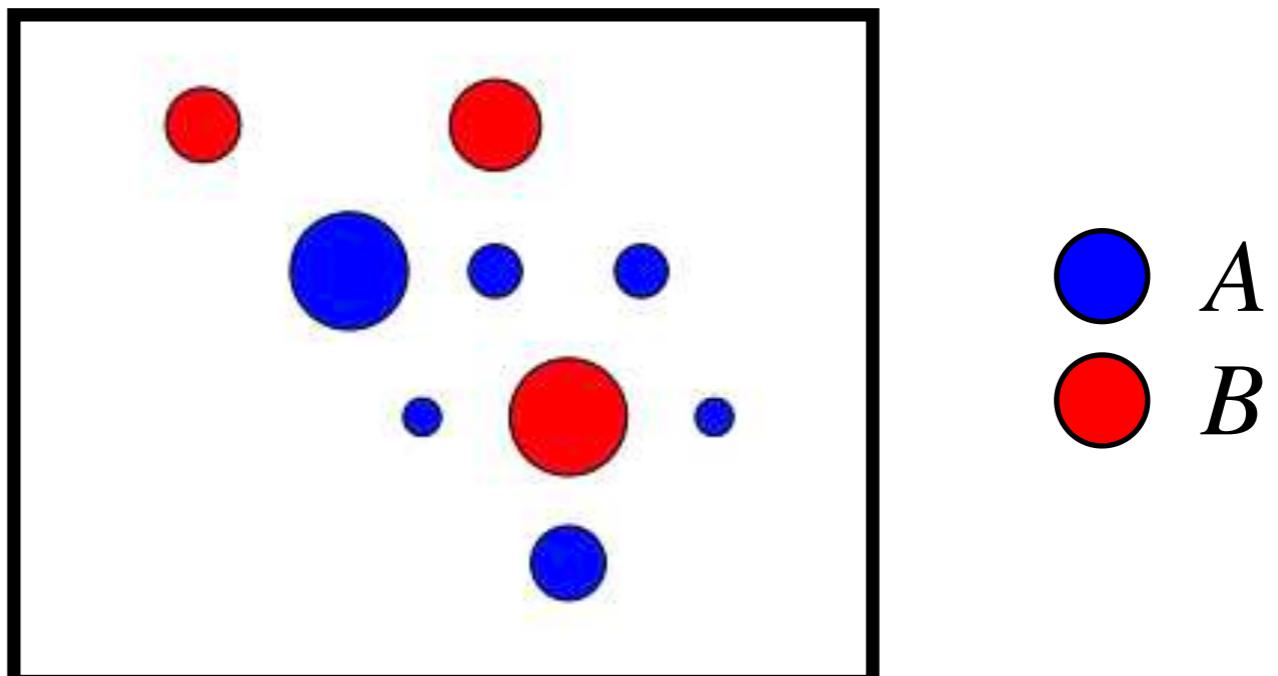
Questions

Appendix

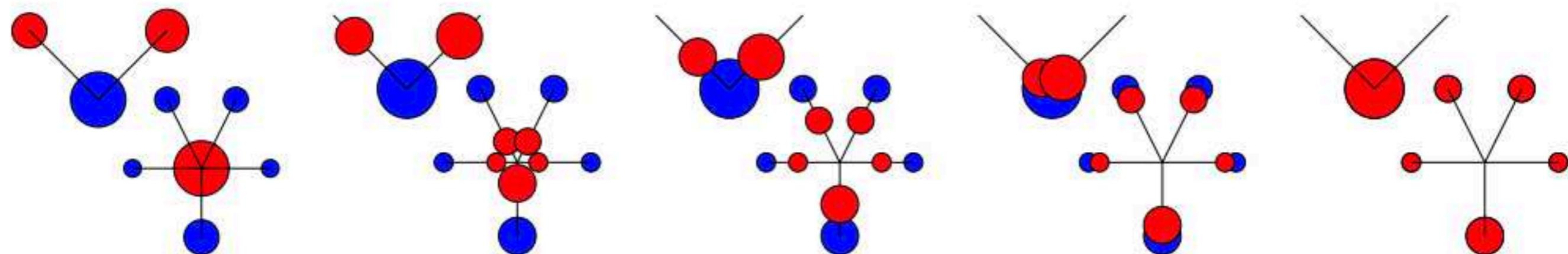
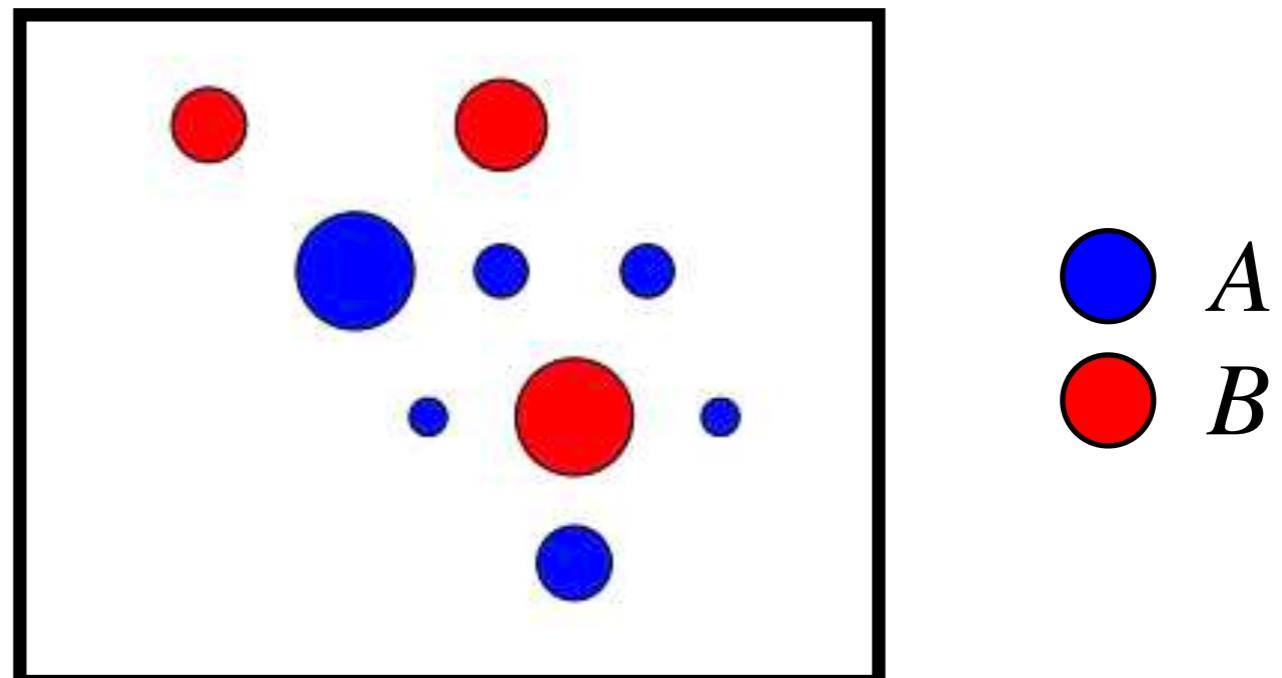
Score Functions

- Techniques to characterize spatial point patterns generally fall into two categories [Haggett, 1977]
 - Distance-based
 - Area-based
- Use distance-based score functions $\Delta(A, B)$ to quantify the similarity of the points within the sets A and B
- Incorporate area-based information via weights Ω^a, Ω^b

Earth-mover's distance
 $EMD(B, A \mid \Omega^b, \Omega^a)$



Earth-mover's distance
 $EMD(B, A \mid \Omega^b, \Omega^a)$



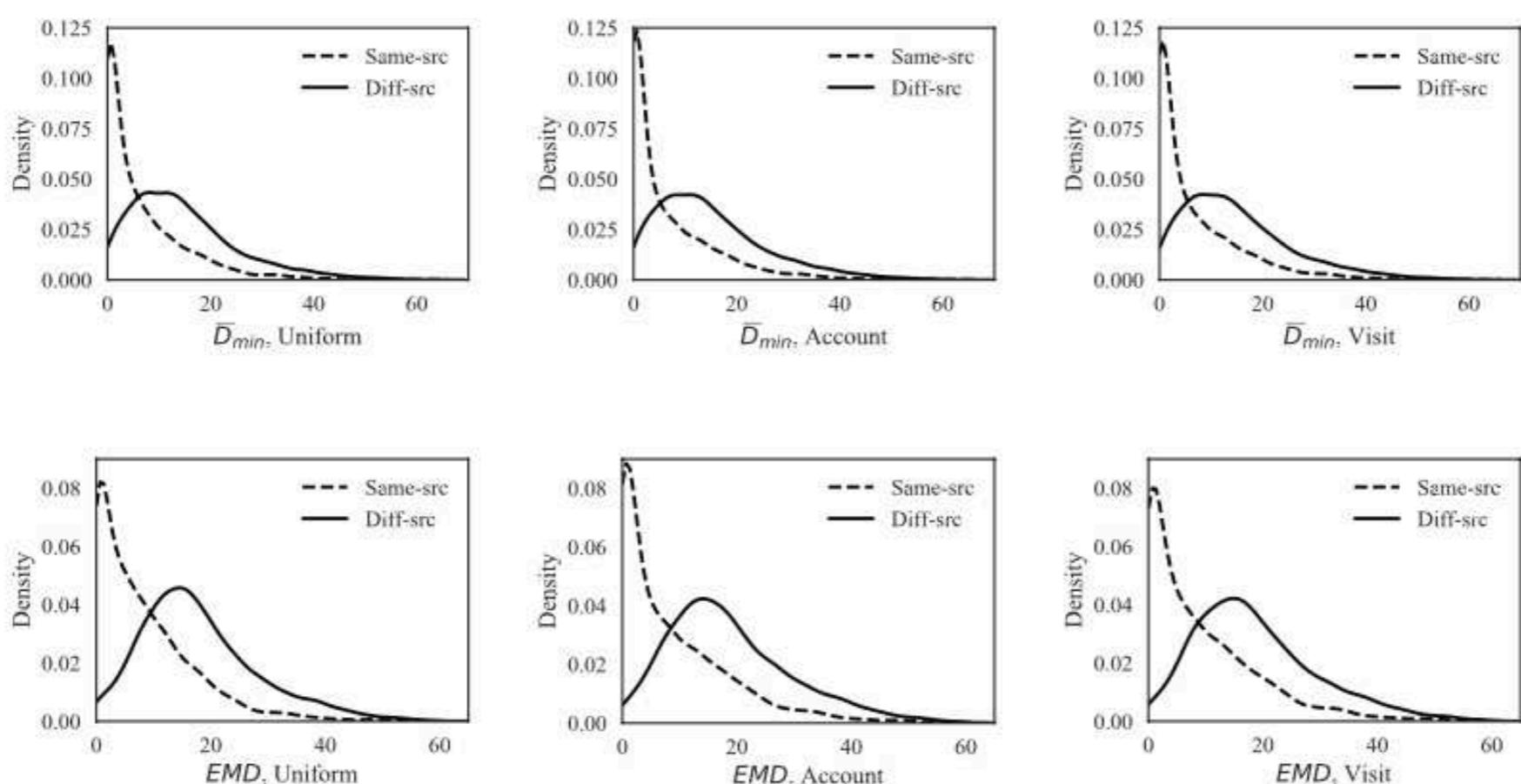


[Lichman, 2017]

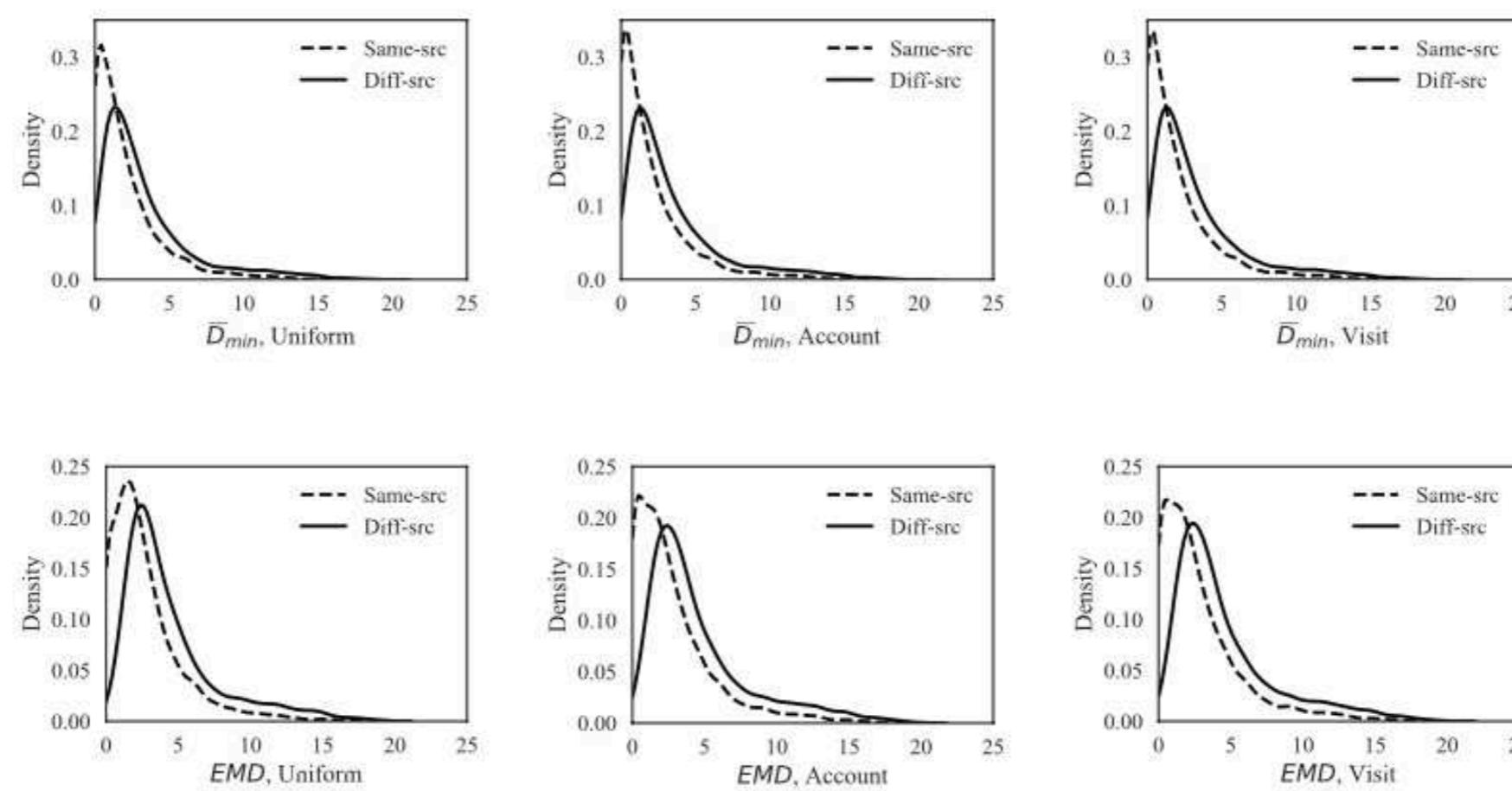
Number of
Points ↑

Weight ↓

OC



NY



Region	Weight	TP@1	FP@1	AUC
OC	0.80	0.340	0.026	0.787
	$\alpha(n_a)$	0.380	0.038	0.845
	$\alpha(n_a \gamma, \rho, \phi)$	0.375	0.037	0.817
NY	0.80	0.251	0.067	0.711
	$\alpha(n_a)$	0.285	0.089	0.768
	$\alpha(n_a \gamma, \rho, \phi)$	0.282	0.088	0.734

LR

Region	Δ	Weights	TP@1	FP@1	AUC
OC	\bar{D}_{min}	Uniform	0.628	0.202	0.768
	\bar{D}_{min}	Account	0.610	0.171	0.774
	\bar{D}_{min}	Visit	0.611	0.180	0.768
	EMD	Uniform	0.654	0.197	0.790
	EMD	Account	0.614	0.162	0.783
	EMD	Visit	0.602	0.169	0.774
NY	\bar{D}_{min}	Uniform	0.508	0.287	0.656
	\bar{D}_{min}	Account	0.494	0.254	0.666
	\bar{D}_{min}	Visit	0.493	0.257	0.663
	EMD	Uniform	0.530	0.253	0.686
	EMD	Account	0.511	0.235	0.685
	EMD	Visit	0.504	0.234	0.679

SLR

