

Disinformation Defense

AI Inference

Karel Baloun

Keng-Shao (Ken) Chang

Matt Holmes

Master of Information and Cybersecurity (MICS)

University of California - Berkeley

Fall 2019 | Capstone | School of Information

Table of Contents

1	Introduction.....	3
2	Threat Modeling.....	6
3	Design	10
3.1	Attribute Distribution.....	10
3.2	Data Extraction and Generation (Bootstrapping).....	10
3.3	Machine Learning (ML) Model	11
3.3.1	Hypothesis.....	11
3.3.2	Supervised Learning with Classification.....	11
3.3.3	Vote Preference Inference (Prediction).....	11
3.4	Domain Generalization	11
3.5	Data Anonymization	12
3.6	Vote Influence.....	12
3.7	Disinformation Engine.....	12
4	Attack: Machine Learning	13
5	Defense: Privacy Engineering.....	14
6	Recommendation	17
7	Conclusion	18

Disinformation Defense

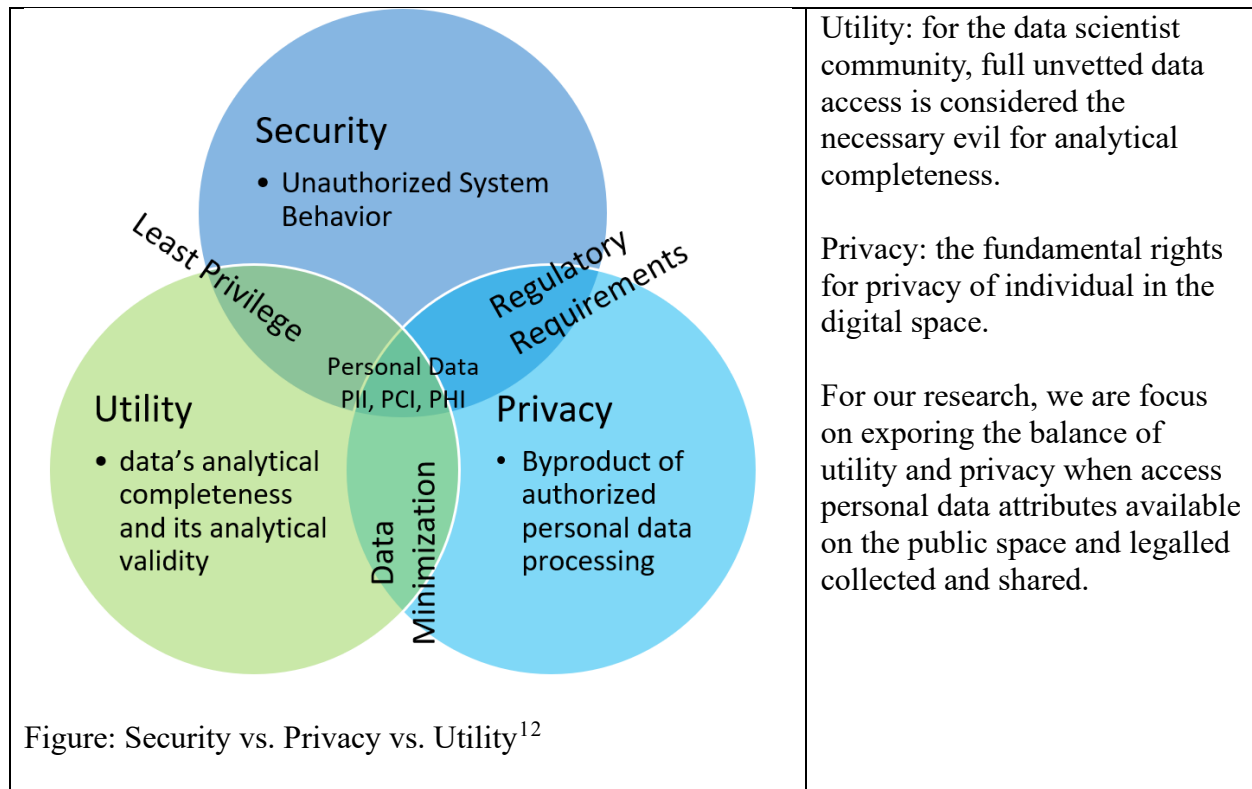
AI Inference

1 Introduction

Internet has no physical boundary on how personal data can be collected across the international community. The personal data collected and generated through the digital transaction of personal life can be considered as an active and passive type of data collection channels. Active user data collection requires a user to spend time entering information and requires that a user feel comfortable with providing information over the Internet. Passive user data can be obtained and unintentionally left behind by the users of the internet and digital devices without explicit consent on how it is being collected, processed and stored.

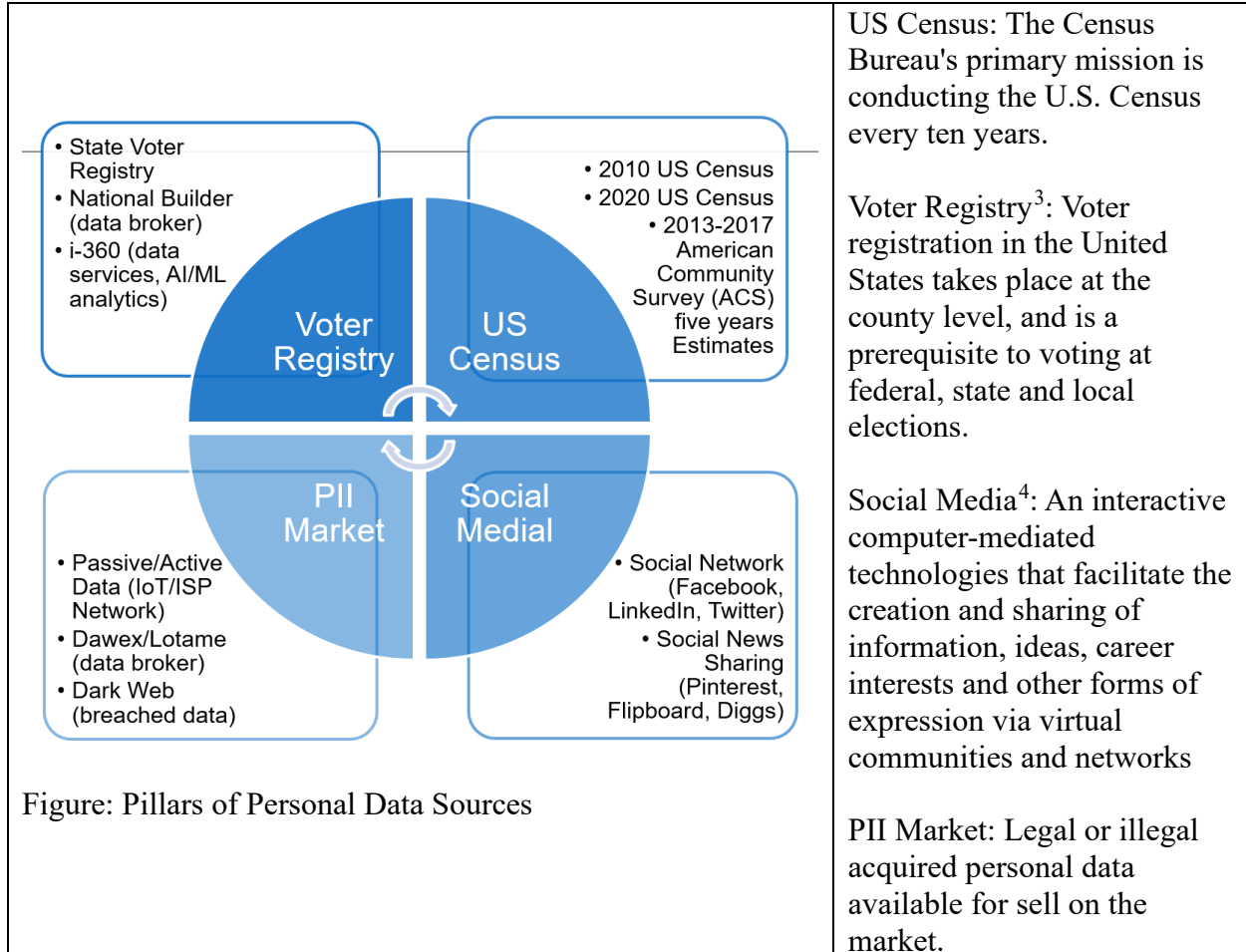
With the advance of Machine Learning and Artificial Intelligence under high-performance computing platform using the vast amount of personal data, the adversaries sponsored by nation-states can perform data linkage and aggregation to analyze and infer the targeted person's personal and professional activities in the critical infrastructure sector which might pose a threat to national security. Furthermore, based on the leaked personal data, the adversary can spread disinformation or misinformation on the targeted person held a critical position at the private or public sector which might eventually lead to the disruption of business and government functions. The Russian government interfered in the 2016 U.S. presidential election with the goal of harming the campaign of Hillary Clinton, boosting the candidacy of Donald Trump, and increasing political and social discord in the United States.

	Security: the traditional information security focus on confidentiality, integrity and availability.
--	--



¹ USENIX Enigma 2019 - Privacy Engineering: Not Just for Privacy Engineers

² NIST 8062



³ https://en.wikipedia.org/wiki/Voter_registration#United_States

⁴ https://en.wikipedia.org/wiki/Social_media

2 Threat Modeling

Collective data breaches leaking personally identifiable information (PII), Payment Card Information (PCI), and Protected Health Information (PHI) in the United States and other developed countries like Taiwan and Singapore. The personal data can be sensitive or not sensitive collected via active and non-active channels. Non-sensitive PII is information that can be transmitted in an unencrypted form without resulting in harm to the individual under normal circumstances. The normal circumstances imply data being collected and processed in isolation without linking other sources of personal data. Non-sensitive PII can be easily gathered from public records, phone books, corporate directories, and websites. Sensitive PII, PCI and PHI data are information which, when disclosed, could result in harm to the individual whose privacy has been breached. Sensitive personal data includes biometric information, medical information, personally identifiable financial information (PIFI⁵) and unique identifiers such as passport used for international travel or Social Security numbers.

The private sector has been under repeated offending attacks in the result of massive data breach events. In October 2017, Yahoo⁶ disclosed 3 billion accounts breached in 2013 incident; up from the original estimate of 1 billion accounts. Specific details of material taken include names, email addresses, telephone numbers, encrypted or unencrypted security questions, and answers, dates of birth, and hashed passwords. In November 2017, Uber⁷ shocked consumers when it admitted that it failed to notify victims for over a year after paying \$100,000 to hackers who had

⁵ <https://www.techopedia.com/definition/14222/personally-identifiable-financial-information-pifi>

⁶ https://en.wikipedia.org/wiki/Yahoo!_data_breaches

⁷ <https://www.theguardian.com/technology/2017/nov/21/uber-data-hack-cyber-attack>

stolen data on 57 million users and drivers. In 2017, one of the largest CRAs, Equifax⁸ Inc. (“Equifax”) announced that it had suffered a data breach that involved the PII of over 145 million Americans in the result of 230 million revenue loss quarter over quarter. In November 2018, Marriott first revealed it had suffered a massive data breach affecting the records of up to 500 million customers. Information accessed included payment information, names, mailing addresses, phone numbers, email addresses, and passport numbers. In February 2015, Anthem⁹, Inc. disclosed that criminal hackers had broken into its servers and potentially stolen over 78.8.5 million records that contain personally identifiable information. The compromised information contained names, birthdays, medical IDs, social security numbers, street addresses, email addresses and employment information, including income data.

The public sector can’t be exempted as the victim of the data breach. In June 2015, the United States Office of Personnel Management (OPM) announced that it had been the target of a data breach targeting the records of as many as twenty-one million people. This includes records of people who had undergone background checks, but who were not necessarily current or former government employees. Same malware, Sakula was used to attack other US companies which led to the conclusion of Chinese adversaries behind the attack.

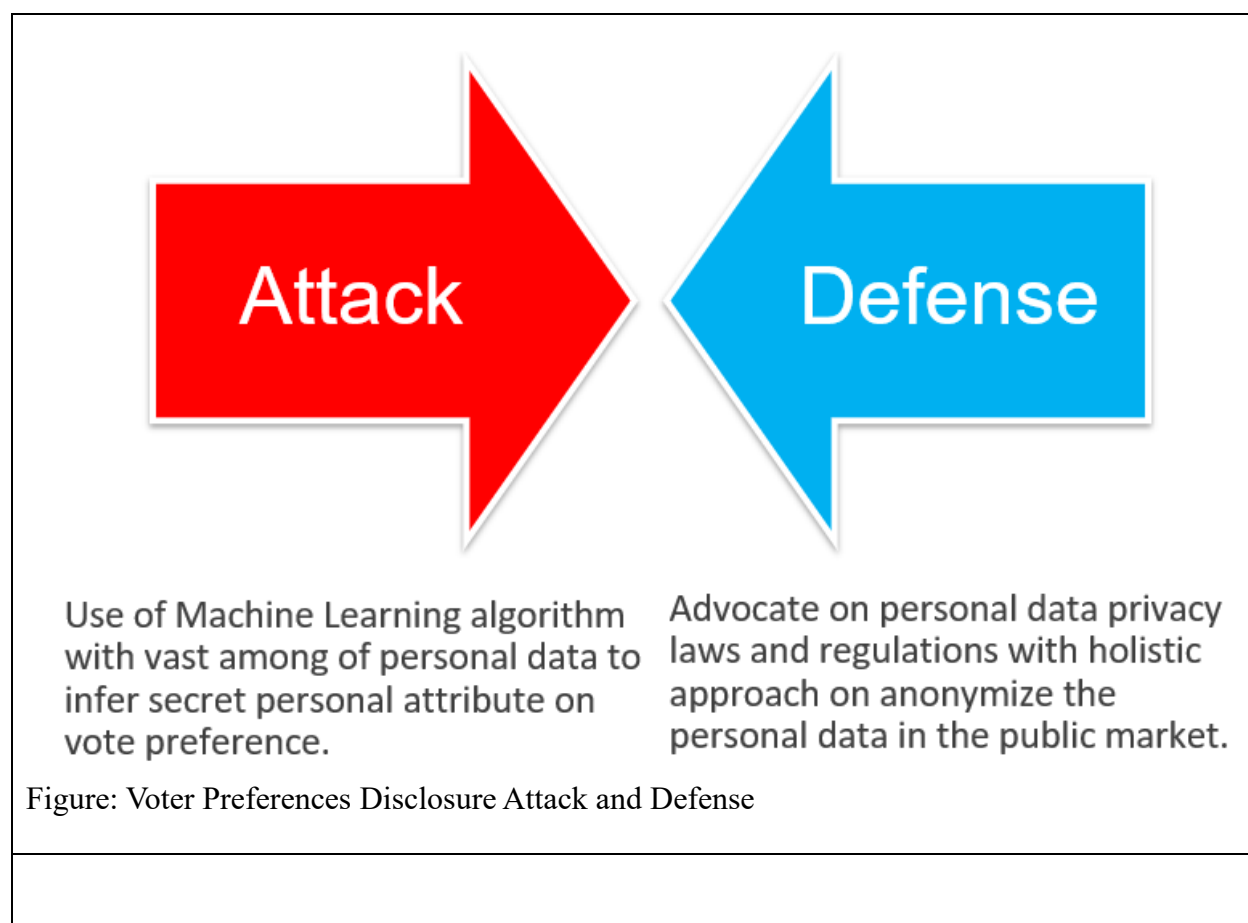
Social Media has been targeted as mega data hubs of personal active and passive data. In 2018, the Facebook¹⁰ data scandal was revealed that Cambridge Analytica had harvested the personal

⁸ <https://investor.equifax.com/~media/Files/E/Equifax-IR/reports-and-presentations/events-and-presentation/investor-relations-presentation-december-4-2018.pdf>

⁹ <http://money.cnn.com/2015/02/04/technology/anthem-insurance-hack-data-security/>

¹⁰ <https://www.theguardian.com/technology/2018/apr/10/facebook-notify-users-data-harvested-cambridge-analytica>

data of millions of people's Facebook profiles without their consent and used it for political advertising purposes.



Assets	Threats	Description
Voter Database	Confidentiality, Integrity	Public static database contains personal attributes such as full legal name, birth date, gender, residential address and registered party affiliation. Data broker provides services for distribution with fee and might subject to data poisoning attack.

US Census Database	Confidentiality	Public statistical database contains summary statistics of person attributes such as ethnicity, income, education, marriage, and number of children. Combination of synthetic data generation, data aggregation and linkage is very likely to infer sensitive personal attributes such as vote preferences.
Social Media	Confidentiality	Data leakage via public post with personal attributes, personal associations with degrees of separation, and subject to targeted marketing advertisements.
Personal Data Market	Confidentiality, Integrity	Collective Data breach events creates black market of personal data on the Web. Consumer service provider collects and resell the personal attributes for profit or in exchange of free service. Personal data exchange market promotes monetized model for individual. The personal data broker is subject to data poisoning attack.
Personal Privacy	Integrity	Personal service right might subject to perpetuate discrimination because machine learning result are trained on biased data. An individual might not receive the service they needed.

3 Design of Experiment

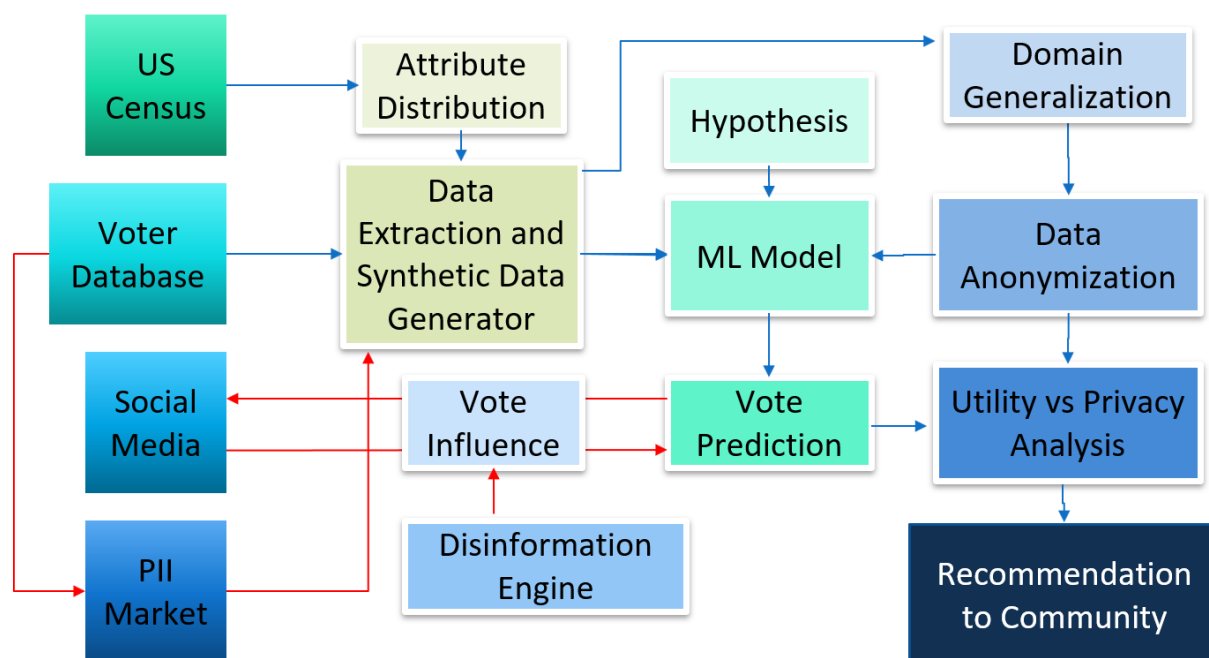


Figure: Design of AI Inference

3.1 Attribute Distribution

2018 estimates available from 2010 census data, for income, education, family size, and ethnicity. These can be filtered by geographic area, gender, and age. Data is difficult to extract from web interface. API exists but the interface for getting an apikey is broken. Margin of error can rapidly overwhelm any information content, as filters lower the number of class members. Due to margin error, not possible to map to voter database, beyond large sample proportional distributions. Real world individual voter profiles exist, with all of these attributes and many more. These exist as purchasable data, as well as advertising API targets.

3.2 Data Extraction and Generation (Bootstrapping)

--	--

```

Input: Voter Data File, vd_file_name, Voter File Table, VFT: vf_attributes = ["state_file_id", "dob", "sex",
"party", "county_registered_address", "zip_registered_address", "state_registered_address",
"federal_district"], Synthetic Attribute Distribution, p_ethnicity, income, education, marriage,
children]

Output: Data Output Table: dt_attributes = ['state', 'zipcode', 'ethnicity', 'income', 'age', 'sex',
'education', 'marriage', 'children', 'party']

Algorithm:
vf_data = load_VF_data(vd_file_name, vf_attributes)

sd_ethnicity = generateEthnicity(len(vf_data), p_ethnicity)
sd_income = generateIncome(len(vf_data), p_income)
sd_education = generateEducation(len(vf_data), p_education)
sd_marriage = generateMarriage(len(vf_data), p_marriage)
sd_children = generateChildren(len(vf_data), p_children)

for i, row_ in enumerate(vf_output):
    for index in range(0, 10):
        dt_row = generate_columns(index, row_, dt_attributes)
        append_row(dt_row)

output(dt_output_file)

```

Figure: Data Bootstrapping Algorithm

3.3 Machine Learning (ML) Model

3.3.1 Hypothesis

Use of Machine Learning Algorithm to infer voter preferences (secret attribute) from personal attributes includes {'ethnicity', 'income', 'age', 'sex', 'education', 'marriage', 'children'} extracted from voter registry and linked with PII data from US Census, PII data market.

3.3.2 Supervised Learning with Classification

Please reference to section 4 for more detail on Attack with Machine Learning Model.

3.3.3 Vote Preference Inference (Prediction)

Please reference to section 4 for more detail on result of vote preference inference on the accuracy of prediction with machine learning model. We will compare the prediction score using the raw data and anonymized data set.

3.4 Domain Generalization

qID = {'ethnicity', 'income', 'age', 'sex', 'education', 'marriage', 'children'};

Ethnicity Generalization Tree:	Sex Generalization Tree:
Income Generalization Tree:	Education Generalization Tree:

Marriage Generalization Tree:	Children Generalization Tree:
Age Generalization Tree:	

3.5 Data Anonymization

Please reference to section 5 for more detail on Defense with Principles of Privacy Engineering.

3.6 Vote Influence

3.7 Disinformation Engine

4 Attack: Machine Learning

5 Defense: Privacy Engineering

In our design for applying privacy engineering, we choose k-anonymity for the simplicity and easy to implement with known effectness on defense against identity disclosure attack. For the voter preference as sensitive attribute, it is not considered as a publishable attribute. Voter registry might contain the voter party but it does not necessary reflect the true vote. To achieve k-anonymity, we choose datafly algorithm with domain generalization for categoricial attributes; developed by Sweeney in 2017 is global, bottom-up, and greedy generalization algorithm. We have referenced and enhanced datafly implementation ¹¹ by Alessio Vierti.

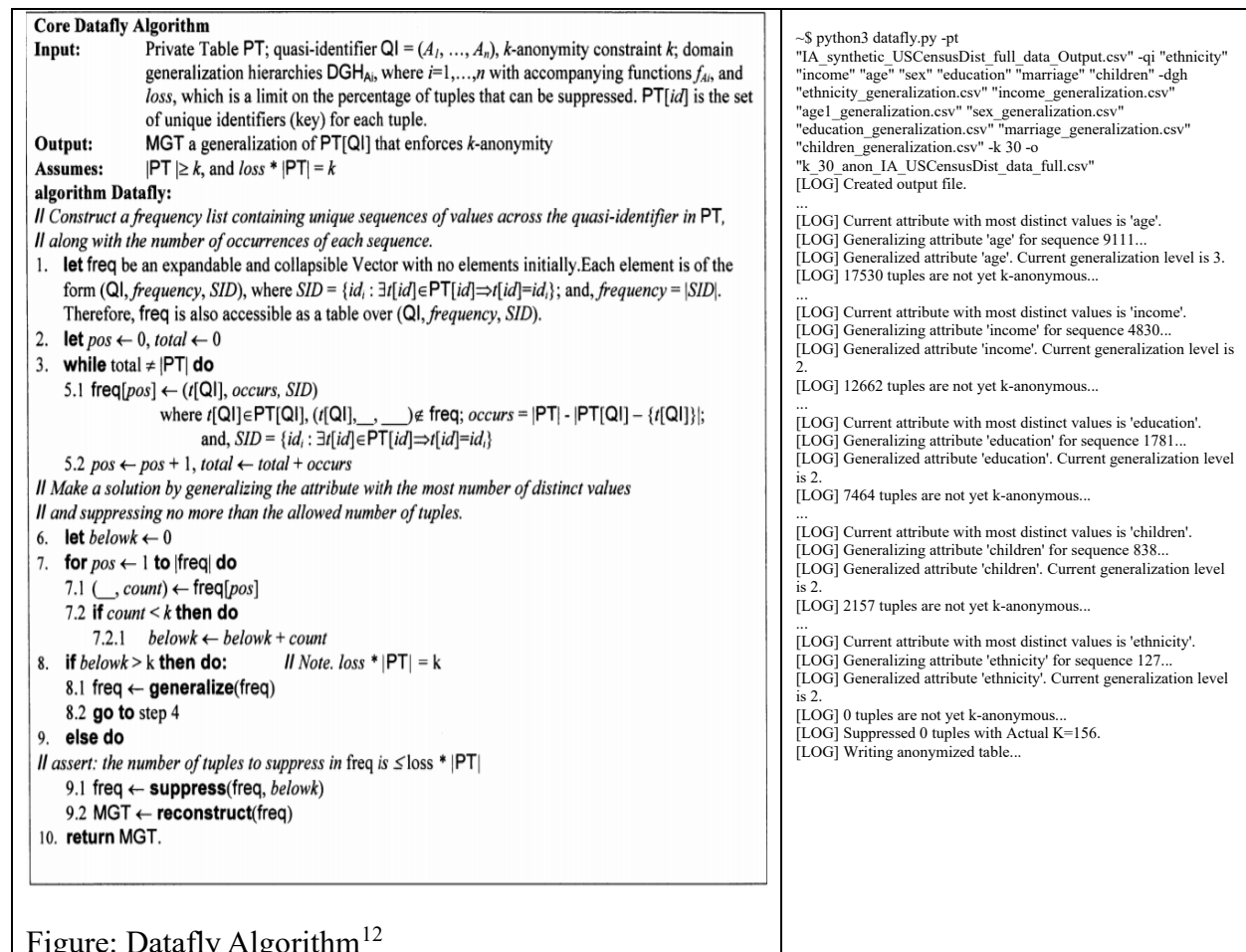


Figure: Datafly Algorithm¹²

¹¹ <https://github.com/alessiovierti/python-datafly>

¹² Computational disclosure control : a primer on data privacy protection - Sweeney, Latanya

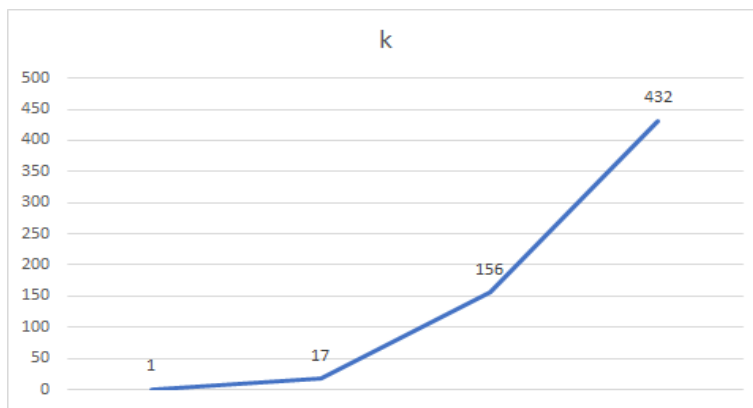


Figure: Data Anonymization

qID = {'ethnicity', 'income', 'age', 'sex', 'education', 'marriage', 'children'};

k=1 on the raw data mixed from voter data and synthetic attributes.

With current design of domain generation on personal features, the Actual $K=\{1, 17, 156, 432\}$ can be anonymized using Datafly algorithm.

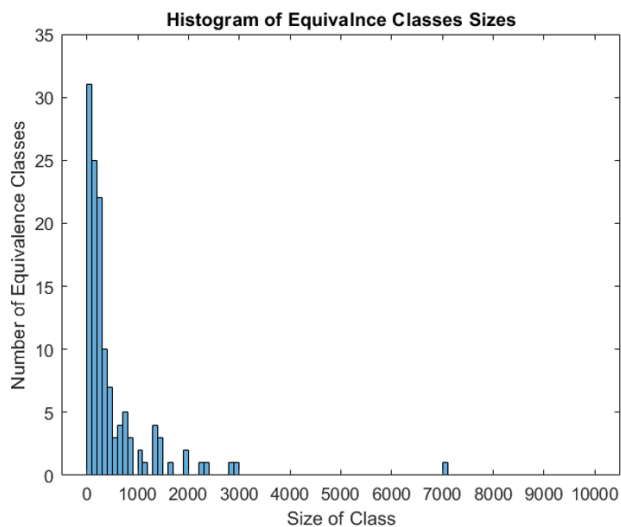
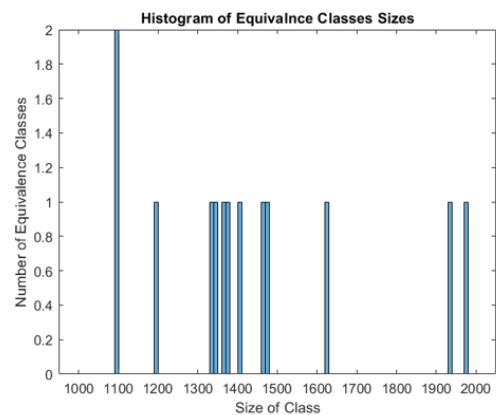
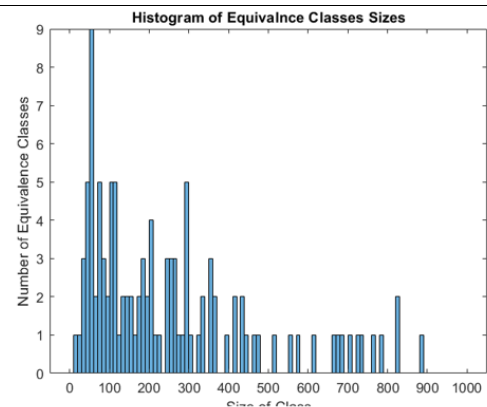
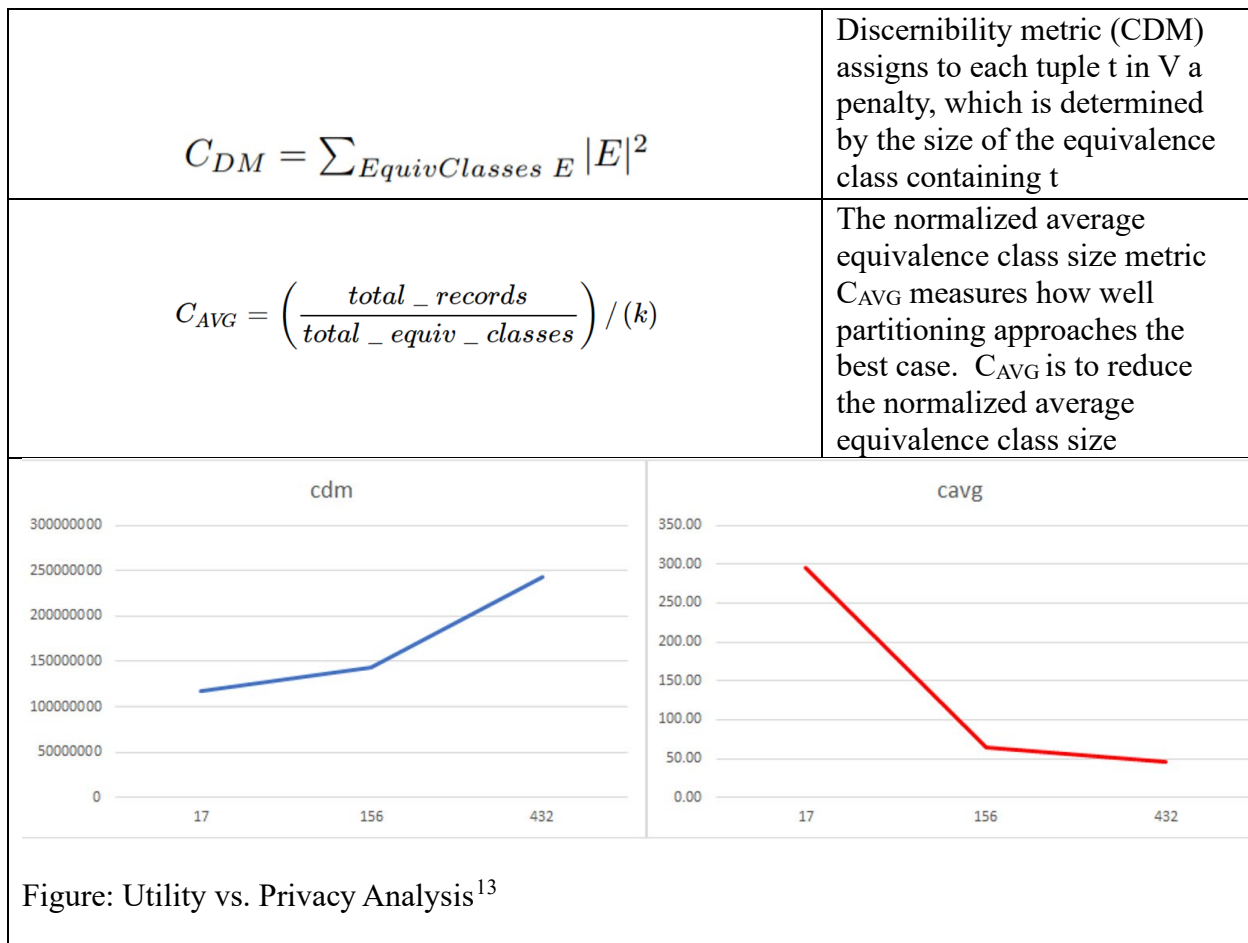


Figure: Equivalence Classes Distribution for Anonymization dat set with K=17



We first verified the raw data set extracted from voter registry and generated from attribute distribution has $k=1$. We have ran data fly algorithm with different desired k outputted with actual k graphed in Data Anonymization figure. On the actual $k=17$ data set, we have graphed the histogram of equivalence classes distribution in different resolution to understand quality of anonymized data for the entropy value. In order to understand the quality of anonymization in respect to different choice of k , we graphed the discernibility metric and normalized average equivalence class size metric to determine the anonymized data set with $k=17$ will be the best choice to maintain best privacy while allow maximum utility.



¹³ <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3243250/>

6 Recommendation

-

7 Conclusion