



HAL
open science

Detection and analysis of drug non-compliance in internet fora using information retrieval approaches

Lise Bigeard, Frantz Thiessard, Natalia Grabar

► To cite this version:

Lise Bigeard, Frantz Thiessard, Natalia Grabar. Detection and analysis of drug non-compliance in internet fora using information retrieval approaches. CICLING 2019, Apr 2019, La Rochelle, France. hal-02430414

HAL Id: hal-02430414

<https://hal.science/hal-02430414>

Submitted on 7 Jan 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Detection and analysis of drug non-compliance in internet fora using information retrieval approaches

lise Bigeard^a, Frantz Thiessard^b Natalia Grabar^a

¹ CNRS, Univ Lille, UMR 8163 STL - Savoirs Textes Langage, F-59000 Lille, France

² U Bordeaux, Inserm, Bordeaux Population Health Research Center, team ERIAS,
UMR 1219, F-33000 Bordeaux, France

Abstract. In the health-related field, drug non-compliance situations happen when patients do not follow their prescriptions and do actions which lead to potentially harmful situations. Although such situations are dangerous, patients usually do not report them to their physicians. Hence, it is necessary to study other sources of information. We propose to study online health fora with information retrieval methods in order to identify messages that contain drug non-compliance information. Exploitation of information retrieval methods permits to detect non-compliance messages with up to 0.529 F-measure, compared to 0.824 F-measure reached with supervised machine learning methods. For some fine-grained categories and on new data, it shows up to 0.70 Precision.

1 Introduction

In the health-related field, drug non-compliance situations happen when patients do not follow indications given by their doctors and by prescriptions. Typically, patients may decide to change the dosage, to refuse to take prescribed drugs, to take drugs without prescriptions, etc.

Misuse of drugs is a specific case of non-compliance, in which the drug is used by patients with a different intent than the one that motivated the prescription. Recreational medication use or suicide attempts are examples of drug misuse, but it has been noticed that the typology of drug misuses includes several other harmful situations [1]. Although these situations are dangerous for patients and their health, they do not inform their doctors that they are not following the prescription instructions. For this reason, the misuse situations become harder to detect and to prevent. Hence, there is the necessity to study other sources of information. We propose to contribute to this research question and to study social media messages, in which patients willingly and without any particular effort talk about their health and their practices regarding drug use [2].

Currently, social media has become an important source of information for various research areas, such as geolocalisation, opinion mining, event extraction, translation, or automatic summarizing [3].

In the medical domain, social media has been efficiently exploited in information retrieval for epidemiological surveillance [4, 5], study of patient’s quality of life [6], or drug adverse effects [7].

Yet, up to now, few works are focused on drug misuse and non-compliance. We can mention non-supervised analysis of tweets about non-medical use of drugs [8], and creation of a semantic web platform on drug abuse [9]. Both of these works are dedicated to one specific case of non-compliance, which is the drug abuse.

In our work, we propose to address a larger set of non-compliance situations. For this, we propose to exploit information retrieval methods. Our objective is to identify messages related to these situations in health fora in French.

In what follows, we first introduce the methods. We will first describe the method designed for the use of information retrieval system for supervised classification of messages and will compare the results obtained with the state of the art machine learning approaches. Then, we will use the information retrieval system on new non-annotated data and on different topics, and assess whether non-compliance messages are correctly detected. The discussion of the results is proposed. Finally, we conclude with directions for future research.

2 Methods

We propose to adapt an existing information retrieval system, using reference data, to perform the categorization of health fora messages in order to detect messages containing non-compliance information. We will compare these results with the state of the art machine learning methods. Then, we will use information retrieval in a fully unsupervised way to explore new non-annotated data on several topics related to non-compliance. The purpose is to assess if non-compliance messages can be found with information retrieval methods.

We first introduce our reference data and then describe the two ways to use information retrieval methods. We also describe the supervised machine learning approach.

2.1 Reference and Test Data

The reference and test data are built from corpora collected from several health fora written in French:

- Doctissimo³ is a well known health website to French-speaking people. It is widely used by people with punctual health questions. The main purpose of this website is to provide the platform to patients and to enable them to have discussions with other patients or their relatives. We collected messages from several Doctissimo fora: pregnancy, general drug-related questions, back pain, accidents in sport activities, diabetes. We collected messages written between 2010 and 2015;

³ <http://forum.doctissimo.fr>

- AlloDocteur⁴ is a question/answer health service, in which patients can ask questions which are answered by real doctors;
- masante.net⁵ is another question/answer service, in which the answers are provided by real doctors;
- Les diabétiques⁶ is a specialized forum related to diabetes.

In all these fora, the contributors are mainly sick persons and their relatives, who join the community to ask questions or provide accounts on their disorders, treatments, etc. These people may be affected by chronic disorders or present punctual health problems.

To build the reference data we use two fora from Doctissimo (pregnancy and general drug-related questions). We keep only messages that mention at least one drug. This gives a total of 119,562 messages (15,699,467 words). For the test data, we collect 145,012 messages from other corpora. Messages longer than 2,500 characters are excluded because they provide heterogeneous content difficult to categorize and process, both manually and automatically. The drug names are detected with specific vocabulary containing French commercial drug names built from several sources: *base CNHIM Thriaque*⁷, *base publique du médicament*⁸, and *base Medic'AM* from *Assurance Maladie*⁹. Each drug name is associated with the corresponding ATC code [10].

The reference data is manually annotated. Three annotators were asked to assign each message to one of the two categories:

- *non-compliance* category contains messages which report on drug misuse or non-compliance. When this category is selected, the annotators are also asked to shortly indicate what type of non-compliance is concerned (overuse, dosage change, brutal quitting...). This indication is written as free text with no defined categories. For instance, the following example shows non-compliance situation due to the forgotten intake of medication: *"bon moi la miss boulette et la tete en l'air je devais commencer mon "utrogestran 200" a j16 bien sur j'ai oublier! donc je l'ai pris ce soir!!!!" (well me miss blunder and with the head in the clouds I had to start the "utrogestran 200" at d16 and I forgot of course! so I took it this evening!!!!)*
- *compliance* category contains messages reporting normal drug use (*"Mais la question que je pose est 'est ce que c'est normal que le loxapac que je prends met des heures agir ???" (Anyway the question I'm asking is whether it is normal that loxapac I'm taking needs hours to do someting???)*) and messages without use of drugs (*"ouf boo, repose toi surtout, il ne t'a pas prescrit d'aspegic nourisson???" (ouch boo, above all take a break, he didn't prescribe aspegic for the baby???)*)

⁴ <http://www.allodocteur.fr>

⁵ <http://ma-sante.net>

⁶ <http://www.lesdiabetiques.com>

⁷ <http://www.theriaque.org>

⁸ <http://base-donnees-publique.medicaments.gouv.fr>

⁹ <https://www.ameli.fr/l-assurance-maladie/statistiques-et-publications/donnees-statistiques/medicament/medic-am/medic-am-mensuel-2017.php>

When annotators are unable to decide, they can mark up the corresponding messages accordingly. The categorization of these messages, as well as the categorization of annotation disagreements, are discussed later. The three annotators involved in the process are: one medical expert in pharmacology and two computer scientists familiar with medical texts and annotation tasks. Because this kind of annotation is a complicated task, especially concerning the decision on drug non-compliance, all messages annotated as non-compliant are additionally verified by one of the annotators.

The manual annotation process permitted to double-annotate 1,850 messages, among which we count 1,717 messages in the *compliance* category and 133 messages in the *non-compliance* category. These numbers indicate the natural distribution of *non-compliance* messages.

Concerning the annotation into *compliance* and *non-compliance* categories, the inter-annotator agreement [11] is 0.46, which is a moderate agreement [12]. This indicates that is difficult categorisation task.

Within the *non-compliance* category, we count 16 types of non-compliance: they contain between 1 and 29 messages. As example, the *change of weight* type contains 2 messages, *recreational use of drugs* 2 messages, *suicide attempt* 2 messages, and *overuse* 20 messages.

We see that due to the small number of non-adherence messages available and the multiplicity of the types of non-compliance, it may be difficult to obtain sufficient reference data to be used with supervised methods. For these reasons, we propose to use non-supervised methods for the detection of messages with sub-categories of drug non-compliance.

The corpus is pre-processed using Treetagger [13] for tokenization, POS-tagging, and lemmatization. The corpus is used in three versions: (1) in the *forms* corpus, the messages are only tokenized and lowercased; (2) in the *lemmas* corpus, the messages are also lemmatized, the numbers are replaced by a unique placeholder, and diacritics are removed such as in *anxit/anxiete (anxiety)*; (3) in the *lexical lemmas* corpus, we keep only lemmas of lexical words (verbs, nouns, adjectives, and adverbs). Besides, in each message, the drugs are indexed with the three first characters of the ATC categories [10].

2.2 Categorization with Information Retrieval

We use the Indri information retrieval system[14] to detect non-compliance messages following two ways:

- First, we use the annotated reference data to distinguish between drug compliant and non-compliant messages. The corpus is split in two sets:
 1. 44 *non-compliance* messages (one third of the whole *non-compliance* category) are used for the creation of queries, and the query lexicon is weighted proportionally to its frequency in the messages;
 2. All 98 *compliance* messages and 89 *non-compliance* messages (two thirds of the whole *non-compliance* category) are used for the evaluation.

The question we want to answer is whether the existing subset of *non-compliance* messages permits to retrieve other *non-compliance* messages. The evaluation is done automatically, using the reference data and computing Precision, Recall, and F-measure with each version of the corpus (forms, lemmas and lexical lemmas). This experiment may give an idea of the performance of this method when searching information in new non-annotated data;

- Then, we look for specific types of non-compliance without using the manual annotations. This exploitation of information retrieval system is fully unsupervised. At the previous steps, we discovered several types of non-compliance that can be found in our corpus. Now, we propose to exploit the existing reference data and our knowledge of the corpus gained at previous steps in order to create the best queries for the detection of similar messages. These queries are applied to a larger corpus with 20,000 randomly selected messages that contain at least one mention of drugs. The results are evaluated manually computing the Precision.

2.3 Supervised Categorization with Machine Learning

We also perform automatic detection of non-compliant messages with machine learning algorithms. The goal of this step is to provide a state of the art evaluation against which the results provided by the information retrieval methods can be compared.

Supervised machine learning algorithms learn a language model from manually annotated data, which can then be applied to new and unseen data. The categories are drug *compliance* and *non-compliance*. The unit processed is the message. The features are the vectorized text of messages (forms, lemmas and lexical lemmas) and the ATC indexing of drugs. The train set contains 94 non-compliant messages and 93 compliant messages. The test set contains 39 non-compliant messages and 40 compliant messages.

We use the Weka [15] implementation of several supervised algorithms: Naive-Bayes [16], Bayes Multinomial [17], J48 [18], Random Forest [19], and Simple Logistic [20]. These algorithms are used with their default parameters and with the string to wordvector function.

3 Results and Discussion

We present three sets of results obtained when detecting messages with the drug non-compliant information: (1) supervised categorization and evaluation of messages using the information retrieval system (Section 3.1); (2) supervised categorization of messages using the machine learning algorithms (Section 3.2) (3) unsupervised retrieval of messages using the information retrieval system (Section 3.3). We also indicate some limitations of the current work (Section 3.4).

3.1 Supervised Categorization and Evaluation with Information Retrieval

Table 1. Information retrieval results for the categorization of messages into the *non-compliance* category

	<i>Precision Recall F-measure</i>		
	<i>Top 10 results</i>		
<i>forms</i>	0.100	0.011	0.020
<i>lemmas</i>	0.400	0.045	0.081
<i>lexical lemmas</i>	0.400	0.045	0.081
	<i>Top 100 results</i>		
<i>forms</i>	0.480	0.539	0.508
<i>lemmas</i>	0.480	0.539	0.508
<i>lexical lemmas</i>	0.500	0.561	0.529

The results obtained with the information retrieval system Indri are presented in Table 1. The evaluation values are computed against the reference data for the top 10, 20, 50 and 100 results. With lower cut-off (10, 20, 50) the Recall is limited, since it is impossible to find all 88 non-compliance messages when only 50 messages are processed. With this experiment, the best results are obtained with the *lexical lemmas* corpus. The lemmatization shows an important improvement over the *forms* corpus, which means that lemmatization is important for the information retrieval applications. As expected, the values of Recall and Precision are improved with a larger sample of data: there is more probability that the 89 relevant messages are found among the top 100 messages. Overall, we can see that this information retrieval system can find non-compliant messages although the results are noisy and incomplete.

3.2 Supervised Categorization with Machine Learning

The results of the categorization of messages into the *non-compliance* category obtained with supervised machine learning algorithms are presented in Table 2. We tested several algorithms but show only the results for the best two algorithms, Naive Bayes and Naive Bayes Multinomial. We can observe that the best results (up to 0.824 F-measure) are obtained on the *lexical lemmas* corpus. In all the experiments, Recall is higher than or equal to Precision.

Among the errors observed with NaiveBayes, 12 messages are wrongly categorized as non-compliant and 9 as compliant. Within these 12 messages, four contain terms associated with excess and negative effects (such as "*Je n'imaginais pas que c'tait si grave*" (*I didn't imagine it was that bad*) or "*s'il vous plait ne faites pas n'importe quoi*" (*please don't make a mess*)), usually specific to non-compliance messages.

Table 2. Machine learning results obtained for the categorization of messages into the *non-compliance* category

	<i>Precision Recall F-measure</i>		
	<i>NaiveBayes</i>		
<i>forms</i>	0.769	0.769	0.769
<i>lemmas</i>	0.786	0.846	0.815
<i>lexical lemmas</i>	0.761	0.897	0.824
	<i>NaiveBayesMultinomial</i>		
<i>forms</i>	0.732	0.769	0.750
<i>lemmas</i>	0.795	0.795	0.795
<i>lexical lemmas</i>	0.786	0.846	0.815

We conclude that, although information retrieval systems can be adapted to perform categorization of messages thanks to the existing reference data, their results are less competitive comparing to the results obtained with machine learning algorithms. We assume that with larger reference data, information retrieval systems may be more competitive.

3.3 Unsupervised Detection with Information Retrieval

Here, we report a more standard exploitation of the information retrieval system. The experiments are done at a finer-grained level and focus on precise types of non-compliance. Several queries are tested. We will present queries and their results related to important drug non-compliance situations : gain and loose of weight, recreational drug use, suicide attempts or ideation, overdoses, and alcohol consumption. The descriptors used for the creation of queries are issued from the reference data. Their selection is based on their frequency and TFIDF scores [21]. The top 20 results are analyzed for each query.

Gaining/losing weight. The keywords used are *poids, kilo, grossir, maigrir (weight, kilo, gain weight, lose weight)*, such as suggested by the manually built reference data. This query is applied to the lemmatized corpus. We expected to find mainly messages related to the use of drugs with the purpose to intentionally lose or gain weight, as well as messages related to weight changes due to side effects of drugs. These expectations are partly verified. Thus, among the top 20 messages, 17 are about weight change as side effects of drugs, one message is about the use of drugs to lose weight intentionally, and 2 messages are about weight loss but with no relation to the consumption of drugs. Overall, among the top 20 messages returned by this query, one new non-compliance message is found. This situation may correspond to the reality (misuse of drugs for weight changes is less frequent than weight change due to drug side effects) or to the corpus used (several messages are concerned with anti-depressant drugs that have as a common side effect weight change). Overall, this gives 0.05 Precision.

Recreational drug use. The goal of this set of queries is to retrieve messages where prescription drugs are used for recreational purpose: be "high", reach hallucinations, feel happy, etc. We tried several queries:

- First, the descriptors *drogue*, *droguer* (*non-medical drug*, *to take non-medical drugs*) are used. In French, the word *drogue* usually refers to street drugs, but not to prescription drugs. Yet, in the corpus, people use this word for neuroleptic medication in order to illustrate their feeling that these drugs open the way to addictions and have the same neuroleptic effects as street drugs. Hence, we can find messages such as *J'ai t drogu pendant 3 ans au xanax* (*I was drugged with xanax for 3 years*) or *Sa soulage mais ses une vrai drogue ce truc !!!* (*It helps but this stuff is really a drug!!!*) These queries find interesting results (15 out of 20 messages), but may provide different insights than those expected;
- Then, the descriptors *hallu*, *allu*, *hallucination* (*hallucination*) are used. Among the top 20 messages, 2 messages report intentional seeking of the hallucination effects caused by some drugs, 7 messages are about people experiencing hallucinations but as unwanted side effects, 11 messages about people suffering from hallucinations and taking drugs to reduce them. This means that this query provided 2 new messages with non-compliance.
- Finally, the descriptor *planer* (*to be high from drugs*) is used. Among the top 20 messages, 19 are about the "high" feeling from drugs, be it intentional (9 messages) or non-intentional (10 messages). The 9 intentional messages correspond to non-compliance situations. Like in this message, *J'ai dj post quelques sujets propos de ce flau qu'est le stilnox (...) je prends du stilnox, pour m'vader, pour planer*" (*I already posted a few topics about this plague that is stilnox (...) I take stilnox, to escape, to get high*).

This set of queries illustrate how it is possible to apply an iterative process for the construction of appropriate queries quickly and with little reference data. Overall, the set of queries shows 0.35 Precision.

Suicide. Here, the descriptor used is *suicide* (*suicide*). The query is applied to the lemmatized corpus. We expected to find messages in which people report on taking or planning to take drugs (like anti-depressants) with suicidal intentions. Among the top 20 results, 9 messages are about drugs and suicide with no particular relation between them, 5 other messages are about the fact that some drugs may increase the risk of suicide, 5 other messages are critical about the fact that drugs may increase the risk of suicide, and one message reported on a real suicide attempt caused by drug withdrawal. We consider that discussions on relation between drugs and suicide, and of course reporting on suicide attempts, may be important for our research because they represent the importance of these topics in the analyzed fora. This gives 0.7 Precision.

Overuse. The descriptor used is *boites* (*boxes*) because it often represents the quantity of drugs taken in the reference data. This query is applied to the *forms*

corpus because, for this query, it is important to preserve the plural form. Among the top 20 messages, 6 messages are directly related to drug overuse, 3 messages are related to high dosage that may correspond to overuses, 2 messages are related to suicide attempts by ingestion of large amounts of drugs, 2 messages in which people propose to share unused prescription drugs, and, finally, 7 messages are unrelated to drug overuse. This gives 0.65 Precision. Another advantage of this query is that it retrieved various non-compliance situations.

Alcohol. The descriptor used is simply *alcohol* (*alcohol*). The query is applied to the lemmatized corpus. We expected to find discussions about adverse side effects of alcohol consumption while taking medication. These expectations are partly verified. Hence, among the top 20 results, 12 messages reported about interactions between alcohol and medication, 3 messages discussed about medication prescribed for alcohol withdrawal, while 5 messages were unrelated to the use of medication. Overall, this query gives 0.6 Precision. Additionally, we can notice that 8 out of the 12 messages regarding alcohol-drug interactions are dedicated to the mood disorder medications. This situation may reflect the fact that neuroleptic drugs interact with the effects of alcohol. But this situation may also be the artifact of the corpus containing several messages dedicated to mood disorder drugs.

3.4 Comparison of the Two Categorization Approaches

For the classification task, supervised machine learning shows better results. Yet it requires a too large amount of reference data to be usable to find specific types of non-adherence. For this task we found that information retrieval is able to find non-compliance messages from topics associated with a specific type of non-compliance, such as suicide or overuse. With this method we found 0.45 of average Precision with up to 0.70 for some queries. Besides in the examples of the *boxes* query we were able to discover a new type of non-compliance : people sharing their unused medication with others.

We conclude that these two approaches are complementary: combination of their results may provide an efficient way to enrich the data. The approaches can also be combined: information retrieval queries can quickly provide varied non-compliance messages and thus help supervised machine learning to perform better.

3.5 Limitations of the Current Work

The main limitation regarding our work is the reduced size of the reference data. It contains indeed only 133 messages in the *non-compliance* category. This may limit the performances of the supervised models. Yet, these reference data allow to create quite efficient categorization models, which reach up to 0.824 F-measure in the case of machine learning. We assume that availability of larger reference data will improve the performance of the methods. One of the main motivations

to exploit information retrieval methods is the possibility to enrich the reference data with this unsupervised approach.

Another limitation of the work is that messages detected as cases of non-compliance are not currently fully analyzed by medical doctors, pharmacists and pharmacovigilants. Further work will be needed to make the results exploitable by the medical community, who may not be familiar with the methods used in our experiments.

4 Conclusions

This work presents the exploitation of information retrieval methods in two different ways to detect drug non-compliance in Internet fora. We mainly exploit the French forum *Doctissimo*. The messages are first manually assigned into *compliance* and *non-compliance* categories and then used for designing automatic methods and their evaluation.

We adapt Indri, an information retrieval system, for supervised categorisation and evaluation, and obtain up to 0.60 Precision at top 10 results and up to 0.34 Precision at top 50 results. This method can be used to detect non-compliance messages but the noise prevents it from being competitive by comparison with machine learning approaches. Indeed, machine learning algorithms reach up to 0.786 Precision and 0.824 F-measure.

The information retrieval system is also used for a more fine-grained categorization of messages at the level of individual types of non-compliance, where the small number of messages in each category makes supervised learning impossible. This approach is fully unsupervised. Five topics are addressed with different queries suggested by the very few messages available in the reference data. This approach provides on average 0.42 Precision, and up to 0.70 Precision for some queries, such as computed among the top 20 results. This second approach exploiting the information retrieval system can also help in discovering new kinds of non-compliance situations and provide additional insight on topics of concern among patients.

The information gathered thank to these methods can be used by concerned experts (pharmaceutical industry, public health, general practitioners...) to provide prevention and education actions to patients and their relatives. For instance, packaging of drugs can be further adapted to their real use, dedicated brochures and discussions can be done with patients on known and possible drug side effects, on necessary precautions, etc.

The main perspective of the current work is to enrich the reference data and to work more closely with health professionals for the exploitation of the results.

Acknowledgments

This work has been performed as part of the DRUGSSAFE project funded by the ANSM, France and of the MIAM project funded by the ANR, France within

the reference ANR-16-CE23-0012. We thank both programs for their funding. We would like also to thank the annotators who helped us with the manual annotation of misuses, Bruno Thiao Layel for extracting the corpus, Vianney Jouhet and Bruno Thiao Layel for building the list with drugs names, and The-Hien Dao for the set of disorders exploited. Finally, we thank the whole ERIAS team for the discussions.

References

1. Bigeard, E., Grabar, N., Thiessard, F.: Typology of drug misuse created from information available in health fora. In: MIE 2018. (2018) 1–5
2. Gauducheau, N.: La communication des émotions dans les échanges médiatisés par ordinateur : bilan et perspectives. *Bulletin de psychologie* (2008) 389–404
3. Louis, A.: Natural language processing for social media. *Computational Linguistics* **42** (2016) 833–836
4. Collier, N.: Towards cross-lingual alerting for bursty epidemic events. *J Biomed Semantics* **2** (2011)
5. Lejeune, G., Brixstel, R., Lecluze, C., Doucet, A., Lucas, N.: Added-value of automatic multilingual text analysis for epidemic surveillance. In: *Artificial Intelligence in Medicine (AIME)*. (2013)
6. Tapi Nzali, M.: Analyse des mdias sociaux de sant pour valuer la qualit de vie des patientes atteintes dun cancer du sein. Thèse de doctorat, Universit de Montpellier, Montpellier, France (2017)
7. Morlane-Hondre, F., Grouin, C., Zweigenbaum, P.: Identification of drug-related medical conditions in social media. In: *LREC*. (2016) 1–7
8. Kalyanam, J., Katsuki, T., Lanckriet, G.R.G., Mackey, T.K.: Exploring trends of nonmedical use of prescription drugs and polydrug abuse in the twittersphere using unsupervised machine learning. *Addictive Behaviors* **65** (2017) 289–295
9. Cameron, D., Smith, G.A., Daniulaityte, R., Sheth, A.P., Dave, D., Chen, L., Anand, G., Carlson, R., Watkins, K.Z., Falck, R.: PREDOSE: a semantic web platform for drug abuse epidemiology using social media. **46** (2013) 985–997
10. Skrbo, A., Begović, B., Skrbo, S.: Classification of drugs using the atc system (anatomic, therapeutic, chemical classification) and the latest changes. *Med Arh* **58** (2004) 138–41
11. Cohen, J.: A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* **20** (1960) 37–46
12. Landis, J., Koch, G.: The measurement of observer agreement for categorical data. *Biometrics* **33** (1977) 159–174
13. Schmid, H.: Probabilistic part-of-speech tagging using decision trees. In: *ICNMLP*, Manchester, UK (1994) 44–49
14. Strohman, T., Metzler, D., Turtle, H., Croft, W.B.: Indri: a language-model based search engine for complex queries. In: *Proceedings of the International Conference on Intelligent Analysis*. (2005)
15. Witten, I., Frank, E.: *Data mining: Practical machine learning tools and techniques*. Morgan Kaufmann, San Francisco (2005)
16. John, G.H., Langley, P.: Estimating continuous distributions in bayesian classifiers. In Kaufmann, M., ed.: *Eleventh Conference on Uncertainty in Artificial Intelligence*, San Mateo (1995) 338–345

17. McCallum, A., Nigam, K.: A comparison of event models for naive bayes text classification. In: AAAI workshop on Learning for Text Categorization, Madison, Wisconsin (1998)
18. Quinlan, J.: C4.5 Programs for Machine Learning. Morgan Kaufmann, San Mateo, CA (1993)
19. Breiman, L.: Random forests. *Machine Learning* **45** (2001) 5–32
20. Landwehr, N., Hall, M., Frank, E.: Logistic model trees. *Machine Learning* **95** (2005) 161–205
21. Salton, G., McGill, M.J. In: Retrieval refinements. (1983) 199–206