

Birds of a Feather Flock Together: How Set Bias Helps to Deanonimize You via Revealed Intersection Sizes

(USENIX Security Symposium, August 2022)

Xiaojie Guo¹, **Ye Han**¹, Zheli Liu¹, Ding Wang¹, Yan Jia¹, Jin Li²

¹Nankai University, ²Guangzhou University



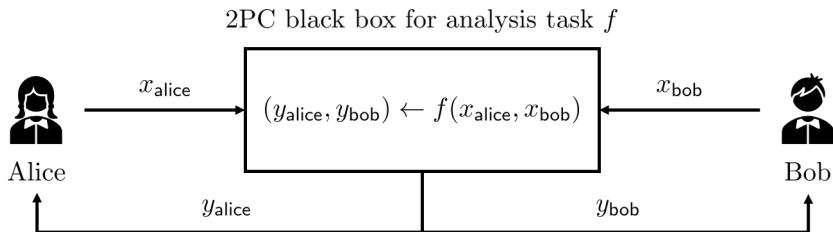
Introduction

Set Membership Inference

Evaluation

Discussion & Conclusion

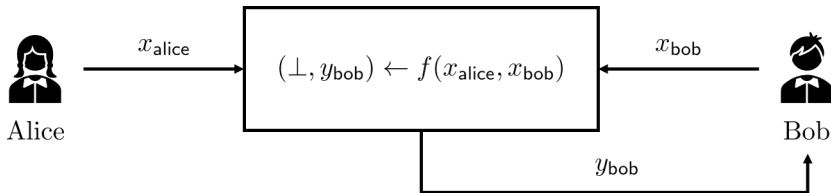
SECURE TWO-PARTY COMPUTATION (2PC)



INTERSECTION-RELATED DATA ANALYSIS FROM 2PC

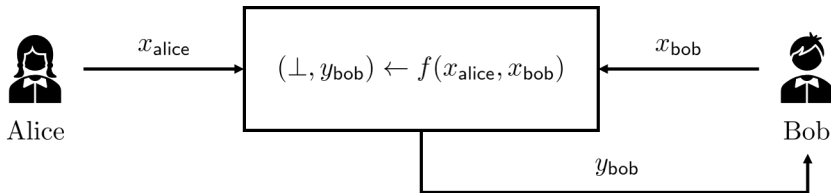
► COVID-19 contact tracing

- Alice: Health authority
- Bob: Client
- x_{alice} : A set Y of (tokens of) infected patients
- x_{bob} : A set X of (tokens of) individuals in contact with
- y_{bob} : $|X \cap Y|$



INTERSECTION-RELATED DATA ANALYSIS FROM 2PC

- ▶ Measurement of ad conversion revenue/lift
 - ▶ Alice: Publisher
 - ▶ Bob: Advertiser
 - ▶ x_{alice} : A table indexed by a set Y of (tokens of) individuals that click/view the ad
 - ▶ x_{bob} : A set X of (tokens of) converted individuals
 - ▶ y_{bob} : Ad conversion revenue/lift from intersecting converted individuals, and $|X \cap Y|$



WHY HIDE INTERSECTION?

- ▶ From Bob's view

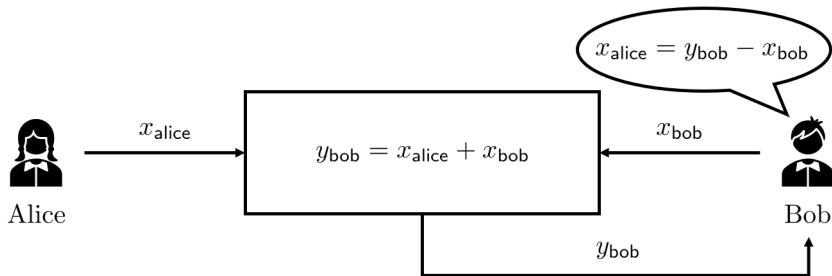
Intersection \Rightarrow Each token's membership regarding Alice's set
 \Rightarrow Each individual's relationship with Alice (*)

- ▶ (*): Linking a token to an individual is possible
 - ▶ COVID-19 contact tracing: Physical contact & token exchange logs¹²
 - ▶ Measurement of ad conversion revenue/lift: Sensitive tokens (e.g., email addresses, IMEI numbers)

¹Yaron Gvili. [Security analysis of the COVID-19 contact tracing specifications by Apple Inc. and Google Inc. \(IACR\).](#)

²<https://www.wired.com/story/apple-google-contact-tracing-strengths-weaknesses/>.

- ▶ 2PC does **NOT** protect what can be deduced from one party's input and output!



- ▶ Inference attacks in intersection-related analysis tasks
 - ▶ COVID-19 contact tracing
 - ▶ Measurement of ad conversion revenue
 - ▶ Measurement of ad conversion lift
- ▶ Our observations
 - ▶ Existing 2PC protocols for these tasks **reveal intersection size**
 - ▶ These tasks need to be performed regularly
- ▶ More severe leakage if there is a non-weak **set bias**
 - ▶ Set bias: Alice's set tends to include individuals with certain features

Introduction

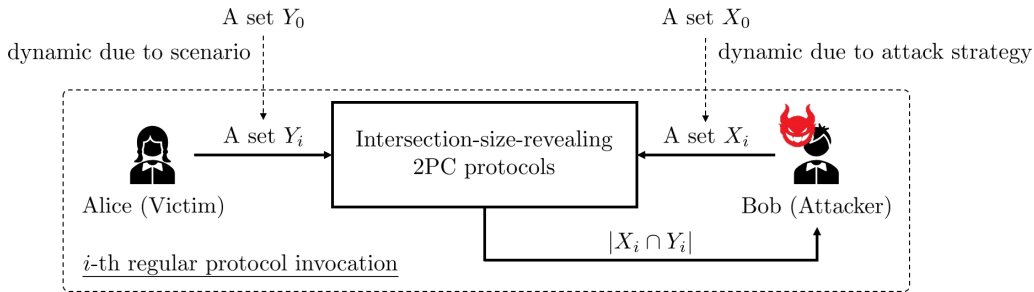
Set Membership Inference

Evaluation

Discussion & Conclusion

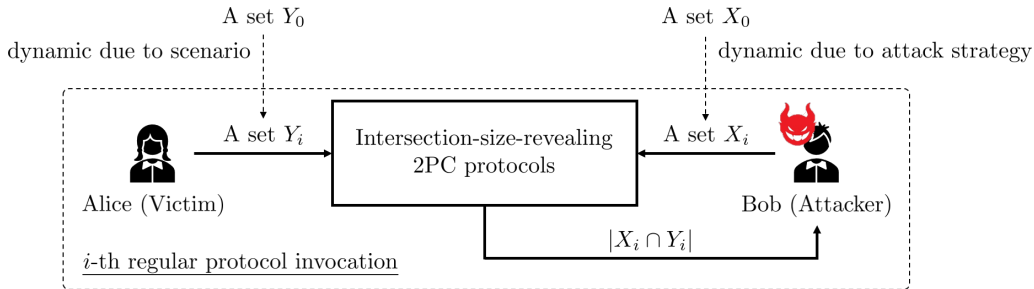
SET MEMBERSHIP INFERENCE: PROBLEM DEFINITION

- ▶ Alice is the victim with a dynamic set $\{Y_0, \dots, Y_i, \dots\}$
- ▶ Bob is the attacker with a fixed set X of **target individuals**
- ▶ Bob can choose its set X_i in each protocol invocation
- ▶ Bob aims to determine whether a target individual has been in $Y = \cup_i Y_i$



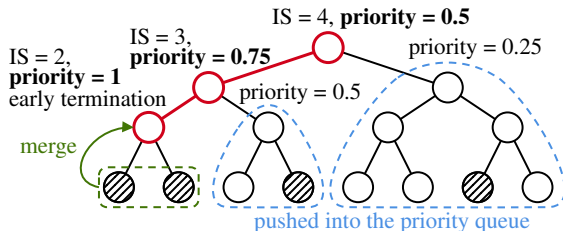
A TOY INFERENCE ATTACK

- ▶ Each X_i contains only **one** target individual
- ▶ Learn its set membership from $|X_i \cap Y_i|$
- ▶ Require $|X|$ protocol invocations



OUR BASELINE ATTACK

- ▶ Choose X_i with binary-tree-based strategy
 - ▶ Set X as the root, and randomly divide a node into two equal-size child nodes
 - ▶ Visit nodes via **priority-based** depth-first search (DFS), and set X_i as the current node
 - ▶ $\text{priority} = \text{intersection size (IS)} / \# \text{ target individuals in the node}$
 - ▶ $\text{IS in right child} = \text{IS in parent} - \text{IS in left child}$



OUR FEATURE-AWARE ATTACK: USING SET BIAS

- ▶ A stronger attacker with some features regarding set bias
- ▶ Same as baseline attack, except that a node is divided using **feature-based clustering**
- ▶ Intuition
 - ▶ Clustering can put target individuals with similar features in the same sub-tree
 - ▶ A non-weak set bias \Rightarrow many member individuals are with similar features
- ▶ Implement clustering with **k-means**

- ▶ No set membership change of target individuals \Rightarrow Perfectly correct result
- ▶ Otherwise, there may be **false positives** and **false negatives**

CONTENTS

Introduction

Set Membership Inference

Evaluation

Discussion & Conclusion

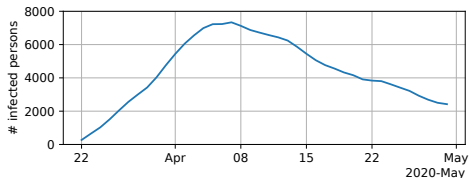
- ▶ Public data sources
 - ▶ COVID-19 contact tracing: COVID-19 dataset of tested individuals in Israel
 - ▶ Measurement of ad conversion revenue: Taobao's dataset of ad display/click records
 - ▶ Measurement of ad conversion lift: Tencent's dataset of ad display records
- ▶ Frequency of protocol invocation
 - ▶ COVID-19 contact tracing: 5 / day
 - ▶ Measurement of ad conversion revenue: 1 / hour
 - ▶ Measurement of ad conversion lift: 1 / day

- Set bias (higher mutual information \Rightarrow stronger set bias regarding a feature)

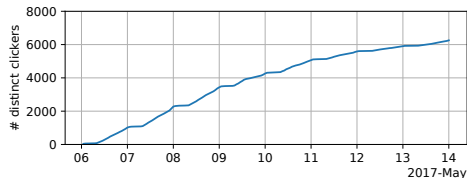
Scenario	feature name (mutual information)
COVID-19 contact tracing	fever (0.0168), cough (0.0099), gender (0.0004)
Measurement of ad conversion revenue	age (0.0010), gender (0.0007), shopping_level (0.0002), work (0.0002), consumption_ability (0.0001), city_level (0.0001)
Measurement of ad conversion lift	marriage_status (0.0013), education (0.0012), consumption_ability (0.0012), age (0.0009), work (0.0005), gender (0.0001)

SUMMARIES OF THE VICTIM'S SET

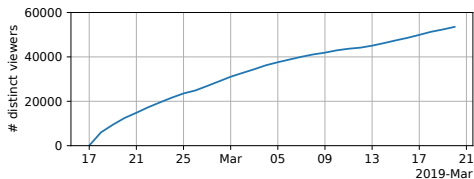
► Set size change



(a) COVID-19 contact tracing.



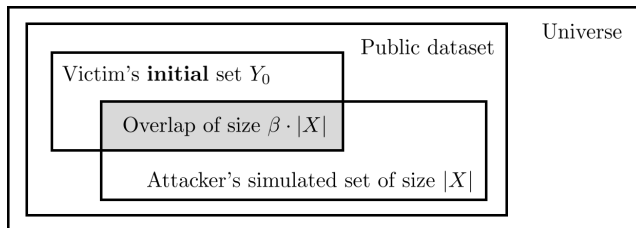
(b) Measurement of ad conversion revenue.



(c) Measurement of ad conversion lift.

SIMULATION OF THE ATTACKER'S SET

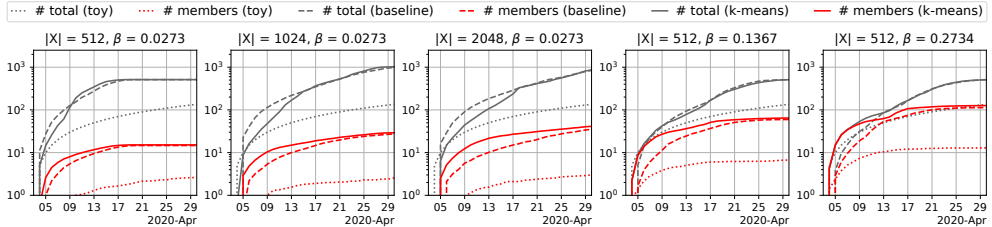
- ▶ Simulation parameters
 - ▶ Size $|X|$ of the attacker's set
 - ▶ Ratio β of # target individuals in the victim's **initial** set Y_0 to $|X|$



- ▶ Feature-aware attacker can only use easy-to-collect features

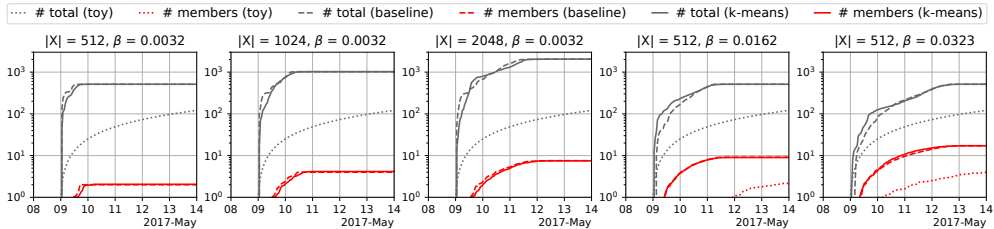
EXPERIMENTAL RESULTS

► Set membership leakage in COVID-19 contact tracing



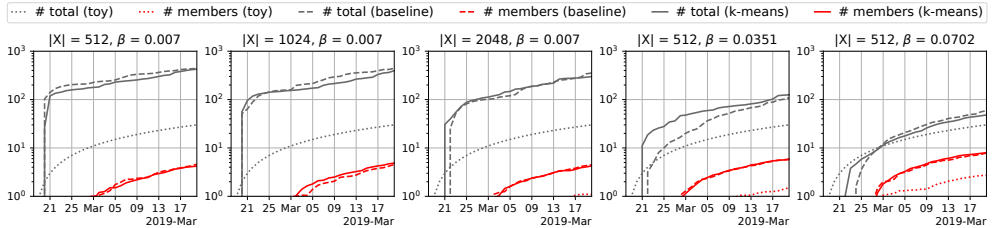
EXPERIMENTAL RESULTS

- Set membership leakage in measurement of ad conversion revenue



EXPERIMENTAL RESULTS

► Set membership leakage in measurement of ad conversion lift



CONTENTS

Introduction

Set Membership Inference

Evaluation

Discussion & Conclusion

- ▶ Set membership leakage does exist in the three scenarios
- ▶ COVID-19 contact tracing
 - ▶ Set membership regarding health authority \Rightarrow Whether a target token belongs to a COVID-19 patient
 - ▶ Can be combined with known linkage attacks in non-2PC settings \Rightarrow Risk of patient deanonymization
- ▶ Measurement of ad conversion revenue
 - ▶ Set membership regarding publisher \Rightarrow Whether a target individual has clicked the ad \Rightarrow Personal interest
- ▶ Measurement of ad conversion lift
 - ▶ Set membership regarding publisher \Rightarrow Whether a target individual has or would have seen the ad \Rightarrow Personal interest

- ▶ Limiting the number of 2PC protocol invocations
 - ▶ Auditing intersection size
 - ▶ Auditing the size of the attacker's set
 - ▶ Applying differential privacy
-
- ▶ But there are also some challenges when using these defenses ...

- ▶ Set membership inference problem in intersection-size-revealing 2PC protocols
- ▶ A baseline attack and a feature-aware attack, where the latter outperforms the former given a non-weak set bias
- ▶ Evaluation in three scenarios with public datasets

Thank You

Contact xiaojie.guo@mail.nankai.edu.cn for any questions