

# Data Mining Methodologies for Pharmacovigilance

Mei Liu

Computer Science  
Department, New Jersey  
Institute of Technology,  
Newark, NJ, 07102, USA  
mei.liu@njit.edu

Michael E. Matheny<sup>1</sup>

Department of Biomedical  
Informatics, Vanderbilt  
University, School of  
Medicine, Nashville, TN  
37203, USA  
michael.matheny@vanderbilt.edu

Yong Hu

Institute of Business  
Intelligence, Guangdong  
University of Foreign Studies,  
Sun Yat-sen University,  
Guangzhou, 510006, China  
henryhu200211@163.com

Hua Xu

Department of Biomedical  
Informatics, Vanderbilt  
University, School of  
Medicine, Nashville, TN  
37203, USA  
hua.xu@vanderbilt.edu

## ABSTRACT

Medicines are designed to cure, treat, or prevent diseases; however, there are also risks in taking any medicine - particularly short term or long term adverse drug reactions (ADRs) can cause serious harm to patients. Adverse drug events have been estimated to cause over 700,000 emergency department visits each year in the United States. Thus, for medication safety, ADR monitoring is required for each drug throughout its life cycle, including early stages of drug design, different phases of clinical trials, and post-marketing surveillance. Pharmacovigilance (PhV) is the science that concerns with the detection, assessment, understanding and prevention of ADRs. In the pre-marketing stages of a drug, PhV primarily focuses on predicting potential ADRs using preclinical characteristics of the compounds (e.g., drug targets, chemical structure) or screening data (e.g., bioassay data). In the post-marketing stage, PhV has traditionally involved in mining spontaneous reports submitted to national surveillance systems. The research focus is currently shifting toward the use of data generated from platforms outside the conventional framework such as electronic medical records (EMRs), biomedical literature, and patient-reported data in online health forums. The emerging trend of PhV is to link preclinical data from the experimental platform with human safety information observed in the post-marketing phase. This article provides a general overview of the current computational methodologies applied for PhV at different stages of drug development and concludes with future directions and challenges.

## Keywords

Pharmacovigilance, drug safety surveillance, adverse drug reaction detection.

## 1. INTRODUCTION

Every year the US public spends billions of dollars on prescription drugs for the cure, treatment, or prevention of diseases. However, caution should be taken seriously when taking any medication because severe adverse drug reactions (ADRs) can lead to patient morbidity. ADRs are often referred to as “any unintended and undesirable effects of a drug beyond its anticipated therapeutic effects occurring during clinical use” [1]. According to the national surveillance study of emergency department visits for outpatient adverse drug events by Budnitz et al. [2], there were total 21,298 adverse drug events from January 1, 2004 through December 31, 2005, yielding weighted annual estimates of 701,547 individuals or 2.4 individuals per 1000 population treated in emergency departments. On the other hand,

Lazarou et al. [3] estimated that each year 6-7% of hospitalized patients experience severe ADRs, which can lead to a potential of 100,000 deaths, making it the fourth largest cause of death in US. Over the past 10 years, both reported ADRs and related deaths have increased ~2.6 times and we have seen a number of drugs withdrew from the US market after presenting unexpected severe ADRs [4, 5]. As a consequence, ADRs not only expose patients to higher risks of mortality and loss of quality of life, but also presents a huge burden on the national economy with an estimated \$136 billion annual cost in the US, which is higher than cardiovascular or diabetic care [6, 7].

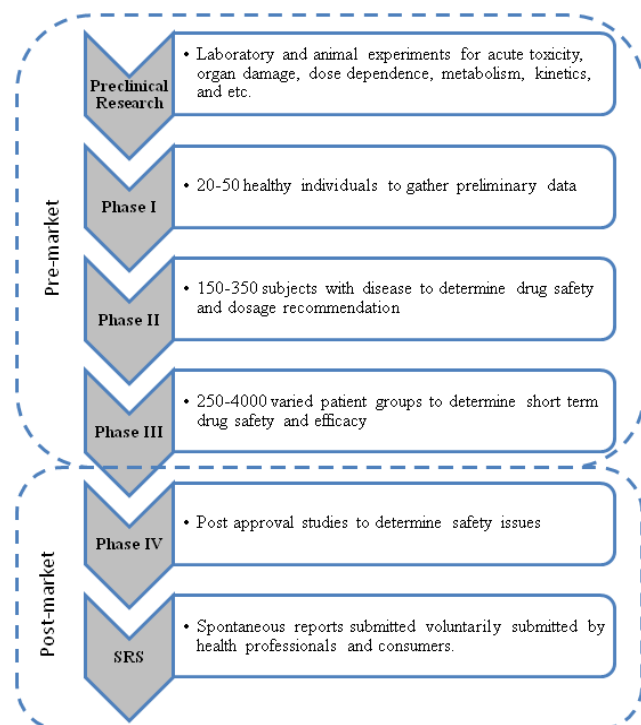
ADRs are also a big concern for the pharmaceutical industry. Drug discovery is a long and expensive process. To bring a new drug to market, it can take at least 10 years and billions of dollars [8]. The main cause is the high failure rate of drug candidates in clinical trials. Unacceptable toxicities account for approximately 30% of the failures [9]. Thus, early prediction of potential ADRs is essential to reduce risks of the costly failures. Additionally, even after a drug is approved to market, undiscovered severe ADRs may lead to withdrawals which can be detrimental financially for the manufacturers. Hence, it is critical to predict and monitor a drug's ADRs throughout its life cycle, from preclinical screening phases to post-market surveillance.<sup>1</sup>

Pharmacovigilance (PhV), also known as drug safety surveillance, is the science to enhance patient care and patient safety regarding the use of medicines by collecting, monitoring, assessing, and evaluating information from healthcare providers and patients. Broadly speaking, PhV can be divided into two stages: (1) pre-marketing surveillance – information regarding ADRs is collected from pre-clinical screening and phases I to III clinical trials; and (2) post-marketing surveillance – data accumulated in the post-approval stage and throughout a drug's market life (Figure 1).

Historically, PhV has relied on biological experiments or manual review of case reports; however, due to the vast quantities and complexity of data to be analyzed, computational methods that can accurately detect ADRs in a timely fashion have become a critical component in PhV. Large-scale compound databases containing structure, bioassay, and genomic information, such as NIH's Molecular Libraries Initiative [10], as well as

<sup>1</sup> Dr. Michael E. Matheny has a joint appointment in the Department of Biostatistics and Division of General Internal Medicine at Vanderbilt University. He is also affiliated with the Geriatric Research Education and Clinical Care in the Veterans Health Administration at Nashville, Tennessee.

comprehensive clinical data sets such as electronic medical record (EMR) databases, have become the enabling resources for computerized ADR detection methods.



**Figure 1.** Pharmacovigilance at different stages of drug development

In this paper, we will cover a broad spectrum of the current computational methodologies for PhV at both pre-marketing and post-marketing stages. The methodologies can be classified along different axes depending on the data sources applied with respect to each PhV stage.

## 2. PRE-MARKETING SURVEILLANCE

Much effort of PhV at the pre-marketing stage has been devoted to predict or assess potential ADRs early in the drug development pipeline. One of the fundamental methods is the application of preclinical *in vitro* Safety Pharmacology Profiling (SPP) by testing compounds with biochemical and cellular assays [11]. The hypothesis is that if a compound binds to a certain target, then its effect may translate into possible occurrence of an ADR in humans. However, experimental detection of ADRs remains challenging in terms of cost and efficiency [11]. There has been large amount of research activities devoted to developing computational approaches to predict potential ADRs using preclinical characteristics of the compounds or screening data. Most of the existing research can be categorized into protein target-based and chemical structure-based approaches. Others have also explored integrative approach.

### 2.1 Protein Target-based Approach

Drugs typically work by activating or inhibiting the function of a protein, which in turn results in therapeutic benefits to a patient. Thus, drug design essentially involves the design of small molecules that have complementary shapes and charges to the protein target with which they can bind and interact. ADRs are

complex phenomenological observations of drugs that have been attributed to a variety of molecular scenarios such as unexpected interaction with the primary or off-targets, downstream pathway perturbations, and kinetics [12]. Many believe direct interaction with proteins to be one of the most important scenarios [11, 13].

Fliri et al. [14] have shown that drugs with similar *in vitro* protein binding profiles tend to exhibit similar side-effects through hierarchical clustering of biological activity spectra and adverse event data of 1045 prescription drugs and 92 ligand-binding assays. This concept was further illustrated by Campillos et al. [15] where they extrapolated new drug targets by analyzing the likelihood of sharing protein targets for 277,885 pairs of 746 marketed drugs using their side-effect similarities.

Scheiber et al. [16] demonstrated the concept by comparing pathways affected by toxic compounds vs. those affected by non-toxic compounds. Fukuzaki et al. [17] proposed a method to predict ADRs using sub-pathways that share correlated modifications of gene-expression profiles in the presence of the drug of interest. To find the “cooperative pathways” (pathways that function together), they developed an algorithm called CoopeRativE Pathway Enumerator (CREPE) to select combinations of sub-pathways that have common activation conditions. Their work depends on the availability of gene-expression data observed under chemical perturbations by a drug.

Xie et al. [18] developed a chemical systems biology approach to identify off-targets of a drug by docking the drug into binding pockets of proteins that are similar to its primary target. Then the drug-protein interaction pair with the best docking score was mapped to known biological pathways to identify potential off-target binding networks of the drug. Unfortunately, scalability of the method is hindered by its requirement for protein 3D structures and known biological pathways.

More recently, Brouwers et al. [19] quantified the contribution of protein interaction network neighborhood on the observed side-effect similarity of drugs. Their fundamental idea is that side-effect similarity of drugs could be attributed to their target proteins being close in a molecular network. They proposed a pathway neighborhood measure to assess the closest distance of drug pairs according to their target proteins in the human protein-protein interaction network and found network neighborhoods to only account for 5.8% of the side-effect similarities compared to 64% by shared drug targets.

Pouliot et al. [20] applied logistic regression (LR) models to identify potential ADRs manifesting in 19 specific system organ classes (SOCs), as defined by the Medical Dictionary for Regulatory Activities [21], across 485 compounds in 508 BioAssays in the PubChem database [22, 23]. The models were evaluated using leave-one-out-cross-validation. The mean AUCs (area under the receiver operating characteristic curve) ranged from 0.60 to 0.92 across different SOCs.

### 2.2 Chemical Structure-based Approach

The chemical structure-based approach attempts to link ADRs to their chemical structures. Most notably, as a proof-of-concept, Bender et al. [24] explored the chemical space of drugs and established its correlation for ADR prediction; however, the positive predictive value was quite low, under 0.5. Thereafter, Scheiber et al. [25] presented a global analysis that identified chemical substructures associated with ADRs, but the method was

not designed to predict ADRs for any specific drug molecule. Yamanishi et al. [26] proposed a method that predicted pharmacological effects from chemical structures and then used the effect similarity to infer drug-target interactions.

Hammann et al. [27] employed decision tree to determine the chemical, physical, and structural properties of compounds that predispose them to causing ADRs. They focused on ADRs in the central nervous system (CNS), liver, and kidney as well as allergic ADRs for 507 compounds. The features used were numerical attributes computed from a compound's structure, which included elemental analysis (e.g., atom count), charge analysis (e.g., polarizability, ion charge, topological polar surface area), and geometry (e.g., number of aromatic rings, rotatable bonds), as well as partitioning coefficients and miscellaneous other characteristics (e.g., indicators of hydrogen bonding) [27]. Their decision tree model was shown to produce predictive accuracies ranging from 78.9 to 90.2% for allergic, renal, CNS, and hepatic ADRs.

Pauwels et al. [28] developed a sparse canonical correlation analysis (SCCA) method to predict high-dimensional side-effect profiles of drug molecules based on the chemical structures. They demonstrated the usefulness of SCCA by predicting 1385 side-effects in the SIDER database [29] from the chemical structures of 888 approved drugs. They compared five methods: random assignment (Random) as a baseline, nearest neighbor (NN), support vector machine (SVM), ordinary canonical correlation analysis (OCCA), and SCCA for their abilities to predict known side-effect profiles through 5-fold cross validation. The best resulting AUC scores are 0.6088, 0.8917, 0.8930, 0.8651, and 0.8932 for Random, NN, SVM, OCCA, SCCA, respectively. Their results suggest that the proposed method, SCCA, outperforms OCCA and its performance is comparable to SVM and NN. The main advantage of OCCA and SCCA over other algorithms is their biological interpretability to understand relationships between the chemical substructures and ADRs.

## 2.3 Integrative Approach

In the past year, approaches integrating various types of data relating to drugs for ADR prediction have gained many interests. Huang et al. [30] proposed a new computational framework to predict ADRs by integrating systems biology data that include protein targets, protein-protein interaction network, gene ontology (GO) annotation [31], and reported side effects. The SVM was applied as the predictive model to predict heart-related ADRs (i.e. cardio toxicity), which resulted in the highest AUC of 0.771. Soon after, Cami et al. [32] developed another ADR prediction framework by combining network structure formed by drug-ADR relationships (809 drugs and 852 ADRs) and information regarding specific drugs and adverse events. LR model was used as the predictive model and achieved an AUC of 0.87.

Despite the success of using chemical and biological information of drugs for ADR prediction, few studies have investigated the use of phenotypic information (e.g., indication and other known ADRs). Existing resources, such as the SIDER database [29], contain comprehensive drug phenotypic information, which has been demonstrated to be useful for other drug related studies [15]. Recently, Liu et al. [33] investigated the use of phenotypic information, together with chemical and biological properties of drugs, to predict ADRs. Similar to the work by Pauwels et al. [28], they conducted a large-scale study to develop and validate the ADR prediction model on 1385 known ADRs for 832 FDA (US

Food and Drug Administration) approved drugs in SIDER using five machine learning algorithms: LR, Naïve Bayes (NB), K-Nearest Neighbor (KNN), Random Forest (RF), and SVM. Evaluation results showed that the integration of chemical, biological, and phenotypic properties outperforms the chemical structured-based method (from 0.9054 to 0.9524 with SVM) and has the potential to detect clinically important ADRs at both preclinical and post-market phases for drug surveillance.

## 3. POST-MARKETING SURVEILLANCE

Although a drug undergoes extensive screening (Figure 1) before its approval by the FDA, many ADRs may still be missed because the clinical trials are often small, short, and biased by excluding patients with comorbid diseases. Premarketing trials do not mirror actual clinical use situations for diverse (e.g. inpatient) populations, thus it is important to continue the surveillance post-market. Several unique data sources are available for post-marketing PhV.

### 3.1 Spontaneous Reports

Spontaneous reporting systems (SRSs) have served as the core data-collection system for post-marketing drug surveillance since 1960. Some of the prominent SRSs are the Adverse Event Reporting System (AERS) maintained by the US FDA and the VigiBase managed by the World Health Organization (WHO). Although the SRSs may differ in structure and content, most of them rely on healthcare professionals and consumers to identify and report suspected cases of ADRs. Information collected usually include the drugs suspected to cause the ADR, concomitant drugs, indications, suspected events, and limited demographic information. Many post-marketing surveillance analyses are based on these reports voluntarily submitted to the national SRSs, which include disproportionality analysis and data mining algorithms.

#### 3.1.1 Disproportionality Analysis

Disproportionality analysis (DPA) has been the driving force behind most PhV methods involving SRS data. The first time use of DPA for drug safety can be dated back to the early 1980s [34]. It is not our intention to exhaustively list and examine all relevant work. Rather, we aim to present the basic concepts and highlight some representative work here. DPA involves frequency analyses of 2x2 contingency tables to quantify the degree to which a drug and ADR co-occurs "disproportionally" compared with what would be expected if there were no association (Table 1) [35].

|         | ADR       | No ADR | Total             |
|---------|-----------|--------|-------------------|
| Drug    | a         | b      | n = a + b         |
| No Drug | c         | d      | c + d             |
| Total   | m = a + c | b + d  | t = a + b + c + d |

**Table 1.** Contingency table used in DPA

Straightforward DPA methods involve the calculation of frequentist metrics. Some of the widely applied frequentist measures (Table 2) include the relative reporting ratio (RRR) [36], proportional reporting ratio (PRR) [37] adopted by the Medicines and Healthcare products regulatory Agency (MHRA) in UK and reporting odds ratio (ROR) [38] adopted by the Netherlands Pharmacovigilance Center. Hypothesis tests of independence (i.e., Chi-square test or Fisher's exact test) are typically used along

with the above association estimates as extra precautionary measures.

| Association Measures               | Definition                |
|------------------------------------|---------------------------|
| Relative Reporting Ratio (RRR)     | $(t * a) / (m * n)$       |
| Proportional Reporting Ratio (PRR) | $(a * (t - n)) / (c * n)$ |
| Reporting Odds Ratio (ROR)         | $(a * d) / (c * b)$       |

**Table 2.** Definitions of the frequentist measures of association

In addition to the frequentist approaches, more complex algorithms based on Bayesian statistics were developed such as the gamma-Poisson shrinker (GPS) [39], the multi-item gamma-Poisson shrinker (MGPS) [40, 41], and empirical Bayesian geometric means (EBGMs) [42, 43]. The GPS and MGPs methods are currently utilized by the FDA. Moreover, Bayesian Confidence Propagation Neural Network (BCPNN) [44-46] analysis was proposed based on Bayesian logic where the relation between the prior and posterior probability was expressed as the “information component (IC)”. The IC given by the BCPNN is applied by the WHO Uppsala Monitoring Center (UMC) to monitor safety signals in their SRSs.

Other groups have also investigated James-Stein type shrinkage estimation strategies in a Bayesian logistic regression model to analyze spontaneous adverse event reporting data [47]. More recently, Ahmed et al. [48, 49] proposed false discovery rate (FDR) estimation for the frequentist methods to address the limitation of arbitrary thresholds. As of now, there is no consensus on which DPA method is better because there is no gold standard dataset available to evaluate the performances of the methods.

### 3.1.2 Data Mining Algorithms

The above mentioned DPA methods are effective in detecting single Drug-ADR associations, but multi-item ADR associations are also important because they could suggest possible drug-drug interactions. A typical SRS database contains thousands of drugs and ADRs, so it is impractical to enumerate all combinations for statistical analysis. Thus, data mining algorithms have been employed to address this problem.

Harpaz et al. [50] applied the association rule mining algorithm to identify multi-item ADRs. Using a set of 162,744 reports submitted to the FDA in 2008, they identified 1167 multi-item ADR associations. Among those identified multi-item associations, 67% were validated by a domain expert. Later, Harpaz et al. [51] applied the biclustering algorithm to identify drug groups that share a common set of ADRs in SRS data. Tatonetti et al. [52] proposed an algorithm to mine drug-drug interactions from the adverse event reports by analyzing latent signals that indirectly provide evidence for ADRs. They discovered that co-administration of pravastatin and paroxetine had a synergistic effect on blood glucose. In contrast, neither drug individually was found to be associated with such change in the glucose levels.

## 3.2 Electronic Medical Records

Electronic medical records (EMRs) have emerged as a prominent resource for observational research as they contain not only detailed patient information but also copious longitudinal clinical data. Recently, investigators have begun to explore the use of EMRs for PhV. EMR databases consist of data in two types of

formats: (1) structured (e.g., laboratory data) and (2) narrative clinical notes.

### 3.2.1 Structured Data

Several groups have employed computational methods on structured or coded data in EMRs to identify specific ADR signals [53, 54]. Jin et al. [55] proposed a new interestingness measure called residual-leverage for association rule mining to identify ADR signals from healthcare administrative databases. Ji et al. [56] introduced potential causal association rules to generate potential causal relationships between a drug and ICD-9 coded signs or symptoms in EMRs. Schildcrout et al. [57] analyzed the relationship between insulin infusion rates and blood glucose levels in patients in an intensive care unit (ICU). Yoon et al. [58] demonstrated laboratory abnormality to be a valuable source for PhV by examining the odds ratio of laboratory abnormalities between a drug-exposed and a matched unexposed group using 10 years of EMR data. Evaluation of their algorithm on 470 randomly selected drug-and-abnormal-lab-event pairs produced a positive predictive value of 0.837 and negative predictive value of 0.659.

### 3.2.2 Unstructured Data

Data in narrative clinical notes is not readily accessible for data mining, thus natural language processing (NLP) technique is required to extract the needed information. Wang et al. [59] first employed NLP techniques to extract drug-ADR candidate pairs from narrative EMRs and then applied the Chi-square test with adjusted volume test to detect ADR signals. Evaluation on 7 selected drugs and their known ADRs produced an overall precision and recall of 0.31 and 0.75 respectively.

Similarly, Wang et al. [60] developed other methods based on mutual information (MI) and data processing inequality (DPI) to characterize drug-and-ADR pairs extracted from EMRs. Evaluation on a random sample of two drugs and two diseases indicated an overall precision of 81%. Furthermore, Wang et al. [61] investigated the use of filtering by sections of reports to improve the performance of NLP extraction for clinically meaningful drug-and-ADR relations. Their evaluation indicated that applying filters improved recall from 0.43 to 0.75 and precision from 0.16 to 0.31.

## 3.3 Non-conventional Data Sources

### 3.3.1 Biomedical Literature

Biomedical literature can be used as a complementary resource for prioritizing drug-ADR associations generated from SRSs. Shetty and Dalal [62] retrieved articles (published between 1949 and 2009) that contain mentions of a pre-defined list of drug-and-ADR pairs (38 drugs and 55 ADRs) from PubMed. The authors then constructed a statistical document classifier to remove irrelevant articles with mentions of treatment relations. Finally, DPA was applied to identify statistically significant pairs from the thousands of pairs in the remaining articles. Evaluation showed that the method identified true associations with 0.41 and 0.71 in precision and recall, respectively.

### 3.3.2 Health Forums

Data posted by users on health-related websites may also contain valuable drug safety information. Leaman et al. [63] described a system to mine drug-and-ADR relationships as reported by consumers in user comments to health-related websites like

DailyStrength (<http://www.dailystrength.org/>). System evaluation was conducted on a manually annotated set of 3600 user posts corresponding to 6 drugs. The system was shown to achieve 0.78 in precision and 0.70 in recall.

Chee et al. [64] explored the use of ensemble classifier over data from online health forums to identify potential watchlist drugs that have an active FDA safety alert. The authors aggregated individuals' opinions and review of drugs and used NLP technique to group drugs that were discussed in similar ways. Interestingly, withdrawn drugs were successfully identified based on messages even before they were removed from market.

## 4. FUTURE PERSPECTIVES

In this paper, we have provided a general overview of the rich and diverse applications of computational approaches with respect to different perspectives of PhV. More and more opportunities have emerged as a result of new data generated from various platforms including EMRs, literature, and self-reported health forums.

It is evident that a new trend of computational approaches for PhV is to link preclinical data from the experimental platform with human safety information observed in the post-marketing phase [65]. From the systems biology perspective, drugs are considered as molecules that induce perturbations to biological systems, which involve various molecular interactions such as protein-protein interactions, signaling pathways, and pathways of drug action and metabolism. When a drug is absorbed into the body and interacts with its intended targets, favorable effects are expected. However, a drug often binds to other protein pockets with varying affinities (off-target interactions), leading to observed side-effects. Thus, the body's response to a drug is a complex phenomenological observation that includes both the favorable and unfavorable reactions. Hence, it is desirable to incorporate various data sources into one framework to understand ADRs.

Moreover, it is essential to identify multi-item ADR associations as they may suggest drug interactions. Drug interactions are extremely important. For example, if a patient is taking two drugs and one of them increases the effect of the other, then the patient may have an overdose. Similarly, if the action of a drug is inhibited, it may reduce therapeutic effect. Drug interactions may also increase the risk of ADRs. Statistical analysis works well with the identification of single drug-and-ADR signals, but not suitable for drug interaction identification. Alternatively, data mining algorithms such as *a priori* algorithm and clustering algorithms are applicable and useful. It provides an excellent opportunity for computer scientists to develop new algorithms for drug interaction detection.

Furthermore, EMRs have become an obvious data choice for PhV. Many challenges exist in mining EMR for ADR prediction. Much detailed and useful information is embedded in the narrative notes making data extraction difficult. There have been studies using NLP techniques to extract drug and ADR concepts from narrative notes for association analysis. Wang et al. [61] have shown that filtering information based on note sections improves the identification of drug-and-ADR relations. Despite the current success, further investigation of other methods, for example more sophisticated statistical methods and temporal models, is needed.

As yet, few studies have explored the automatic construction of large cohort or case-control studies from EMRs for ADR

prediction. There are many issues to consider in the NLP-based cohort/case-control study construction. For instance, how to extract event concepts from narrative notes? It is common for multiple concepts to describe the same outcome/phenotype. Since most current practices focus on single outcome at a time, phenotype is usually defined manually by experts. However, for large-scale ADR studies, how to automatically define the phenotypes? Also, how to accurately determine the time-relations between events in the narrative text? Each of these questions is an active area of research.

After overcoming the above hurdles in study design, one must keep in mind of the confounding problem during analysis. For instance, the basic concept behind the cohort design is to partition a population into those who are "exposed" (taking a specific drug) and "unexposed" (taking a comparator drug or not taking a specific drug). A drug is determined to be associated with a specified outcome when the outcome occurs more often in exposed group than in the unexposed group. Since the group assignment is not random, increased attention must be given when selecting the 'unexposed' group. A common technique to minimize the issues caused by confounding and bias is to match patient groups based on a set of basic covariates such as gender, age, and comorbidities. On the other hand, case-control designs divide the study population into those who experienced the outcome ("case") and those who did not experience the outcome ("control"). If the drug exposure occurs more frequently in the cases than in the controls, the drug is said to be associated with the outcome. The same issues with confounding apply to the case-control studies. Matching two groups before analysis is usually a good idea.

Lastly, it is important to note that most of the existing methodologies for PhV involve assessment of association between a drug and ADR. However, association does not necessarily imply causation. Intuitively, causation not only requires correlation but also a counterfactual dependence. Inferring cause-and-effect relationships is an intrinsically hard problem in data mining and need to be further investigated for the PhV application.

## 5. ACKNOWLEDGMENTS

Dr. Michael Matheny is supported by a Veterans Administration HSR&D Career Development Award (CDA-08-020). Dr. Yong Hu is partly supported by the National Natural Science Foundation of China (NSFC, project no.: 70801020) and the Science and Technology Planning Project of Guangdong Province, China (project no.: 2010B010600034). Dr. Hua Xu is supported by grants from NLM R01-LM007995 and NCI R01CA141307.

## 6. REFERENCES

- [1] Pirmohamed, M., Breckenridge, A.M., Kitteringham, N.R. and Park, B.K. Adverse drug reactions. *BMJ*, 316, 7140 (Apr 25 1998), 1295-1298.
- [2] Budnitz, D.S., Pollock, D.A., Weidenbach, K.N., Mendelsohn, A.B., Schroeder, T.J. and Anest, J.L. National surveillance of emergency department visits for outpatient adverse drug events. *JAMA*, 296, 15 (Oct 18 2006), 1858-1866.
- [3] Lazarou, J., Pomeranz, B.H. and Corey, P.N. Incidence of adverse drug reactions in hospitalized patients: a meta-

- analysis of prospective studies. *JAMA*, 279, 15 (Apr 15 1998), 1200-1205.
- [4] Moore, T.J., Cohen, M.R. and Furberg, C.D. Serious adverse drug events reported to the Food and Drug Administration, 1998-2005. *Arch Intern Med*, 167, 16 (Sep 10 2007), 1752-1759.
- [5] Giacomini, K.M. Krauss, R.M., Roden, D.M., Eichelbaum, M., Hayden, M.R. and Nakamura, Y. when good drugs go bad. *Nature*, 446, 7139 (Apr 26 2007), 975-977.
- [6] Leone, R., Sottosanti, L., Luisa Iorio, M. Santuccio, C., Conforti, A., Sabatini, V., Moretti, U. and Venegoni, M. Drug-related deaths: an analysis of the Italian spontaneous reporting database. *Drug Saf*, 31, 8 (2008), 703-713.
- [7] van der Hooft, C.S. Sturkenboom, M.C., van Grootheest, K., Kingma, H.J. and Stricker, B.H. Adverse drug reaction-related hospitalizations: a nationwide study in the Netherlands. *Drug Saf*, 29, 2 (2006), 161-168.
- [8] Paul, S.M., Mytelka, D.S., Dunwiddie, C.T., Persinger, C.C., Munos, B.H., Lindborg, S.R. and Schacht, A.L. How to improve R&D productivity: the pharmaceutical industry's grand challenge. *Nat Rev Drug Discov*, 9, 3 (Mar 2010), 203-214.
- [9] Hopkins, A.L. Network pharmacology: the next paradigm in drug discovery. *Nat Chem Biol*, 4, 11 (Nov 2008), 682-690.
- [10] Austin, C. P., Brady, L. S., Insel, T. R. and Collins, F. S. NIH Molecular Libraries Initiative. *Science*, 306, 5699 (Nov 12 2004), 1138-1139.
- [11] Whitebread, S., Hamon, J., Bojanic, D. and Urban, L. Keynote review: in vitro safety pharmacology profiling: an essential tool for successful drug development. *Drug Discov Today*, 10, 21 (Nov 1 2005), 1421-1433.
- [12] Liebler, D.C. and Guengerich, F.P. Elucidating mechanisms of drug-induced toxicity. *Nat Rev Drug Discov*, 4, 5 (May 2005), 410-420.
- [13] Blagg, J. Structure-activity relationships for in vitro and in vivo toxicity. *Annu Rep Med Chem*, 41 (2006), 353-368.
- [14] Fliri, A.F., Loging, W.T., Thadeio, P.F. and Volkmann, R.A. Analysis of drug-induced effect patterns to link structure and side effects of medicines. *Nat Chem Biol*, 1, 7 (Dec 2005), 389-397.
- [15] Campillos, M., Kuhn, M., Gavin, A.C., Jensen, L.J. and Bork, P. Drug target identification using side-effect similarity. *Science*, 321, 5886 (Jul 11 2008), 263-266.
- [16] Scheiber, J., Chen, B., Milik, M., Sukuru, S.C., Bender, a., Mikhailov, D., Whitebread, S., Hamon, J., Azzaoui, K., Urban, L., Glick, M., Davies, J.W. and Jenkins, J.L. Gaining insight into off-target mediated effects of drug candidates with a comprehensive systems chemical biology analysis. *J Chem Inf Model*, 49, 2 (Feb 2009), 308-317.
- [17] Fuzuzaki, M., Seki, M., Kashima, H. and Sese, J. Side effect prediction using cooperative pathways. In *Proceedings of IEEE International Conference on Bioinformatics and Biomedicine '09* (Washington DC, 2009), 142-147.
- [18] Xie, L., Li, J. and Bourne, P.E. Drug discovery using chemical systems biology: identification of the protein-ligand binding network to explain the side effects of CETP inhibitors. *PLoS Comput Biol*, 5, 5 (May 2009), e1000387.
- [19] Brouwers, L., Iskar, M., Zeller, G., van Noort, V. and Bork, P. Network neighbors of drug targets contribute to drug side-effect similarity. *PLoS One*, 6, 7 (Jul, 2011), e22187.
- [20] Pouliot, Y., Chiang, A. P. and Butte, A. J. Predicting adverse drug reactions using publicly available PubChem BioAssay data. *Clin Pharmacol Ther*, 90, 1 (Jul 2011), 90-99.
- [21] Brown, E. G., Wood, L. and Wood, S. The medical dictionary for regulatory activities (MedDRA). *Drug Safety*, 20, (1999), 109-117.
- [22] Chen, B., Wild, D. and Guha, R. PubChem as a source of polypharmacology. *J Chem Inf Model*, 49, 9 (Sep 2009), 2044-2055.
- [23] Bolton, E., Wang, Y., Thiessen, P. A. and Bryant, S. H. PubChem: integrated platform of small molecules and biological activities. American Chemical Society, City, 2008.
- [24] Bender, A., Scheiber, J., Glick, M., Davies, J.W., Azzaoui, K., Hamon, J., Urban, L., Whitebread, S. and Jenkins, J.L. Analysis of pharmacology data and the prediction of adverse drug reactions and off-target effects from chemical structure. *ChemMedChem*, 2, 6, (Jun 2007), 861-873.
- [25] Scheiber, J., Jenkins, J.L., Sukuru, S.C., Bender, A., Mikhailov, D., Milik, M., Azzaoui, K., Whitebread, S., Hamon, J., Urban, L., Glick, M. and Davies, J.W. Mapping adverse drug reactions in chemical space. *J Med Chem*, 52, 9 (May 14 2009), 3103-3107.
- [26] Yamanishi, Y., Kotera, M., Kanehisa, M. and Goto, S. Drug-target interaction prediction from chemical, genomic and pharmacological data in an integrated framework. *Bioinformatics*, 26, 12 (Jun 15 2010), i246-254.
- [27] Hammann, F., Gutmann, H., Vogt, N., Helma, C. and Drewe, J. Prediction of adverse drug reactions using decision tree modeling. *Clin Pharmacol Ther*, 88, 1 (Jul 2010), 52-59.
- [28] Pauwels, E., Stoven, V. and Yamanishi, Y. Predicting drug side-effect profiles: a chemical fragment-based approach. *BMC Bioinformatics*, 12(2011), 169.
- [29] Kuhn, M., Campillos, M., Letunic, I., Jensen, L. J. and Bork, P. A side effect resource to capture phenotypic effects of drugs. *Mol Syst Biol*, 6(2010), 343.
- [30] Huang, L. C., Wu, X. and Chen, J. Y. Predicting adverse side effects of drugs. *BMC Genomics*, 12 Suppl 5(Dec 23 2011), S11.
- [31] Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., Davis, A. P., Dolinski, K., Dwight, S. S., Eppig, J. T., Harris, M. A., Hill, D. P., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese, J. C., Richardson, J. E., Ringwald, M., Rubin, G. M. and Sherlock, G. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet*, 25, 1 (May 2000), 25-29.
- [32] Cami, A., Arnold, A., Manzi, S. and Reis, B. Predicting adverse drug events using pharmacological network models. *Sci Transl Med*, 3, 114 (Dec 21 2011), 114ra127.
- [33] Liu, M., Wu, Y., Chen, Y., Sun, J., Zhao, Z., Chen, X. W., Matheny, M. E. and Xu, H. Large-scale Prediction of



Adverse Drug Reactions Using Chemical, Biological, and Phenotypic Properties of Drugs. *J Am Med Inform Assoc*, 19, (2012), e28-e35.

- [34] Montastruc, J. L., Sommet, A., Bagheri, H. and Lapeyre-Mestre, M. Benefits and strengths of the disproportionality analysis for identification of adverse drug reactions in a pharmacovigilance database. *Br J Clin Pharmacol*, 72, 6 (Dec 2011), 905-908.
- [35] Bate, A. and Evans, S. J. Quantitative signal detection using spontaneous ADR reporting. *Pharmacoepidemiol Drug Saf*, 18, 6 (Jun 2009), 427-436.
- [36] Hauben, M., Madigan, D., Gerrits, C. M., Walsh, L. and Van Puijenbroek, E. P. The role of data mining in pharmacovigilance. *Expert Opin Drug Saf*, 4, 5 (Sep 2005), 929-948.
- [37] Evans, S. J., Waller, P. C. and Davis, S. Use of proportional reporting ratios (PRRs) for signal generation from spontaneous adverse drug reaction reports. *Pharmacoepidemiol Drug Saf*, 10, 6 (Oct-Nov 2001), 483-486.
- [38] Szarfman, A., Machado, S. G. and O'Neill, R. T. Use of screening algorithms and computer systems to efficiently signal higher-than-expected combinations of drugs and events in the US FDA's spontaneous reports database. *Drug Saf*, 25, 6 (2002), 381-392.
- [39] Ahmed, I., Haramburu, F., Fourrier-Reglat, A., Thiessard, F., Kreft-Jais, C., Miremont-Salame, G., Begaud, B. and Tubert-Bitter, P. Bayesian pharmacovigilance signal detection methods revisited in a multiple comparison setting. *Stat Med*, 28, 13 (Jun 15 2009), 1774-1792.
- [40] DuMouchel, W. Bayesian data mining in large frequency tables, with an application to the FDA spontaneous reporting system. *The American Statistician*, 53, 3 (1999), 177-202.
- [41] Almenoff, J. S., Pattishall, E. N., Gibbs, T. G., DuMouchel, W., Evans, S. J. and Yuen, N. Novel statistical tools for monitoring the safety of marketed drugs. *Clin Pharmacol Ther*, 82, 2 (Aug 2007), 157-166.
- [42] DuMouchel, W., Smith, E. T., Beasley, R., Nelson, H., Yang, X., Fram, D. and Almenoff, J. S. Association of asthma therapy and Churg-Strauss syndrome: an analysis of postmarketing surveillance data. *Clin Ther*, 26, 7 (Jul 2004), 1092-1104.
- [43] Gould, A. L. Accounting for multiplicity in the evaluation of "signals" obtained by data mining from spontaneous report adverse event databases. *Biom J*, 49, 1 (Feb 2007), 151-165.
- [44] Bate, A., Lindquist, M., Edwards, I. R., Olsson, S., Orre, R., Lansner, A. and De Freitas, R. M. A Bayesian neural network method for adverse drug reaction signal generation. *Eur J Clin Pharmacol*, 54, 4 (Jun 1998), 315-321.
- [45] Lindquist, M., Edwards, I. R., Bate, A., Fucik, H., Nunes, A. M. and Stahl, M. From association to alert--a revised approach to international signal analysis. *Pharmacoepidemiol Drug Saf*, 8 Suppl 1 (Apr 1999), S15-25.
- [46] Lindquist, M., Stahl, M., Bate, A., Edwards, I. R. and Meyboom, R. H. A retrospective evaluation of a data mining approach to aid finding new adverse drug reaction signals in the WHO international database. *Drug Saf*, 23, 6 (Dec 2000), 533-542.
- [47] An, L., Fung, K. Y. and Krewski, D. Mining pharmacovigilance data using Bayesian logistic regression with James-Stein type shrinkage estimation. *J Biopharm Stat*, 20, 5 (Sep 2010), 998-1012.
- [48] Ahmed, I., Dalmasso, C., Haramburu, F., Thiessard, F., Broet, P. and Tubert-Bitter, P. False discovery rate estimation for frequentist pharmacovigilance signal detection methods. *Biometrics*, 66, 1 (Mar 2010), 301-309.
- [49] Ahmed, I., Thiessard, F., Miremont-Salame, G., Begaud, B. and Tubert-Bitter, P. Pharmacovigilance data mining with methods based on false discovery rates: a comparative simulation study. *Clin Pharmacol Ther*, 88, 4 (Oct 2010), 492-498.
- [50] Harpaz, R., Chase, H. S. and Friedman, C. Mining multi-item drug adverse effect associations in spontaneous reporting systems. *BMC Bioinformatics*, 11 Suppl 92010), S7.
- [51] Harpaz, R., Perez, H., Chase, H. S., Rabadan, R., Hripcsak, G. and Friedman, C. Biclustering of adverse drug events in the FDA's spontaneous reporting system. *Clin Pharmacol Ther*, 89, 2 (Feb 2011), 243-250.
- [52] Tatonetti, N. P., Denny, J. C., Murphy, S. N., Fernald, G. H., Krishnan, G., Castro, V., Yue, P., Tsao, P. S., Kohane, I., Roden, D. M. and Altman, R. B. Detecting drug interactions from adverse-event reports: interaction between paroxetine and pravastatin increases blood glucose levels. *Clin Pharmacol Ther*, 90, 1 (Jul 2011), 133-142.
- [53] Brown, J. S., Kulldorff, M., Chan, K. A., Davis, R. L., Graham, D., Pettus, P. T., Andrade, S. E., Raebel, M. A., Herrinton, L., Roblin, D., Boudreau, D., Smith, D., Gurwitz, J. H., Gunter, M. J. and Platt, R. Early detection of adverse drug events within population-based health networks: application of sequential testing methods. *Pharmacoepidemiol Drug Saf*, 16, 12 (Dec 2007), 1275-1284.
- [54] Berlowitz, D. R., Miller, D. R., Oliveria, S. A., Cunningham, F., Gomez-Caminero, A. and Rothendler, J. A. Differential associations of beta-blockers with hemorrhagic events for chronic heart failure patients on warfarin. *Pharmacoepidemiol Drug Saf*, 15, 11 (Nov 2006), 799-807.
- [55] Jin, H. D., Chen, J., He, H. X., Williams, G. J., Kelman, C. and O'Keefe, C. M. Mining unexpected temporal associations: Applications in detecting adverse drug reactions. *Ieee T Inf Technol B*, 12, 4 (Jul 2008), 488-500.
- [56] Ji, Y. Q., Ying, H., Dews, P., Mansour, A., Tran, J., Miller, R. E. and Massanari, R. M. A Potential Causal Association Mining Algorithm for Screening Adverse Drug Reactions in Postmarketing Surveillance. *Ieee T Inf Technol B*, 15, 3 (May 2011), 428-437.
- [57] Schildcrout, J. S., Haneuse, S., Peterson, J. F., Denny, J. C., Matheny, M. E., Waitman, L. R. and Miller, R. A. Analyses of longitudinal, hospital clinical laboratory data with application to blood glucose concentrations. *Stat Med*, 30, 27 (Nov 30 2011), 3208-3220.
- [58] Yoon, D., Park, M. Y., Choi, N. K., Park, B. J., Kim, J. H. and Park, R. W. Detection of Adverse Drug Reaction Signals

Using an Electronic Health Records Database: Comparison of the Laboratory Extreme Abnormality Ratio (CLEAR) Algorithm. *Clin Pharmacol Ther*(Jan 11 2012).

- [59] Wang, X., Hripcsak, G., Markatou, M. and Friedman, C. Active computerized pharmacovigilance using natural language processing, statistics, and electronic health records: a feasibility study. *J Am Med Inform Assoc*, 16, 3 (May-Jun 2009), 328-337.
- [60] Wang, X., Hripcsak, G. and Friedman, C. Characterizing environmental and phenotypic associations using information theory and electronic health records. *BMC Bioinformatics*, 10 Suppl 9(2009), S13.
- [61] Wang, X., Chase, H., Markatou, M., Hripcsak, G. and Friedman, C. Selecting information in electronic health records for knowledge acquisition. *J Biomed Inform*, 43, 4 (Aug 2010), 595-601.
- [62] Shetty, K. D. and Dalal, S. R. Using information mining of the medical literature to improve drug safety. *J Am Med Inform Assoc*, 18, 5 (Sep-Oct 2011), 668-674.
- [63] Leaman, R., Wojtulewicz, L., Sullivan, R., Skariah, A., Yang, J. and Gonzalez, G. Towards internet-age pharmacovigilance: extracting adverse drug reactions from user posts to health-related social networks. *Proceedings of the 2010 workshop on Biomedical Natural Language Processing*, (2010), 117-125.
- [64] Chee, B. W., Berlin, R. and Schatz, B. Predicting adverse drug events from personal health messages. *AMIA Annu Symp Proc* (2011), Washington DC, 217-226.
- [65] Harpaz, R., Dumouchel, W., Shah, N. H., Madigan, D., Ryan, P. and Friedman, C. Novel Data-Mining Methodologies for Adverse Drug Event Discovery and Analysis. *Clin Pharmacol Ther*(May 2 2012).