

**PRIVACY PRESERVING DATA PUBLISHING  
FRAMEWORK FOR UNSTRUCTURED TEXTUAL  
SOCIAL MEDIA DATA**

Peruma Baduge Prasadi Apsara Abeywardana

(189302A)

Dissertation submitted in partial fulfillment of the requirements for the degree Master  
of Science in Computer Science

Department of Computer Science and Engineering

University of Moratuwa  
Sri Lanka

July 2020

## **DECLARATION**

I declare that this is my own work and this dissertation does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works.

---

P.B.P.A. Abeywardana

Date

I certify that the declaration above by the candidate is true to the best of my knowledge. The above candidate has carried out research for the Masters dissertation under my supervision.

---

T. Uthayasanker (PhD)

Date

## **Abstract**

Privacy has become an essential part of data science and analytics due to the potential of personal data misuse. As a result of privacy breaches reported in various analytical studies privacy preservation has become a legal responsibility rather than a simple social responsibility. Preserving privacy of unstructured data is more challenging compared to structured data. Social media has become largely popular over the past couple of decades and they are pumping a huge amount of data at a high velocity into analytical systems. Social media profiles contain a wealth of personal and sensitive information, creating enormous opportunities for third parties to analyze them with different algorithms, draw conclusions and use in disinformation campaigns and micro targeting based dark advertising. The primary goal of this study is to provide a mitigation mechanism for privacy breaches happening via disinformation campaigns that are done based on the insights extracted from personal/sensitive data analysis. Specifically, this research is aimed at building a privacy preserving data publishing framework for unstructured and textual social media data without compromising the true analytical value of those data. A novel way is proposed to apply traditional structured privacy preserving techniques on unstructured data. Creating a comprehensive twitter corpus annotated with privacy attributes is another objective of this research, especially because the research community is lacking one.

An easily extensible framework that can be adopted by many domains is implemented here, integrating different concepts from the literature. A comprehensive set of experiments are also performed in order to assess the capabilities of the machine learning models, algorithms as well as to simulate some real-world privacy preserving data publishing use cases.

## ACKNOWLEDGMENTS

I would like to express my gratitude to my advisor, Dr. T. Uthayasanker, for his invaluable support, inspiration, supervision, and useful suggestions throughout this research work. He was never reluctant to guide me through composing this report in a successful manner.

Also, I would like to thank my family who was very supportive throughout the process, specially my parents and my husband. And also, I am grateful to my friends and colleagues with whom I discussed the concepts and who gave me back valuable inputs and feedback.

I must thank the independent data annotators Sriyoukan Sriranjan and Lavanaraj Sivarasa, who helped me annotating the data corpus.

## TABLE OF CONTENTS

DECLARATION .....	i
Abstract.....	ii
ACKNOWLEDGMENTS .....	iii
TABLE OF CONTENTS.....	iv
LIST OF FIGURES .....	vii
LIST OF TABLES .....	viii
LIST OF ABBREVIATIONS.....	ix
CHAPTER 1: INTRODUCTION .....	1
1.1    Personal data .....	1
1.2    Personal data in social media .....	2
1.2.1    Social threats of personal data analysis.....	3
1.3    Data protection regulations .....	3
1.3.1    General Data Protection Regulation (GDPR) .....	4
1.3.2    Russian Federal Law on Personal Data.....	4
1.3.3    German Bundesdatenschutzgesetz (BDSG).....	4
1.4    Motivation.....	4
1.5    Problem statement.....	5
1.6    Research objectives.....	5
1.7    Outline.....	6
CHAPTER 2: LITERATURE REVIEW .....	7
2.1    Existing PPDp techniques.....	7
2.1.1    Suppression .....	7
2.1.2    Generalization .....	8
2.1.3    Swapping.....	9
2.1.4    Anatomization.....	9
2.1.5    Permutation .....	9
2.1.6    Perturbation.....	10
2.2    Existing privacy models.....	11
2.2.1    k-anonymity .....	12
2.2.2    l-diversity .....	14
2.2.3    t-closeness .....	16
2.3    Real world applications of PPDp .....	20
2.3.1    Mobile data .....	20
2.3.2    Health care data.....	20

2.3.3	Social media data .....	22
2.4	Privacy metrics.....	24
2.4.1	Confidence level .....	24
2.4.2	Average conditional entropy .....	24
2.4.3	Hidden failure .....	25
2.5	Utility metrics.....	25
2.5.1	Generalization/Suppression counting.....	25
2.5.2	Loss Metric (LM).....	25
2.5.3	Discernibility Metric (DM).....	26
2.5.4	KL divergence.....	26
2.5.5	Bivariate measures .....	26
2.5.6	Workload – aware metrics .....	26
2.6	Future directions .....	27
2.7	Challenges with unstructured data .....	27
2.7.1	Difficulty in identifying sensitive attributes .....	28
2.7.2	Volume of data.....	28
2.7.3	Quality of data.....	28
2.8	Unstructured privacy preserving data publishing scenarios.....	28
2.9	Unstructured data anonymizing techniques .....	29
2.10	Summary .....	30
	<b>CHAPTER 3: METHODOLOGY .....</b>	<b>31</b>
3.1	High-level architecture.....	31
3.2	Technologies adopted .....	33
3.3	Privacy attribute extractor.....	34
3.3.1	Twitter corpus .....	34
3.3.2	Data preprocessing.....	38
3.3.3	Data transformation.....	38
3.4	Privacy attribute anonymizer .....	44
3.4.1	Simple anonymization.....	45
3.4.2	K anonymization.....	46
3.5	Utility evaluator .....	51
3.5.1	Discernibility metrics .....	51
3.5.2	Loss metrics .....	51
3.5.3	Generalization counting .....	51
3.6	Sample web client .....	52

3.7	Summary .....	55
CHAPTER 4: EXPERIMENTAL DESIGN AND RESULTS .....		56
4.1	Evaluating privacy attribute annotator accuracy.....	56
4.1.1	k-fold cross validation.....	56
4.1.2	Confusion matrix.....	57
4.1.3	Experimental setup.....	59
4.1.4.	Experimental results.....	59
4.2	Mocking real world PPDP workflows .....	62
4.2.1	Single tweet experiment.....	63
4.2.1	Multiple tweets experiment.....	63
4.2.3	Twitter live search experiment.....	65
4.2.4	Usability evaluation of anonymized dataset .....	65
4.3	Summary .....	66
CHAPTER 5: CONCLUSION.....		67
5.1	Summary .....	67
5.2	Research outcomes.....	68
5.3	Research limitations.....	69
5.4	Future work.....	69
5.5	Discussion.....	70
REFERENCES .....		72

## **LIST OF FIGURES**

Figure 2.1: Overall Architecture Proposed by Gardner et al. [30].....	21
Figure 2.2: Graphical Representation of Social Network Data.....	23
Figure 3.1: High Level Architecture .....	33
Figure 3.2: API Documentation.....	34
Figure 3.3: Tweets Transformation Process .....	38
Figure 3.4: Data Partitioning Process .....	47
Figure 3.5: Partitioning function.....	48
Figure 3.6: Textual data converted to structured format.....	49
Figure 3.7: K-Anonymized data frame .....	50
Figure 3.8: Snapshots from the web client .....	53
Figure 3.9: Snapshots from the web client .....	54
Figure 3.10: Snapshots from the web client .....	54
Figure 3.11: Snapshots from the web client .....	55
Figure 4.1: Confusion Matrix.....	58
Figure 4.2: Results from the experimental data set annotation.....	64
Figure 4.3: Results from the experimental data set k-anonymization.....	64
Figure 4.4: Results from the live twitter data set annotation .....	65

## LIST OF TABLES

Table 2.1: Health Records .....	11
Table 2.2: Anonymized Health Records .....	11
Table 2.3: 4-Anonymized Health Records .....	13
Table 2.4: 3-Diverse Health Records .....	15
Table 2.5: 3-Diverse Anonymized Health/Salary Records .....	16
Table 2.6: 10000 Records of a Virus that affects only 1% of the Population .....	17
Table 2.7: Summary of Privacy Models .....	20
Table 3.1: Annotation Scheme .....	35
Table 3.2: Corpus Statistics .....	36
Table 3.3: Annotation Statistics.....	36
Table 3.4: Feature List .....	40
Table 3.5: spaCy Named Entities .....	41
Table 3.6: Classifiers Attempted .....	42
Table 3.7: Aggregated Accuracy of the models .....	43
Table 3.8: Libraries used for Implementing Automatic Tagger.....	44
Table 3.9: Anonymization Scheme .....	45
Table 3.10: Anonymization Techniques Applied .....	48
Table 3.11: Libraries used for the data anonymization module .....	50
Table 4.1: Average Accuracy Values for Each Algorithm .....	60
Table 4.2: Accuracy Details per Class.....	62
Table 4.3: Utility comparison after anonymization .....	66

## **LIST OF ABBREVIATIONS**

API	Application Programming Interface
EU	European Union
EEA	European Economic Area
BDSG	Bundesdatenschutzgesetz
GDPR	General Data Protection Regulation
PPDP	Privacy Preserving Data Publishing
GSR	Global Science Research
PPDM	Privacy Preserving Data Mining
MinGen	Minimum Generalization Algorithm
PPDC	Privacy Preserving Data Collection
GPS	Global Positioning System
LM	Loss Metric
DM	Discernibility Metrics
IoT	Internet of Things
UK DA	United Kingdom Data Archive
KL Distance	Kullback-Leibler Distance