

## PROBLEM STATEMENT

Consider a pair of user-generated event series

$$M = (A, B) = \{(t_j, m(t_j)) : j = 1, \dots, n\}$$

where  $t_j \in \mathbb{R}^+$  is the time and  $m(t_j) \in \{A, B\}$  is the type of the  $j^{th}$  event. We want to quantify the likelihood that the pair was generated by the same source.

## MEASURES OF ASSOCIATION

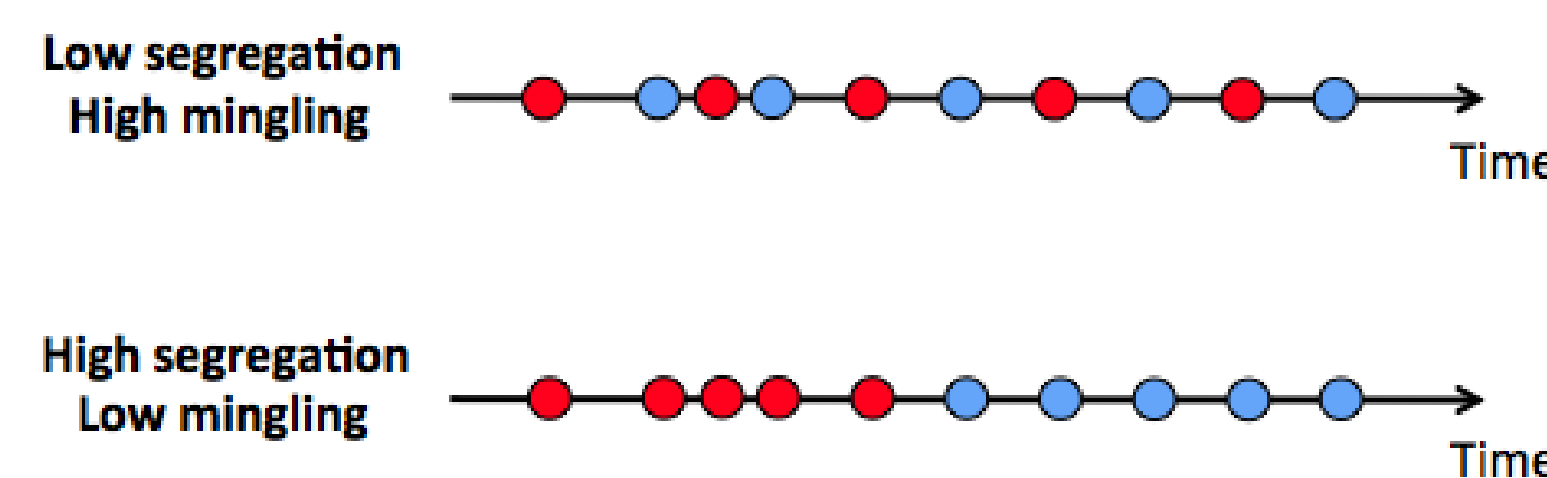
### Score Functions using Nearest Neighbors

- Coefficient of Segregation** [3]: function of the ratio of the probability that a reference point (i.e., a randomly selected event in  $(A, B)$ ) and its nearest neighbor have different marks to the same probability for independent marks.

$$S(A, B) = 1 - \frac{p_{AB} + p_{BA}}{p_{AP \cdot B} + p_{BP \cdot A}} \in [-1, 1]$$

- Mingling Index** [4]: mean fraction of points among the  $k$  nearest neighbors whose type is different than that of the reference point

$$\bar{M}_k(A, B) = \frac{1}{nk} \sum_{j=1}^n \sum_{\ell=1}^k \mathbb{I}[m(t_j) \neq m(z_\ell(t_j))] \in [0, 1]$$



### Score Functions using Inter-Event Times

Assume that  $n_B < n_A$  and fix series  $B$ . We then measure the time from each event in  $B$  to the closest event in series  $A$  in either direction

$$\tau_{BA} \equiv \{\tau_{BA,j} : j = 1, \dots, n_B\}$$

$$\text{where } \tau_{BA,j} = \min_{k \in \{1, \dots, n_A\}} |t_{b,j} - t_{a,k}|$$

- Mean inter-event time from  $B$  to  $A$

$$\bar{\tau}_{BA} = \frac{1}{n_B} \sum_{j=1}^{n_B} \tau_{BA,j} \in (0, \infty)$$

- Median inter-event time from  $B$  to  $A$

$$\text{med}(\tau_{BA}) \in (0, \infty)$$

## POPULATION-BASED APPROACH

### Given

- Pair of interest:  $(A^*, B^*)$
- Score function:  $\Delta$
- Sample of  $N$  pairs of event time series with known sources:  $M_i = (A_i, B_i)$  for  $i = 1, \dots, N$

### Method

- Two competing hypotheses:

$$H_s : (A^*, B^*) \text{ came from the same source}$$

$$H_d : (A^*, B^*) \text{ came from different sources}$$

- Use sample  $M_i = (A_i, B_i)$  for  $i = 1, \dots, N$  to estimate the *score-based likelihood ratio*

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

- Different interpretations of the denominator [1]

## RESAMPLING APPROACH

### Given

- Pair of interest:  $(A^*, B^*)$
- Score function:  $\Delta$

### Method

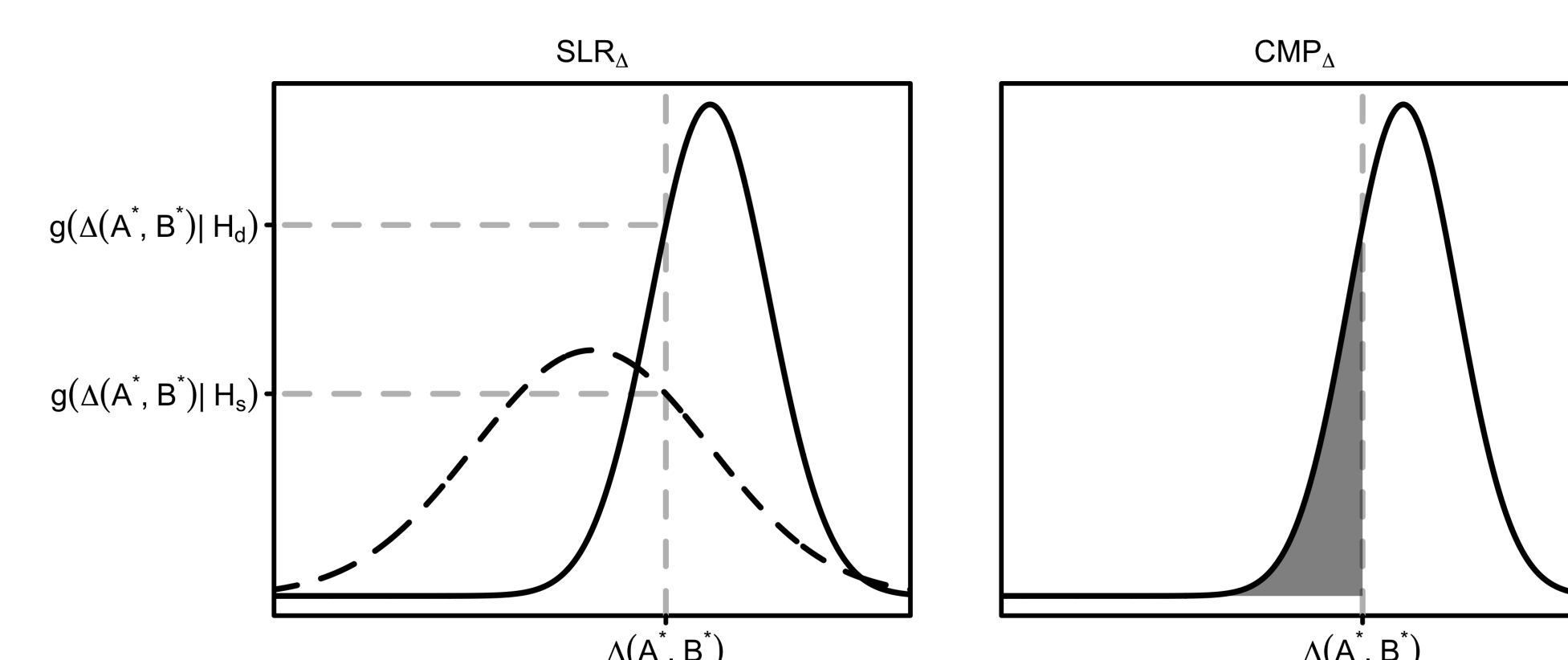
- Focus on the denominator of  $SLR_\Delta$
- Coincidental match probability**: probability that a different-source pair with **observed score**  $\Delta(A^*, B^*)$  exhibits association by chance

$$CMP_\Delta = Pr(\Delta(A, B) < \Delta(A^*, B^*)|H_d)$$

- Use resampling in time to simulate different-source pairs  $(A^{(i)}, B^{(i)})$  and estimate

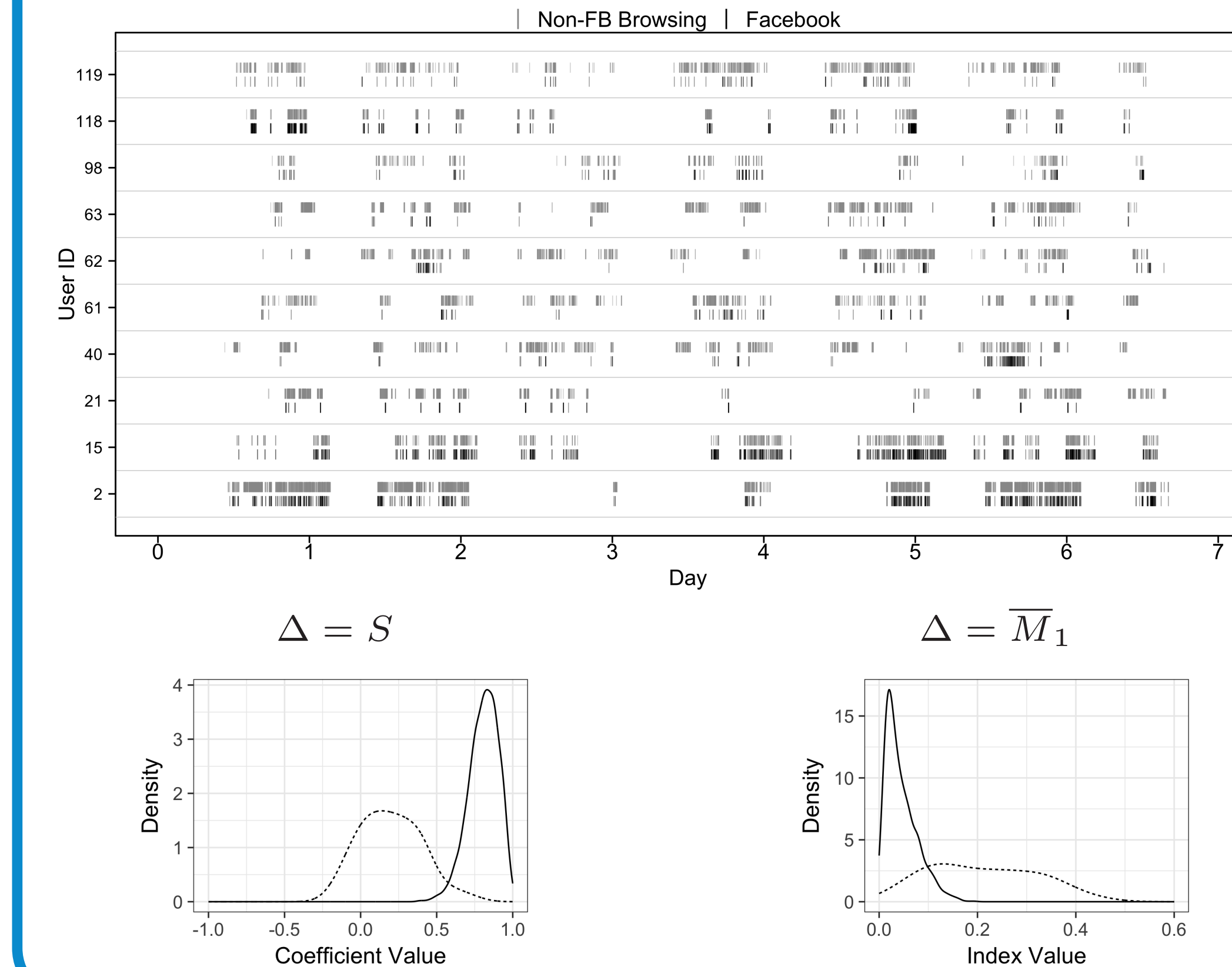
$$\widehat{CMP}_\Delta = \frac{1}{n_{sim}} \sum_{i=1}^{n_{sim}} \mathbb{I}[\Delta(A^{(i)}, B^{(i)}) < \Delta(A^*, B^*)]$$

## COMPARISON OF APPROACHES



## CASE STUDY

- Data from a 2013-2014 study at UCI that placed logging software on 124 students' computers that recorded all browser activity for one week [2]
- Event series created by dichotomizing browsing events to Facebook versus non-Facebook urls
- Only considered 55 students with at least 50 web browsing events of each type



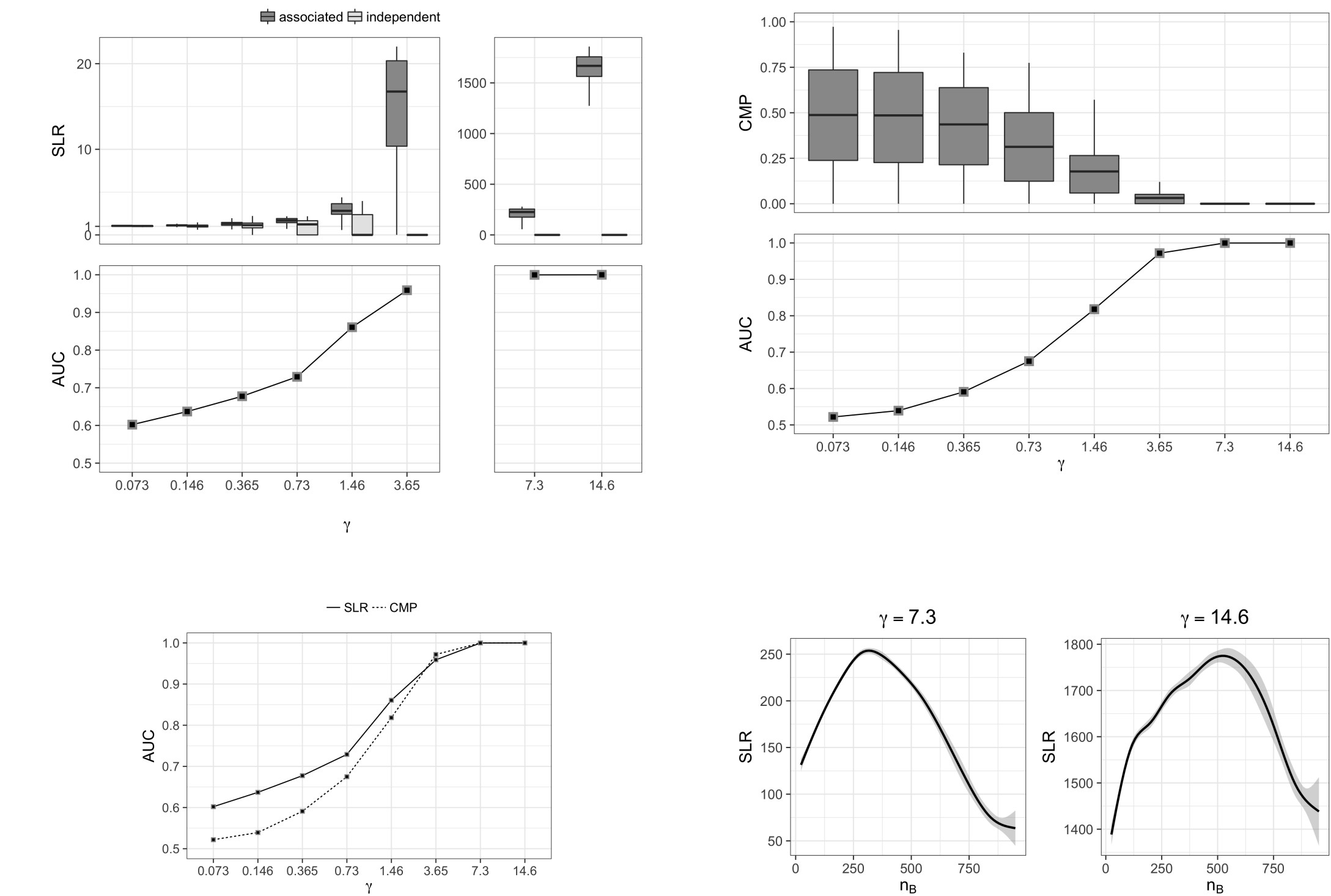
## SIMULATION

- Simulated equivalent of one week of data for pairs of processes with varying degrees of association
  - $A$ : Poisson process with intensity  $\lambda_A$
  - $B$ : independent Poisson process with intensity

$$\lambda_B = p\lambda_A \text{ with } p \in (0, 1)$$

or w.p.  $p$  add Gaussian noise to event in  $A$

- 10,000 independent & 10,000 associated pairs for each combination of  $\langle p, \sigma \rangle$ , distribution of  $\lambda_A$
- Most important factors in detecting associated pairs:
  - Number of events in series  $B$ :  $n_B$
  - Signal-to-noise ratio:  $\gamma = \bar{\tau}_{AA}/\sigma$



## CONCLUSIONS

- Resampling approach shows promise in situations where no reference data is available
- Population-based SLR is preferred, given
  - Better performance for weakly associated pairs
  - Similar performance for strongly associated pairs
  - Well-established in forensic investigation

## REFERENCES

- A. B. Hepler, C. P. Saunders, L. J. Davis, and J. Buscaglia, "Score-based likelihood ratios for handwriting evidence," *Forensic Science International*, vol. 219, no. 1, pp. 129–140, 2012.
- Y. Wang, M. Niya, G. Mark, S. Reich, and M. Warschauer, "Coming of age (digitally): an ecological view of social media use among college students," in *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*, pp. 571–582, 2015.
- E. C. Pielou, *Mathematical Ecology*. John Wiley & Sons, Inc., 1977.
- J. Illian, A. Penttinen, H. Stoyan, and D. Stoyan, *Statistical Analysis and Modelling of Spatial Point Patterns*. England: John Wiley & Sons Ltd, 2008.

## ACKNOWLEDGEMENT

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.

## CONTACT INFORMATION

**Web** [www.ics.uci.edu/~galbraic](http://www.ics.uci.edu/~galbraic)  
**Email** [galbraic@uci.edu](mailto:galbraic@uci.edu)