# Soteria

security made **usable**

i290T Open Data          Ashley   Evie   Shreyas   Victor

# OBJECTIVE

To provide **usable analysis** of **security data** to the

open community

# DATA

- Source: <u>VERIS Community DB</u>

- 3,084 breach incidents

- 321 breach features (cause &effect)

  - **attack.vectors** (270+ booleans, hacking, malware, …)

  - **attack.effects** (victim.count, victim.revenueloss, …)

# DATA CLEANING

- **Handle Missing Values** : NaN, Unknown

- **Recode Variables**

  - Eg: **timeline.discovery** -> **discovery.daycount**

- **Group Variables** (Using **CENSUS** API)

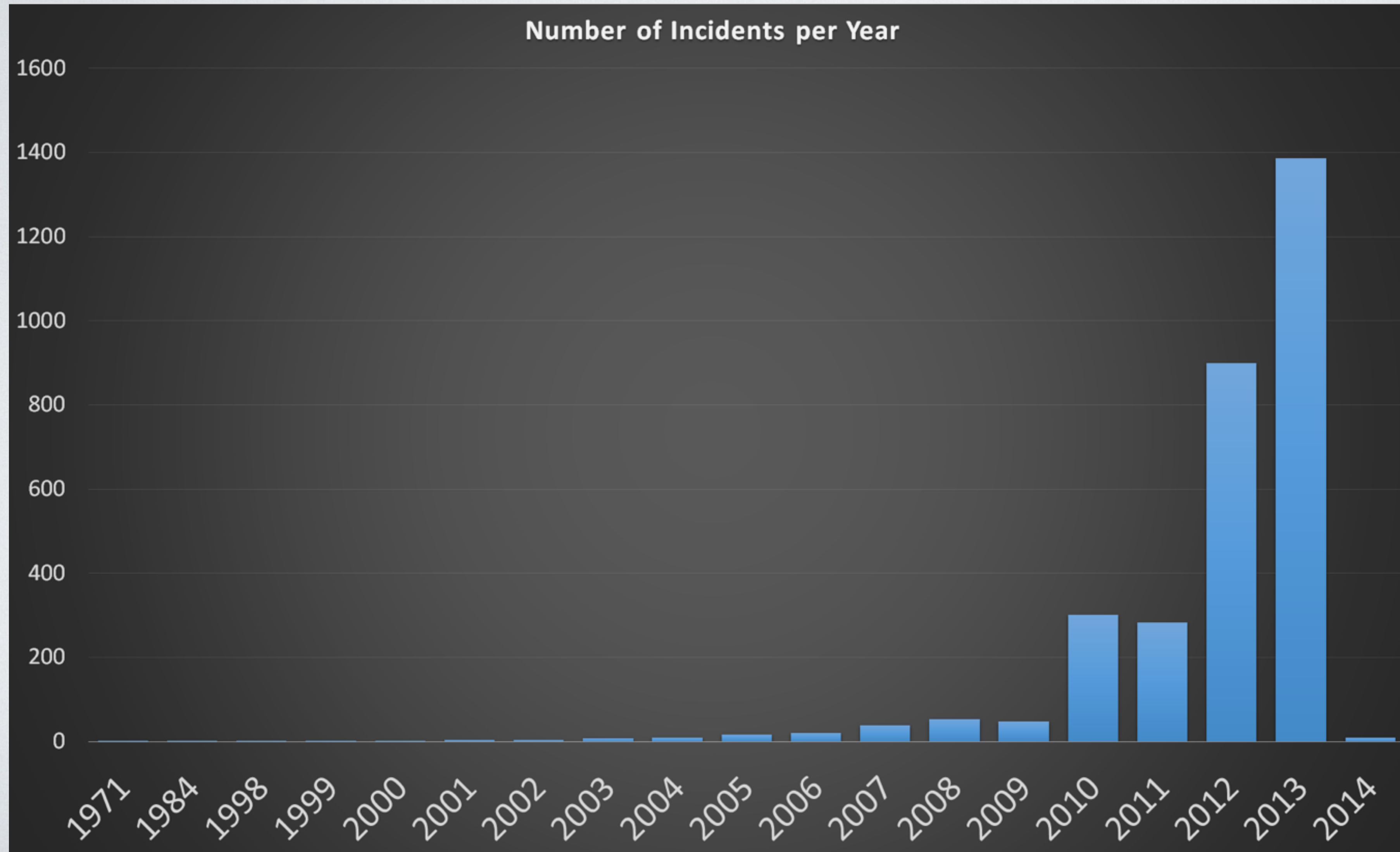  - Eg: **Victim.IndustryName** -> **Victim.Industry.Category**

# DATA EXPLORATION

Soteria

(aggregate statistics)

(one-variable-at-a-time)
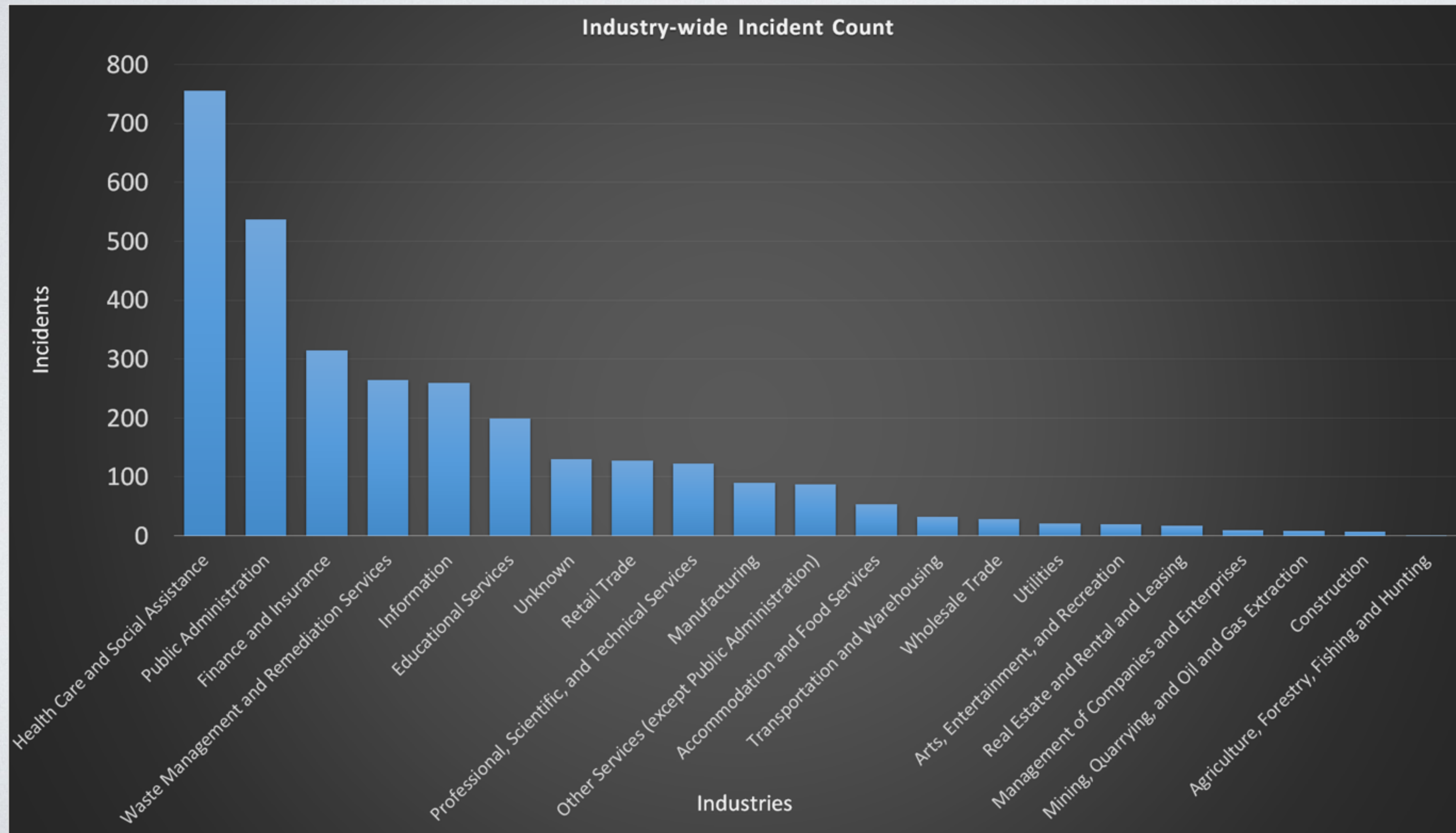
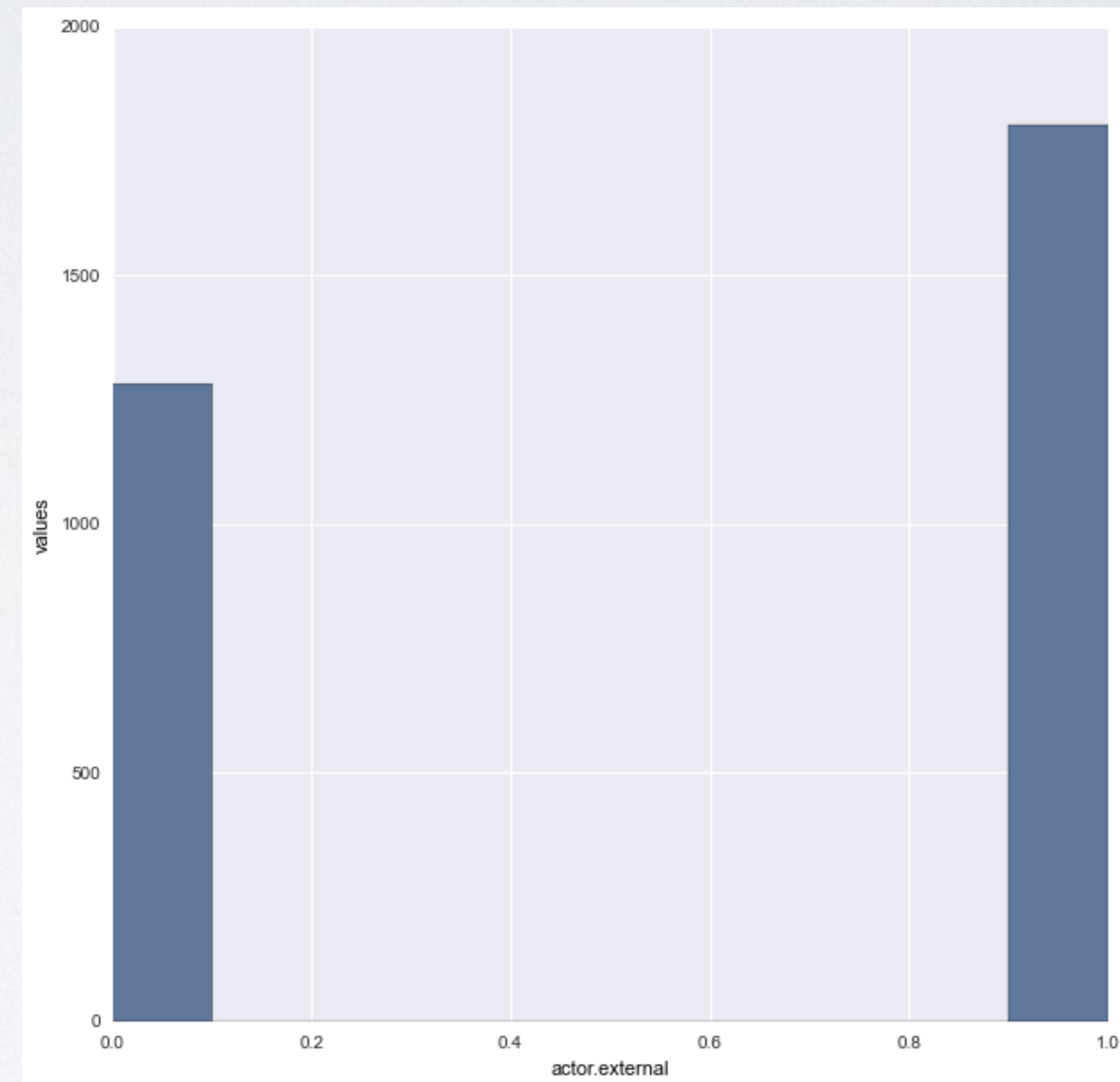(2-variables-at-a-time)

(all-key-variables-at-a-time)

# ONE-VARIABLE

Soteria

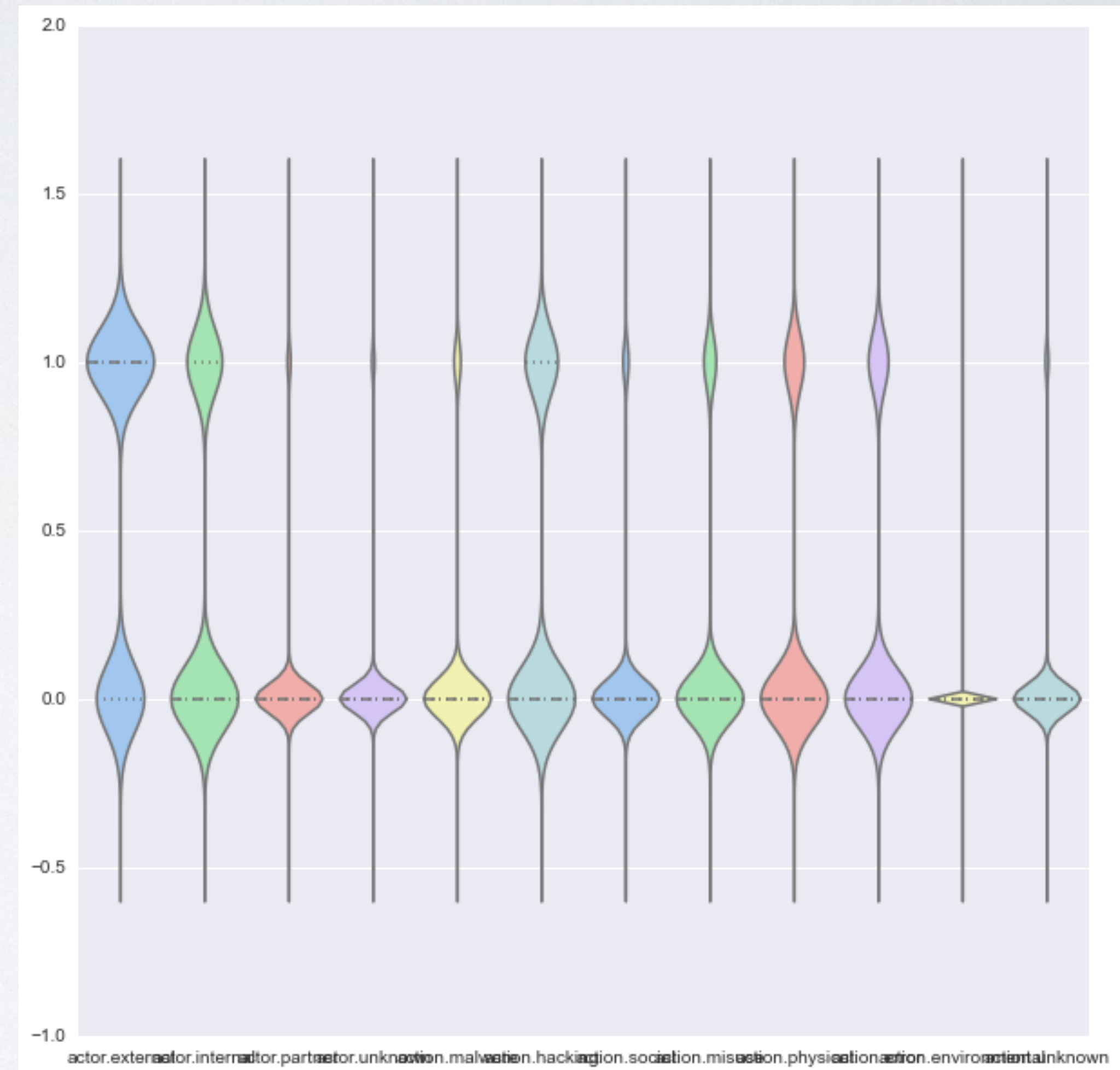# KEY-VARIABLES

Soteria

## ViolinPlot

# DATA SCIENCE

- **cause** (booleans) & **effect** (victim) variables

- usable data exploration through filtering

- find pairs of similar incidents

# USABLE FILTERING

Soteria

**front-end demo**

# SIMILAR INCIDENTS

Soteria

attack_similarity.py

**(MapReduce)**

# CONCLUSIONS (PRELIMINARY)

Soteria

- a **lot** of similarity

- Our inferences:

  - attacks are **reused**

  - breach **vocabulary** is still not **deep enough**

# QUESTIONS ...