

Security and Privacy Issues in Bipolar Disorder Research

Panagiotis Rizomiliotis, Aggeliki Tsohou, Costas Lambrinoudakis, Stefanos Gritzalis

Department of Information and Communication Systems Engineering, University of the Aegean, Samos, Greece.

ABSTRACT

Mental health diseases are common but research to further knowledge and understanding of them is hampered by data privacy and confidentiality regulations that apply to medical records. Centralised databases containing the relevant medical history of thousands of patients with an individual mental disease would be of great value for researchers, enabling techniques such as data mining to be applied. The major challenge in achieving this is anonymising the data to satisfy legal and ethical requirements without removing important clinical information. In this paper we propose a model that can be used to create a central repository of anonymised data for patients with bipolar disease. Knowledge obtained from the database is fed into an expert system which can guide clinicians in patient management. Security requirements are provided by access to the database being controlled by RBAC (Role Based Access Control).

INTRODUCTION

Mental disorders or mental illnesses such as mood disorders, anxiety disorders, psychotic disorders, eating disorders, and personality disorders affect approximately 1 in 4 of the population¹. One relatively common serious mental illness is bipolar disorder (BD), which is also known as manic depression, manic depressive disorder or bipolar affective disorder. BD is characterised by episodes of full-blown mania which is defined as periods of abnormally expanded or irritable mood, along with major depression. These episodes can have devastating consequences on the professional and social life of those affected. Alcohol and drug abuse and dependence, and social and professional isolation, are the most common complications of BD. Around 10–20% of bipolar patients who have been hospitalised at some stage during their

Correspondence and reprint requests: Panagiotis Rizomiliotis, Department of Information and Communication Systems Engineering, University of the Aegean, Samos GR-83200, Greece. E-mail: prizomil@aegean.gr.

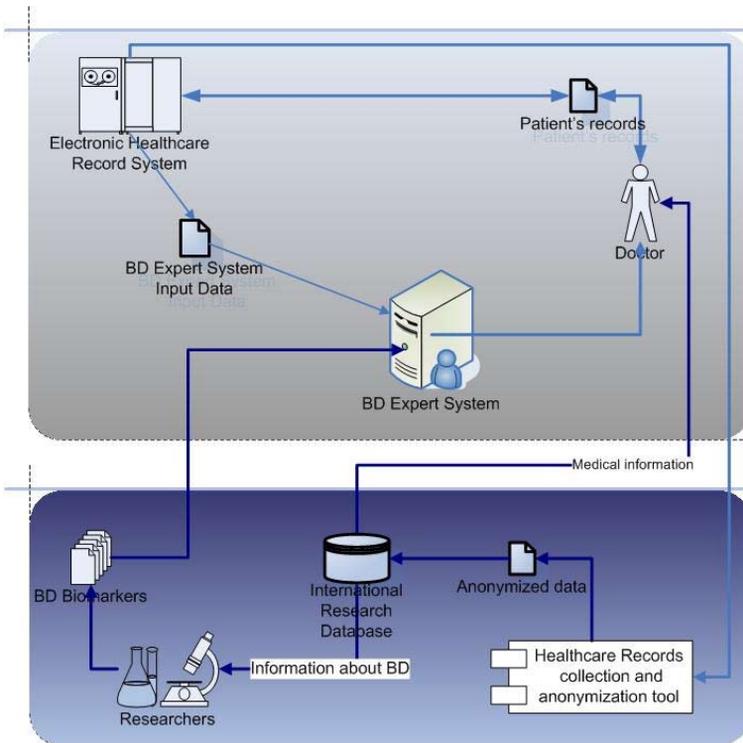
1 treatment, die by committing suicide^{2,3}. BD is a significant health problem worldwide
 2 as its prevalence is approximately 1% across all populations.

3 Despite the fact that BD is fairly common, diagnosis and management of indi-
 4 vidual cases is still frequently sub-optimal. This is partly due to a lack of detailed
 5 knowledge and understanding of BD and this is not helped by the fact that research
 6 in the disease is hampered by patient data being classified as sensitive and confiden-
 7 tial. As a result knowledge that could result from subjecting data from large numbers
 8 of patients to techniques such as data analysis and data mining cannot be obtained.
 9 Achieving this in practice requires preserving anonymity and confidentiality when
 10 pooling data from large numbers of BD patients.

11 In this paper we propose a model to enable centralised collection of anonymised
 12 data from patients with BD. This data can be studied by researchers to further
 13 knowledge and understanding of BD and the new knowledge can be fed into an
 14 expert system that clinicians can use to assist them in patient management.

15
 16 **SYSTEM ARCHITECTURE**

17
 18 The architecture of the proposed system is shown in Figure 1.



19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41 **Figure 1.** Proposed architecture for an anonymised research database and expert
 42 system for bipolar disease (BD)

1 The top half of the figure demonstrates how a doctor obtains help with patient
2 management. The expert system processes data collected from the patient and infor-
3 mation stored in the patient’s electronic health record (EHR), including information
4 relevant to BD. The output of the system is based on state-of-the-art knowledge
5 about BD. The knowledge is codified in a special form known as a BD-biomarker.
6 If necessary the doctor may also consult the research database, for example if he
7 wants to obtain details on the case management of another patient similar to the
8 one he is treating.

9 The bottom half of the figure shows how the model supports research into BD.
10 A special tool is used to collect and anonymise data from different EHR systems.
11 The collected data is used to update a Research Database. Authorised personnel
12 can access to the Research Database and new knowledge gained from the database
13 is formalised as electronic BD biomarkers. This is made available for interaction
14 through the expert system presented in the top half of the diagram. The key aspects
15 of the system with respect to anonymisation of data and meeting legal and ethical
16 requirements will now be described in detail.

17
18 SECURITY AND PRIVACY REQUIREMENTS

19
20 In order to implement the BD research central repository, the processing of health
21 data stored in geographically spread EHRs is inevitable. The communication between
22 the different EHR systems and the repository of accumulated anonymised BD-
23 related data system has to be carefully designed in order to guarantee the confiden-
24 tiality, authenticity and integrity of data. The provided solution must be generic and
25 flexible in the sense that it should address different systems ranging from current
26 healthcare systems to legacy IT systems to allow interoperability between the EHRs
27 and the repository and also with the BD research community systems. Finally access
28 to the research data should be restricted only to authorised users.

29 Processing health data, by definition, raises several security and privacy issues,
30 such as the protection of data integrity and confidentiality and the preservation of
31 the patient’s privacy. Health data belongs to a special category of personal data, com-
32 monly known as sensitive data. Legally all data contained in medical documenta-
33 tion such as electronic health records is considered as sensitive data. Consequently
34 according to legislation it cannot be shared in a way that identifies a patient without
35 the patient’s explicit permission. Data can be anonymised by removing elements
36 from it, but the challenge then becomes of how to best remove data to ensure that
37 the patient cannot be identified whilst at the same time ensuring that the remaining
38 data contains all the key elements necessary for research purposes including data
39 mining.

40
41
42

1 ANONYMISING DATA FOR DATA MINING

2
3 A number of techniques have been proposed for modifying or transforming data in
4 such a way so as to preserve privacy, but leaving it suitable for data mining. Some
5 examples of these are:

6
7 *The randomisation method:* The randomisation technique uses data distortion meth-
8 ods in order to create private representations of the records^{4,5}. In most cases, the
9 individual records cannot be recovered, but only aggregate distributions can be
10 recovered. These aggregate distributions can be used for data mining purposes.
11 The randomisation approach is particularly well suited to privacy-preserving data
12 mining of streams, since the noise added to a given record is independent of the
13 rest of the data. The most common methods of randomisation are those of additive
14 perturbations and multiplicative perturbations.

15
16 *The k-anonymity model and l-diversity:* The k-anonymity method is based on tech-
17 niques such as generalisation and suppression according to which any given record
18 maps to at least k other records in the data. The k-anonymity model was developed
19 in order to prohibit the indirect identification of records from public databases, since
20 combinations of record attributes can be used to exactly identify individual records.
21 The l-diversity model was designed to handle some weaknesses in the k-anonymity
22 model. Protecting identities to the level of k-individuals is not the same as protecting
23 the corresponding sensitive values, especially when there is homogeneity of sensitive
24 values within a group. Thus, the concept of intra-group diversity of sensitive values
25 is promoted within the anonymisation scheme⁶.

26
27 *Distributed privacy preservation:* A partition is a division of a logical database or
28 its constituting elements into distinct independent parts. The partitioning may be
29 horizontal (when the records are distributed across multiple entities) or vertical
30 (when the attributes are distributed across multiple entities). There are applications
31 where users wish to derive aggregate results from data sets partitioned across other
32 individuals. While the individuals do not desire to share their entire data sets, they
33 consent to limited information sharing. The overall effect of such methods is to
34 maintain privacy for each individual entity, while deriving aggregate results over
35 the entire data⁷⁻⁹.

36
37 *Downgrading Application Effectiveness:* The output of applications such as associa-
38 tion rule mining, classification or query processing may lead to violations of privacy
39 and motivated research into downgrading the effectiveness of applications by either
40 data or application modifications. Such techniques include association rule hiding¹⁰,
41 classifier downgrading¹¹, and query auditing¹².

1 Each one of the above techniques has its own advantages and disadvantages.
2 Depending on the application the system designer has to choose the most adequate
3 method. In this context, it is not a straightforward task to identify the most appropriate
4 techniques for the anonymisation of medical data.

6 ANONYMISING MEDICAL DATA

8 Two specific systems that have been developed to anonymise medical data are:

- 9 • **The Scrub system**¹³
- 10 • **The Datafly system**¹⁴

11 The Scrub system¹³ was designed for de-identification of clinical notes which
12 usually occur in the form of textual data and contain references to patients, patients'
13 family members, addresses etc. The Scrub system uses detection algorithms, based
14 on several local knowledge sources, to determine when a block of text leaks information
15 concerning the name, address or a phone number of a patient or a member of
16 its family. This system was proposed in order to replace the traditional, and in most
17 cases insufficient techniques, based on a simple “search and replace procedure”.

18 The Datafly System¹⁴ is one of the earliest practical systems for anonymisation
19 and one of the first applications of privacy-preserving transformations. The system
20 was designed in response to the concern that the process of removing only directly
21 identifying attributes such as social security numbers was not sufficient to guarantee
22 privacy. This work has a similar motive as the k -anonymity approach of preventing
23 record identification, but it does not formally use a k -anonymity model in order to
24 prevent identification through linkage attacks. The Datafly system, as well as most
25 of its successors, proposes anonymity levels ranging from 0 to 1. An anonymity level
26 of 0 results in Datafly providing the original data, whereas an anonymity level of 1
27 results in the maximum level of generalisation of the underlying data.

28 To enable collection of BD data in a centralised database we propose using the
29 distributed databases model. All data is collected and stored in a local database
30 maintained at the hospital or the clinic that treats the patient. The anonymisation
31 process is applied to every distributed database and the anonymised data is then
32 stored in the research repository. This will give research community members
33 access to a full collection of anonymised clinical data. The major challenge with
34 this approach is updating the centralised database. While it is relatively easy to add
35 new records, special care is required for updating existing ones. One solution to
36 this problem is the application of general purpose secure multiparty computation
37 techniques borrowed from the cryptographic literature¹⁵.

38 The privacy preservation techniques that is the most appropriate for the proposed
39 architecture is the *Distributed Privacy Preservation*. In the proposed scenario, even
40 though several entities (hospitals, clinics, individual doctors) do not desire to share
41 their entire data sets, they are willing to give their consent to limited information
42 sharing. At the same time, there are entities wishing to derive aggregate results from

1 data sets partitioned across other individuals. To achieve this, distributed algorithms
2 for k-Anonymity can be used, combining previous proposed solutions⁷⁻⁹, in order to
3 maintain k-anonymity across different distributed parties. It is assumed that the data
4 record has both sensitive attributes and quasi-identifier attributes. The solution uses
5 encryption on the sensitive attributes which can be decrypted only if there are at
6 least k records with the same values on the quasi-identifiers⁹. Thus, k-anonymity is
7 maintained. The issue of k-anonymity is also important in the context of hiding iden-
8 tification in the context of distributed location databases^{7,8}. In this case, k-anonymity
9 of the user-identity is maintained even when the location information is released.

11 ACCESS CONTROL

13 The access to the research database must follow a role-based access control (RBAC)
14 policy. Ideally it should define specific roles that will be authorised to access the
15 research database, associating specific access privileges to each role. Examples of
16 some of the roles and privileges are as follows:

- 17 • *Health Researchers* (e.g. clinicians, scientists, pharmaceutical companies, etc)
18 who study BD. Their need is met by granting access on all available medical
19 data (patient medical history, treatments, pharmaceutical substances etc)
20 related to their field of interest.
- 21 • *Doctors* providing health care services to BD-patients. In cases where the
22 diagnosis or/and treatment of the patient is not straight forward, the doctor
23 will be able to obtain guidance from the expert system and also may be able
24 to access previous similar cases that have been treated by other clinicians.
- 25 • *BD-Biomarker administrators* will be responsible for 'representing' new knowl-
26 edge about the disease in a structured format, known as a BD-biomarker. The
27 biomarkers are then utilised by the expert systems for supporting clinicians
28 in the optimal management of patients with BD.

30 DISCUSSION

32 In this paper we have proposed a model that can be used to support research to
33 further knowledge and understanding of bipolar disease and also aid clinicians in
34 managing individual patients using the best current available knowledge. The former
35 objective is achieved by providing a centralised system architecture that enhances
36 access to BD related data. To meet the requirements for data privacy and confiden-
37 tiality, legal and ethical requirements the model includes an anonymisation process
38 and a role-based access control policy. The implementation of the proposed system
39 will enhance the secure interoperability and seamless communication of BD health
40 data between clinicians, health researchers, and those responsible for creating the
41 knowledge in the Expert system. This can then be used by clinicians to aid them in
42 patient management.

1 REFERENCES

- 2
- 3 1 Andrews G, Totov N. Depression is very disabling. *Lancet*. 2007; **370**: 808–9.
- 4 2 Harris E, Barraclough B. Suicide as an outcome for mental disorders: a meta-analysis. *The*
5 *British Journal of Psychiatry* 1997; **170**: 205–28.
- 6 3 Bostwick JM, Pankratz V. (2000). Affective disorders and suicide risks: a reexamination. *Am*
7 *J. Psychiatry* 2000; **157**: 1925–32.
- 8 4 Agrawal R, Srikant R. Privacy-preserving data mining. *Proceedings of the ACM SIGMOD*
9 *Conference, 2000*.
- 10 5 Agrawal D, Aggarwal C. On the Design and Quantification of Privacy-Preserving Data Min-
11 ing Algorithms. *ACM PODS Conference, 2002*.
- 12 6 Machanavajjhala A, Gehrke J, Kifer D, Venkitasubramaniam M. l-Diversity: Privacy Beyond
13 k-Anonymity. *ICDE* (2006).
- 14 7 Bettini C, Wang XS, Jajodia S. Protecting privacy against location based personal identifica-
15 tion. *Proc. of Secure Data Management Workshop*, Trondheim, Norway, 2005.
- 16 8 Gedik B, Liu L. A customizable k-anonymity model for protecting location privacy, *ICDCS*
17 *Conference, 2005*.
- 18 9 Zhong S, Yang Z, Wright R. Privacy-enhancing k-anonymization of customer data, In *Pro-*
19 *ceedings of the ACM SIGMOD-SIGACT-SIGART Principles of Database Systems*, Baltimore,
20 2005.
- 21 10 Verykios VS, Elmagarmid A, Bertino E, Saygin Y, Dasseni E. Association Rule Hiding. *IEEE*
22 *Transactions on Knowledge and Data Engineering* 2004; **16**(4).
- 23 11 Moskowitz I, Chang L. A decision theoretic system for information downgrading. *Joint Con-*
24 *ference on Information Sciences, 2000*.
- 25 12 Adam N, Wortmann JC. Security-Control methods for statistical databases: a comparison
26 study. *ACM Computing Surveys*, 1989; **21**: 515–56.
- 27 13 Sweeney L. Replacing personally-identifying information in medical records, the Scrub
28 system. In: Cimino, JJ, (ed.), *Proceedings, Journal of the American Medical Informatics*
29 *Association*. Washington, DC: Hanley & Belfus, Inc, 2006, pp. 333–37.
- 30 14 Sweeney L. Guaranteeing anonymity while sharing data, the Datafly system. *Proc AMIA Annu*
31 *Fall Symp*. 1997, pp. 51–55.
- 32 15 Goldreich O, Micali S, Wigderson A. How to play the mental game. *Proceedings of the 19th*
33 *Annual ACM Symposium on Theory of Computing*, ACM, 1987, pp. 218–29.
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42

