



# Detecting Users from Website Sessions: A Simulation Study

---

Corné de Ruijt  
Sandjai Bhulai

DATA ANALYTICS 2020

Vrije Universiteit Amsterdam  
Contact: [c.a.m.de.ruijt@vu.nl](mailto:c.a.m.de.ruijt@vu.nl)

# Introduction

- PhD student at the Vrije universiteit Amsterdam
- Research is mainly about job recommender systems and recruitment websites



- In this study we wish to address the problem of recognizing unique users from session data.
- Recognizing users from session data is typically an interplay between cookies and log-in.
- 20% of the internet users may delete their cookie at least once a week <sup>1</sup>.
- We study this problem by: 1) proposing a click simulation model, 2) studying how effective (H)DBSCAN\*-type algorithms are to this problem.

---

<sup>1</sup>Dasgupta et al. Overcoming browser cookie churn with clustering. Proceedings of the fifth ACM international conference on Websearch and data mining, pages 83–92. ACM, 2012

- Introduction
- Literature
- Click simulation model
- MS-(H)DBSCAN\* algorithms for session clustering
- Results and conclusion

- Specific type of entity/identity resolution problem <sup>2</sup>
- ICDM<sup>3</sup> and CIKM <sup>4</sup> cross-device matching competitions
- Ambiguity in the literature:
  - Problem is addressed under different names (user stitching, visitor stitching, automatic identity linkage)
  - Problem is addresses from single and multiple website perspective

---

<sup>2</sup>Di et al. "node2bits: Compact Time-and Attribute-aware Node Representations", 2019

<sup>3</sup>ICDM 2015: Drawbridge Cross-Device Connections, 2015.

(<https://www.kaggle.com/c/icdm-2015-drawbridge-cross-device-connections>)

<sup>4</sup>CIKM Cup 2016 Track 1: Cross-Device Entity Linking Challenge, 2016.

(<https://competitions.codalab.org/competitions/11171>)

# Research objective

We assume:

- Single website (a search engine)
- Homogeneous queries, but with different utility functions

Interested in:

- What simulation model should we use to obtain realistic simulated datasets, without making the simulation model overly complex?
- How effective are (H)DBSCAN\*-type algorithms on clustering sessions compared to using cookies, in terms of multiple website statistics?
- Sensitivity of (H)DBSCAN\*-type algorithms on the underlying dataset.

## Why simulation?

- Have an actual ground truth (problem of unary classification).
- Sensitivity of clustering algorithm under different simulation settings.
- Quite some literature on models explaining website behavior on search engines <sup>5</sup>.
- Most public data sets originate from large search engines/large advertisers. How generalizable are these for other websites?

---

<sup>5</sup>Chuklin et al. "Click models for web search." Synthesis lectures on information concepts, retrieval, and services 7.3 (2015): 1-115.

- Introduction
- Literature
- Click simulation model
- MS-(H)DBSCAN\* algorithms for session clustering
- Results and conclusion



- **Users' item preferences:** Fleder-Hosanagar model<sup>6</sup>.
- **Interaction with search engine:** Simplified Dynamic Bayesian Network click model<sup>7</sup>.
- **Cookie churn:** 1) modeling cookie lifetime, 2) allowing users to use multiple devices (with different cookie-IDs).

---

<sup>6</sup>Fleder et al. "Blockbuster culture's next rise or fall: The impact of recommender systems on sales diversity. *Management science*", 55(5):697-712, 2009

<sup>7</sup>Chapelle et al. "A dynamic bayesian network click model for web search ranking." *Proceedings of the 18th international conference on World wide web*. 2009.

# Simulation summary

---

---

```
1 Draw items  $\mathcal{V}$  and users  $\mathcal{U}$  and compute each user's attraction
  and satisfaction over all items using Fleder's model;
2 for  $u \in \mathcal{U}$  do
3   Draw initial device, cookie lifetime;
4   Set  $t$  and  $t_{\text{start}}$  to arrival time of  $u$ , draw user lifetime  $T_u^{\text{user}}$ ;
5    $i \leftarrow 1$ ;
6   while  $t \leq t_{\text{start}} + T_u^{\text{user}}$  do
7     Simulate clicks for query session  $i$  according to SDBN,
      items are presented according to a popularity
      recommender;
8     Draw  $T_{u,i}^{\text{abs}}$  and update  $t$ ;
9     Update device and cookies;
10     $i \leftarrow i + 1$ 
11  end
12 end
```

---

# Simulation, output

User ( $u$ )	Query-session ( $i$ )	Cookie ( $o$ )	Device ( $d$ )	List position	Item ( $v$ )	$t$	Click/skip
1	1	1	1	1	12	12	1
1	1	1	1	2	13	12	0
1	1	1	1	3	44	12	0
...	...	...	...	...	...	...	...
1	1	1	1	10	84	12	0
2	1	1	2	1	5	12	0
...	...	...	...	...	...	...	...
312	32	2	2	1	871	16	0

- Introduction
- Literature
- Click simulation model
- MS-(H)DBSCAN\* algorithms for session clustering
- Results and conclusion

# Results

Model	Dataset	APE unique users	KL-div. session count	KL-div conversion	ARI	New user accuracy
MS-DBSCAN	train	15	<b>0.55</b>	0.13	0.0012	<b>0.56</b>
MS-DBSCAN <sub><math>\rho</math></sub>	train	77	0.74	<b>0.092</b>	<b>0.14</b>	0.5
DBSCAN-RAND	train	<b>0.011</b>	1	0.096	0.0002	0.42
MS-HDBSCAN* <sup>+</sup>	train	10	0.75	0.15	0.00079	0.52
MS-HDBSCAN* <sup>-</sup>	train	10	0.75	0.15	0.00079	0.52
MS-HDBSCAN* <sup>+</sup> <sub><math>\rho</math></sub>	train	<b>0.011</b>	0.9	0.11	0.092	0.46
MS-HDBSCAN* <sup>-</sup> <sub><math>\rho</math></sub>	train	<b>0.011</b>	0.9	0.11	0.1	0.46
OBS	<i>train</i>	15	<i>0.017</i>	<i>0.0032</i>	<i>0.91</i>	<i>0.95</i>
MS-DBSCAN	valid	60	<b>0.11</b>	<b>0.0026</b>	<b>0.0022</b>	<b>0.56</b>
MS-DBSCAN <sub><math>\rho</math></sub>	valid	<b>6.8</b>	1.4	0.13	0.0015	0.4
DBSCAN-RAND	valid	40	0.32	0.015	0.00046	0.5
MS-HDBSCAN* <sup>+</sup>	valid	53	0.16	0.0042	0.002	0.55
MS-HDBSCAN* <sup>-</sup>	valid	53	0.16	0.0042	0.002	0.55
MS-HDBSCAN* <sup>+</sup> <sub><math>\rho</math></sub>	valid	7.2	1.4	0.13	0.0015	0.4
MS-HDBSCAN* <sup>-</sup> <sub><math>\rho</math></sub>	valid	7.2	1.4	0.13	0.0015	0.4
OBS	<i>valid</i>	51	<i>0.1</i>	<i>0.0076</i>	<i>0.91</i>	<i>0.95</i>

# Conclusion

- We presented a click simulation model with cookie censoring and illustrated its usage and advantages on the problem of uncovering users from their web sessions.
- Usage of cookie-IDS as clusters outperforms MS-(H)DBSCAN\*-type algorithms on the homogeneous query case;
- Of the MS-(H)DBSCAN\*-type algorithms, MS-DBSCAN significantly outperformed other MS-(H)DBSCAN\* type algorithms in terms of ARI, including DBSCAN-RAND;
- Increasing ARI or new user accuracy also seems to increase the unique user average percentage error;
- Increasing the signal (e.g., by increasing the number of clicks) improved the results, but slightly.

## Further research

- What happens if we move to a multi-query case?
- Adjust clustering approaches using supervised methods s.t. they are less prone to overfitting.
- Holistic research on cookie (censoring) models.