

This is the peer reviewed version of the following article:

CLAUDETTE: an automated detector of potentially unfair clauses in online terms of service / Lippi, Marco; Pałka, Przemysław; Contissa, Giuseppe; Lagioia, Francesca; Micklitz, Hans-Wolfgang; Sartor, Giovanni; Torroni, Paolo. - In: ARTIFICIAL INTELLIGENCE AND LAW. - ISSN 0924-8463. - 27:2(2019), pp. 117-139. [10.1007/s10506-019-09243-2]

Terms of use:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

02/05/2026 07:27

(Article begins on next page)

CLAUDETTE: an Automated Detector of Potentially Unfair Clauses in Online Terms of Service

Marco Lippi · Przemyslaw Palka ·
Giuseppe Contissa · Francesca Lagioia ·
Hans-Wolfgang Micklitz · Giovanni
Sartor · Paolo Torroni

Received: date / Accepted: date

Abstract Terms of service of on-line platforms too often contain clauses that are potentially unfair to the consumer. We present an experimental study where machine learning is employed to automatically detect such potentially unfair clauses. Results show that the proposed system could provide a valuable tool for lawyers and consumers alike.

1 Introduction

A recent survey on policy-reading behaviour (Obar and Oeldorf-Hirsch 2016) reveals that consumers rarely read the contracts they are required to accept. This resonates with our direct experience and with what has long been said, that the biggest lie on the Internet is “I have read and agree to the terms and conditions.” We use smartphones to gather and share information, connect on social media, entertain ourselves, check our online banking and so on. Virtually

M. Lippi
DISMI – Università di Modena e Reggio Emilia, Italy
Tel.: +390522522660
E-mail: marco.lippi@unimore.it

P. Palka, F. Lagioia, H.-W. Micklitz
Law Department, European University Institute, Florence, Italy
E-mail: {przemyslaw.palka, francesca.lagioia, hans.micklitz}@eui.eu

G. Contissa, G. Sartor
CIRSFID, Alma Mater – Università di Bologna, Italy
E-mail: {giuseppe.contissa, giovanni.sartor}@unibo.it

P. Torroni
DISI, Alma Mater – Università di Bologna, Italy
E-mail: paolo.torroni@unibo.it

This is a pre-print of an article published in Artificial Intelligence and Law. The final authenticated version is available online at: <https://doi.org/10.1007/s10506-019-09243-2>.

every app we install and website we browse have their own Terms of Service (ToS). These are contracts that bind us by the time we switch on the phone or browse a website, although we are not necessarily aware of what we just agreed upon.

There are reasons why many consumers do not read or understand ToS, as well as privacy policies or end-user license agreements (EULA) (Bakos et al 2014). Reports indicate that such documents can be overwhelming to the few consumers who actually venture to read them (Department of Commerce 2010). It has been estimated that actually reading the privacy policies alone would carry costs in time of over 200 hours a year per Internet user (McDonald and Cranor 2008). Another problem is that even if consumers did read the ToS thoroughly, they would have no means to influence their content: the choice is to either agree to the terms offered by a web app or simply not use the service at all.

All this created a need for limitations on traders' contractual freedom, not only to protect consumer interests, but also to enhance the consumers' trust in transnational transactions and improve the common market (Nebbia 2007). European consumer law aims to prevent businesses from using so-called 'unfair contractual terms' in the contracts they unilaterally draft and require consumers to accept (Reich et al 2014). Law regarding such terms applies also to the ToS of on-line platforms (Loos and Luzak 2016). Unfortunately, it turns out such platforms' owners do use in their ToS unfair contractual clauses (Micklitz et al 2017), in spite of the European law, and regardless of consumer protection organizations agencies, which have the competence, but not necessarily the resources, to fight against such unlawful practices.

To address this problem, we propose a machine learning-based method and tool for partially automating the detection of potentially unfair clauses. Such a tool could be useful both for consumer protection organizations and agencies, **GS: by** making their work more effective and efficient, and for consumers, by increasing their understanding of what they agree upon.

This paper builds upon and significantly extends results presented by Lippi et al (2017) in a smaller-scale study where a Support Vector Machine (SVM) was trained on a 20-document corpus. With respect to previous work, the contributions of this study are:

- a larger corpus consisting of 50 contracts (over 12,000 sentences), so as to train and evaluate the proposed approach on a wider and more heterogeneous data set;
- an extensive comparison of several machine learning systems, including some recent deep learning architectures for text categorization, and a structured SVM for collective classification, which takes into account the sequence of sentences within a document;
- a more comprehensive task, which is not restricted to the detection of potentially unfair clauses but also encompasses the classification of such clauses into categories;

- the description of a web server, named CLAUDETTE, which we have made available to the community, so as to allow users to submit query documents and analyze the behavior of the system.

The paper is organized as follows. In Section 2 we first define the problem from the legal point of view. Then, in Section 3 we describe the novel corpus and the document annotation procedure. Section 4 explains the machine learning methodology employed in the system, whereas Section 5 presents experimental results. Section 6 describes the web server. Section 7 discusses related work. Section 8 concludes with a look to future research.

2 Problem Description

In this section we briefly introduce the European consumer law on unfair contractual terms (clauses). We explain what an unfair contractual term is, present the legal mechanisms created to prevent business from employing **GS: unfair terms**, and describe our contribution to these mechanisms.

According to art. 3 of the Directive 93/13 on Unfair Terms in Consumer Contracts, a contractual term is unfair if: 1) it has not been individually negotiated; and 2) contrary to the requirement of good faith, it causes a significant imbalance in the parties rights and obligations, to the detriment of the consumer. This general definition is **GS: further** specified in the Annex to the Directive, containing an indicative and non-exhaustive list of the terms which may be regarded as unfair, as well in a few dozen judgments of the Court of Justice of the EU (Micklitz and Reich 2014). Examples of unfair clauses encompass taking jurisdiction away from the consumer, limiting liability for damages on health and/or gross negligence, imposing obligatory arbitration in a country different from consumers residence etc.

Loos and Luzak (Loos and Luzak 2016) identified five categories of *potentially unfair* clauses often appearing in the terms of on-line services: 1) establishing jurisdiction for disputes in a country different than consumers residence; 2) choice of a foreign law governing the contract; 3) limitation of liability; 4) the provider’s right to unilaterally terminate the contract/access to the service; and 5) the provider’s right to unilaterally modify the contract/the service. Our research has identified three additional categories: 6) requiring a consumer to undertake arbitration before the court proceedings can commence; 7) the provider retaining the right to unilaterally remove consumer content from the service, including in-app purchases; 8) having a consumer accept the agreement simply by using the service, not only without reading it, but even without having to click on “I agree/I accept.”

The 93/13 Directive creates two mechanisms to prevent the use of unfair contractual terms: *individual* and *abstract* control of fairness. The former takes place when a consumer goes to court: if a court finds that a clause is unfair (which it can do on its own motion), it will consider that the clause is not binding on the consumer (art. 6). However, most consumers do not take their disputes to courts. That is why abstract fairness control has been

created. In each EU Member State, consumer protection organizations have the competence to initiate judicial or administrative proceedings, to obtain the declaration that clauses in consumer contracts are unfair. The national implementations of abstract control differ in various ways. For instance, consumer protection agencies and/or consumer organisations may be involved to a different degree, there may or may not be fines for using unfair contractual terms, etc. (Schulte-Nölke et al 2008). One thing that all member states have in common is that if a business uses unfair terms in their contracts, in principle there is always a competent party with the authority to challenge such contracts.

Unfortunately, the legal mechanism for enforcing the prohibition of unfair contract terms have failed to effectively counter this practice so far. As reported by some literature (Loos and Luzak 2016), and as our own research indicates (Micklitz et al 2017), unfair contractual terms are, as of today, widely used in ToS of online platforms.

In our previous research (Micklitz et al 2017), we developed a theoretical model of tasks that human lawyers currently need to carry out, before starting the legal proceedings concerning the abstract control of fairness of clauses. Those include: 1) finding and choosing the documents; 2) mining the documents for potentially unfair clauses; 3) conducting the actual legal assessment of fairness; 4) drafting the case files and beginning the proceedings. Our work aims to automate the second step, enabling a senior lawyer to focus only on clauses that are found by a machine learning classifier to be potentially unfair, thus saving significant time and labor.

We focus on *potentially* unfair clauses for two reasons. First, we may be unsure whether a certain type of clause falls under the abstract legislative definition of an “unfair contractual term”. From a legal standpoint, a given clause can be deemed unfair with absolute certainty only if a competent institution, such as a national court having referred to the European Court of Justice, has ruled in that sense. That is the case for certain kinds of clauses, such as a jurisdiction clause indicating a country different from the consumer’s residence, or limitation of liability for gross negligence (Micklitz et al 2017). In other cases the unfairness of a clause has to be argued for, showing that it creates an unacceptable imbalance in the parties’ rights and obligations. A consumer protection body might want to take the case to a court in order to authoritatively establish the unfairness of that clause, but a legal argument for that needs to be created, and the clause may eventually turn out to be judged fair. As a second point, unfairness may depend not only on a clause’s textual content, but also on the context in which the clause is to be applied. For instance, a mutual right to unilaterally terminate the contract might be fair in some cases, and unfair in others, for example if unilateral termination would entail losing some digital content (purchased apps, email address, etc.) on the side of the consumer.

3 Corpus Annotation

The corpus consists of 50 relevant on-line consumer contracts, i.e. the Terms of Service of on-line platforms. Such contracts were selected among those offered by some of the major players in terms of number of users, global relevance, and time of establishment of the service.¹ Such contracts are usually quite detailed in content, are frequently updated to reflect changes both in the service and in the applicable law, and are often available in different versions for different jurisdictions. In the presence of multiple versions of the same contract, we selected the most recent version available on-line for the European customers. The mark-up was done in XML. A description of the annotation process follows.

3.1 Annotation process

In analyzing the Terms of Service of the selected on-line platforms, we identified eight different categories of unfair clauses. For each type of clause we defined a corresponding XML tag, as shown in Table 1.

Notice that not necessarily all the documents contain all clause categories. For example, Twitter provides two different ToS, the first one for US and non-US residents and the second one for EU residents. The tagged version is the version applicable in the EU and it does not contain any choice of law, arbitration or jurisdiction clauses.

We assumed that each type of clause could be classified as clearly fair, potentially unfair, or clearly unfair. In order to mark the different degrees of (un)fairness we appended a numeric value to each XML tag, with 1 meaning clearly fair, 2 potentially unfair, and 3 clearly unfair. Nested tags were used to annotate text segments relevant to more than one type of clause. If one clause covers more than one paragraph, we chose to tag each paragraph separately, possibly with different degrees of (un)fairness.

The **jurisdiction** clause stipulates what courts will have the competence to adjudicate disputes under the contract. Jurisdiction clauses giving consumers a right to bring disputes in their place of residence were marked as clearly fair, whereas clauses stating that any judicial proceeding takes a residence away (i.e. in a different city, different country) were marked as clearly unfair. This assessment is grounded in ECJ's case law, see for example *Oceano* case no. C-240/98. An example of jurisdiction clauses is the following one, taken from the Dropbox terms of service:

¹ In particular, we selected the ToS offered by: 9gag.com, Academia.edu, Airbnb, Amazon, Atlas Solutions, Betterpoints, Booking.com, Crowdtangle, Deliveroo, Dropbox, Duolingo, eBay, Endomondo, Evernote, Facebook, Fitbit, Google, Headspace, Instagram, Linden Lab, LinkedIn, Masquerade, Microsoft, Moves-app, musically, Netflix, Nintendo, Oculus, Onavo, Pokemon GO, Rovio, Skype, Skyscanner, Snapchat, Spotify, Supercell, SyncMe, Tinder, TripAdvisor, TrueCaller, Twitter, Uber, Viber, Vimeo, Vivino, WhatsApp, World of Warcraft, Yahoo, YouTube and Zynga. The annotated corpus can be downloaded from <http://155.185.228.137/claurette/ToS.zip>.

Table 1 Categories of clause unfairness, with the corresponding symbol used for tagging.

Type of clause	Symbol
Arbitration	<a>
Unilateral change	<ch>
Content removal	<cr>
Jurisdiction	<j>
Choice of law	<law>
Limitation of liability	<ltld>
Unilateral termination	<ter>
Contract by using	<use>

```
<j3> You and Dropbox agree that any judicial proceeding to resolve
claims relating to these Terms or the Services will be brought in the
federal or state courts of San Francisco County, California, subject to
the mandatory arbitration provisions below. Both you and Dropbox consent
to venue and personal jurisdiction in such courts.</j3>
```

```
<j1>If you reside in a country (for example, European Union
member states) with laws that give consumers the right to bring
disputes in their local courts, this paragraph doesn't affect those
requirements.</j1>
```

The second clause introduces an exception to the general rule stated in the first clause, thus we marked the first one as clearly unfair and the second as clearly fair.

The **choice of law** clause specifies what law will govern the contract, meaning also what law will be applied in potential adjudication of a dispute arising under the contract. Clauses defining the applicable law as the law of the consumer's country of residence were marked as clearly fair, as reported in the following examples, taken from the Microsoft services agreements:

```
<law1>If you live in (or, if a business, your principal place of
business is in) the United States, the laws of the state where you live
govern all claims, regardless of conflict of laws principles, except
that the Federal Arbitration Act governs all provisions relating to
arbitration.</law1>
```

```
<law1>If you acquired the application in the United States or Canada,
the laws of the state or province where you live (or, if a business,
where your principal place of business is located) govern the
interpretation of these terms, claims for breach of them, and all other
claims (including consumer protection, unfair competition, and tort
claims), regardless of conflict of laws principles.</law1>
```

```
<law1>Outside the United States and Canada. If you acquired the
application in any other country, the laws of that country apply.</law1>
```

In every other case, the choice of law clause was considered as potentially unfair. This is because the evaluation of the choice of law clause needs to take into account several other conditions besides those specified the clause itself (for example, level of protection offered by the chosen law). Consider the following example, taken from the Facebook terms of service:

<law2>The laws of the State of California will govern this Statement, as well as any claim that might arise between you and us, without regard to conflict of law provisions</law2>

The **limitation of liability** clause stipulates that the duty to pay damages is limited or excluded, for certain kind of losses, under certain conditions. Clauses that explicitly affirm non-excludable providers' liabilities were marked as clearly fair. For example, consider the example below, taken from World of Warcraft terms of use:

<1td1>Blizzard Entertainment is liable in accordance with statutory law (i) in case of intentional breach, (ii) in case of gross negligence, (iii) for damages arising as result of any injury to life, limb or health or (iv) under any applicable product liability act.</1td1>

Clauses that reduce, limit, or exclude the liability of the service provider were marked as potentially unfair when concerning broad categories of losses or causes of them, such as any harm to the computer system because of malware or loss of data or the suspension, modification, discontinuance or lack of the availability of the service. Also those liability limitation clauses containing a blanket phrase like "to the fullest extent permissible by law", where considered potentially unfair. The following example is taken from 9gag terms of service:

<1td2>You agree that neither 9GAG, Inc nor the Site will be liable in any event to you or any other party for any suspension, modification, discontinuance or lack of availability of the Site, the service, your Subscriber Content or other Content.</1td2>

Clause meant to reduce, limit, or exclude the liability of the service provider for physical injuries, intentional damages as well as in case of gross negligence were marked as clearly unfair (based on the Annex to the Directive) as showed in the example below, taken from the Rovio license agreement:

<1td3> In no event will Rovio, Rovio's affiliates, Rovio's licensors or channel partners be liable for special, incidental or consequential damages resulting from possession, access, use or malfunction of the Rovio services, including but not limited to, damages to property, loss of goodwill, computer failure or malfunction and, to the extent permitted by law, damages for personal injuries, property damage, lost profits or punitive damages from any causes of action arising out of or related to this EULA or the software, whether arising in tort (including negligence), contract, strict liability or otherwise and whether or not Rovio, Rovio's licensors or channel partners have been advised of the possibility of such damages.</1td3>

The **unilateral change** clause specifies the conditions under which the service provider could amend and modify the terms of service and/or the service itself. Such clause was always considered as potentially unfair. This is because the ECJ has not yet issued a judgment in this regard, though the Annex to the Directive contains several examples supporting such a qualification. Consider the following examples from the Twitter terms of service:

<ch2>As such, the Services may change from time to time, at our discretion.</ch2>

<ch2>We also retain the right to create limits on use and storage at our sole discretion at any time.</ch2>

<ch2>We may revise these Terms from time to time. The changes will not be retroactive, and the most current version of the Terms, which will always be at twitter.com/tos, will govern our relationship with you.</ch2>

The **unilateral termination** clause gives provider the right to suspend and/or terminate the service and/or the contract, and sometimes details the circumstances under which the provider claims to have a right to do so. Unilateral termination clauses that specify reasons for termination were marked as potentially unfair, whereas clauses stipulating that the service provider may suspend or terminate the service at any time for any or no reasons and/or without notice were marked as clearly unfair. That is the case in the three following examples, taken from the Dropbox and Academia terms of use, respectively:

<ter2> We reserve the right to suspend or terminate your access to the Services with notice to you if: (a) you're in breach of these Terms, (b) you're using the Services in a manner that would cause a real risk of harm or loss to us or other users, or (c) you don't have a Paid Account and haven't accessed our Services for 12 consecutive months.</ter2>

<ter3>Academia.edu reserves the right, at its sole discretion, to discontinue or terminate the Site and Services and to terminate these Terms, at any time and without prior notice.</ter3>

The **contract by using** clause stipulates that the consumer is bound by the terms of use of a specific service, simply by using the service, without even being required to mark that he or she has read and accepted them. We always marked such clauses as potentially unfair. The reason for this choice is that a good argument can be offered for these clauses to be unfair, because they originate an imbalance in rights and duties of the parties, but this argument has no decisive authoritative backing yet, since the ECJ has never assessed a clause of this type. Consider an example taken from the Spotify terms and conditions of use:

<use2>By signing up or otherwise using the Spotify service, websites, and software applications (together, the "Spotify Service" or "Service"), or accessing any content or material that is made available by Spotify through the Service (the "Content") you are entering into a binding contract with the Spotify entity indicated at the bottom of this document.</use2>

The **content removal** gives the provider a right to modify/delete user's content, including in-app purchases, and sometimes specifies the conditions under which the service provider may do so. As in the case of unilateral termination, clauses that indicate conditions for content removal were marked as potentially unfair, whereas clauses stipulating that the service provider may remove content in his full discretion, and/or at any time for any or no reasons and/or without notice nor possibility to retrieve the content were marked as clearly unfair. For instance, consider the following examples, taken from Facebook's and Spotify's terms of use:

<cr2> If you select a username or similar identifier for your account or Page, we reserve the right to remove or reclaim it if we believe it is appropriate (such as when a trademark owner complains about a username that does not closely relate to a user's actual name).</cr2>

<cr2> We can remove any content or information you post on Facebook if we believe that it violates this Statement or our policies.</cr2>

<cr3>In all cases, Spotify reserves the right to remove or disable access to any User Content for any or no reason, including but not limited to, User Content that, in Spotify’s sole discretion, violates the Agreements. Spotify may take these actions without prior notification to you or any third party.</cr3>

The **arbitration** clause requires or allows the parties to resolve their disputes through an arbitration process, before the case could go to court. It is therefore considered a kind of forum selection clause. However, such a clause may or may not specify that arbitration should occur within a specific jurisdiction. Clauses stipulating that the arbitration should (1) take place in a state other than the state of consumer’s residence and/or (2) be based not on law but on arbiter’s discretion were marked as clearly unfair. As an illustration, consider the following clause of the Rovio terms of use:

<j1> <a3>Any dispute, controversy or claim arising out of or relating to this EULA or the breach, termination or validity thereof shall be finally settled at Rovio’s discretion (i) at your domicile’s competent courts; or (ii) by arbitration in accordance with the Rules for Expedited Arbitration of the Arbitration Institute of the Finland Chamber of Commerce. The arbitration shall be conducted in Helsinki, Finland, in the English language.</a3> </j1>

Notice that the above clause concerns both jurisdiction and arbitration (notice the use of nested tags). Clauses defining arbitration as fully optional would have to be marked as clearly fair. However, our corpus does not contain any example of fully optional arbitration clause. Thus, all arbitration clauses were marked as potentially unfair. An example is the following segment of Amazon’s terms of service:

<a2>Any dispute or claim relating in any way to your use of any Amazon Service, or to any products or services sold or distributed by Amazon or through Amazon.com will be resolved by binding arbitration, rather than in court, except that you may assert claims in small claims court if your claims qualify. The Federal Arbitration Act and federal arbitration law apply to this agreement.</a2>

3.2 Corpus statistics

The final corpus contains 12,011 sentences² overall, 1,032 of which (8.6%) were labeled as positive, thus containing a potentially unfair clause. The distribution of the different categories across the 50 documents is reported in Table 2: arbitration clauses are the least common, being present in 28 documents only, whereas all the other categories appear in at least 40 out of 50 documents. Limitation of liability and unilateral termination categories represent more than half of the total potentially unfair clauses. The percentage of potentially unfair clauses in each document is quite heterogeneous, ranging from 3.3% (Microsoft) up to 16.2% (TrueCaller).

² The segmentation into sentences was obtained with Stanford CoreNLP suite.

Table 2 Corpus statistics. For each category of clause unfairness, we report the overall number of clauses and the number of documents they appear in.

Type of clause	# clauses	# documents
Arbitration	44	28
Unilateral change	188	49
Content removal	118	45
Jurisdiction	68	40
Choice of law	70	47
Limitation of liability	296	49
Unilateral termination	236	48
Contract by using	117	48

4 Machine Learning Methodology

In this section we briefly describe the representation and learning methods we used in our study. We address two different tasks: a detection task, aimed at predicting whether a given sentence contains a (potentially) unfair clause, and a classification task, aimed at predicting the category an unfair clause belongs to, which indeed could be a valuable piece of information to a potential user.

4.1 Learning algorithms

We address the problem of detecting potentially unfair contract clauses as a sentence classification task. Such a task could be tackled by treating sentences independently of one another (*sentence-wide* classification). This is the most standard and classic approach in machine learning, traditionally addressed by methods such as Support Vector Machines or Artificial Neural Networks (including recent deep learning approaches).

Alternatively, one could take into account the structure of the document, in particular the *sequence* of sentences, so as to perform a *collective* classification. Because it is not uncommon for unfair clauses to span across consecutive sentences in a document, this approach could also have some advantages.

In sentence-wide classification the problem can be formalized as follows. Given a sentence, the goal is to classify it as *positive* if it contains a potentially unfair clause, or *negative* otherwise. Within this setting, a machine learning classifier is trained with a data set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, which consists of a collection of N pairs, where x_i encodes some representation of a sentence, and y_i is its corresponding (positive or negative) class.

In collective classification, the data set consists of a collection of M *documents*, represented as sequences of sentences:

$$\mathcal{D} = \{d_j = \{(x_1^j, y_1^j), \dots, (x_{k_j}^j, y_{k_j}^j)\}\}_{j=1}^M,$$

where the j -th document contains k_j sentences.

Different machine learning systems can be developed for each classification setup, according to the learning framework and to the features employed

to represent each sentence. As for the learning methodology, for sentence-wide classification in this paper we compare Support Vector Machines (SVMs) (Joachims 1998) with some recent deep learning architectures, namely Convolutional Neural Networks (CNNs) (Kim 2014) and Long-Short Term Memory Networks (LSTMs) (Graves and Schmidhuber 2005). For collective classification, we rely on structured Support Vector Machines, and in particular on SVM-HMMs, which combine SVMs with Hidden Markov Models (Tsochantaridis et al 2005), by jointly assigning a label to each element in a given sequence (in our case, to each sentence in the considered document).

4.2 Sentence representation

As for the features represented to encode sentences, in an effort to make our method as general as possible, we decided to opt for traditional features for text categorization, excluding other, possibly more sophisticated, handcrafted features.

One of the most classic, yet still widely used, set of features for text categorization, is the well-known *bag-of-words* (BoW) model. In such a model, one feature is associated to each word in the vocabulary: the value of such feature is either zero, if the word does not appear in the sentence, or other than zero, if it does. Such a value is usually computed as the TF-IDF score, that is product of the number of occurrences of the word in the sentence (Term Frequency, TF) by a term that strengthen the contribution of infrequent words (Inverse Document Frequency, IDF) (Sebastiani 2002).

The BoW model can be extended to consider also n -grams, i.e., consecutive word combinations, rather than simple words, so that the order of the words in the sentences is (at least locally) exploited. Grammatical information can be included as well, by constructing a bag of part-of-speech tags, i.e., word categories such as nouns, verbs, etc. (Leopold and Kindermann 2002). Despite their simplicity, BoW features are very informative, as they encode the lexical information of a sentence, and thus represent a challenging baseline in those cases where the presence of some keywords and phrases is highly discriminative for the categorization of sentences.

A second approach we consider for the representation of a sentence is to exploit a constituency parse tree, which naturally encodes the *structure* of the sentence (see Figure 1) by describing the grammatical relations between sentence portions through a tree. Similarity between tree structures can be exploited with *tree kernels* (Moschitti 2006) (TK). A TK consists of a *similarity measure* between two trees, which takes into account the number of common substructures or *fragments*. Different definitions of fragments induce different TK functions. In our study we use the SubSet Tree Kernel (SSTK) (Collins and Duffy 2002) which counts as fragments those subtrees of the constituency parse tree terminating either at the leaves or at the level of non-terminal symbols. SSTK have been shown to outperform other TK functions in several argumentation mining sub-tasks (Lippi and Torroni 2016b).

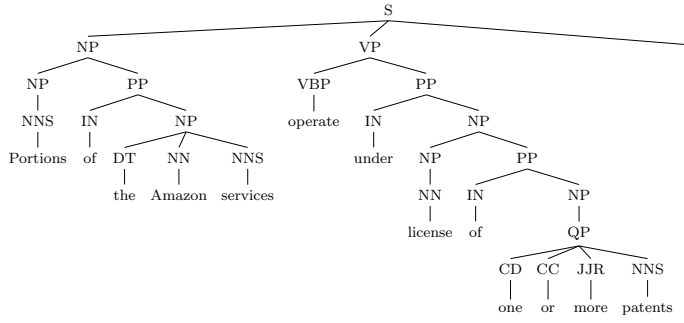


Fig. 1 An example of a constituency parse tree for a sentence in our corpus.

A third approach for sentence representation is based on *word embeddings* (Mikolov et al 2013), a popular technique that has been recently developed in the context of neural language models and deep learning applications. Neural networks such as CNNs and LSTMs can handle textual input, by converting it into a sequence of identifiers (one for each different word): it is the neural network, then, which directly learns a vector representation (named embedding) of words and sentences.

5 Experimental Results

We evaluated and compared several different machine learning systems on the data set presented in Section 3. Each document was segmented into sentences, tokenized and parsed with the Stanford CoreNLP tool.³ Sentences and text fragments shorter than 5 words were discarded. We obtained a total of 9,414 sentences, 1,032 of which (11.0%) were labeled as positive, thus containing a potentially unfair clause. We run experiments following the *leave-one-document-out* (LOO) procedure, in which each document in the corpus, in turn, is used as test set, leaving the remaining documents for training set (4/5) and validation set (1/5) for model selection. To quantitatively evaluate the different systems, we computed precision (P) as the fraction of positive predictions, which are actually labeled as positive, recall (R) as the fraction of positive examples that are correctly detected, and finally F_1 as the harmonic mean between precision and recall ($F_1 = \frac{2PR}{P+R}$). These performance measurements were aggregated using the macro-average over documents (Sebastiani 2002).

For the first task (potentially unfair clause detection) we compared several systems. The problem is formulated as a binary classification task, where the positive class is either the union of all potentially unfair sentences, or the set of potentially unfair clauses of a single category, as described below. We considered the following systems:

³ <https://stanfordnlp.github.io/CoreNLP/>

- C1: a single SVM exploiting BoW (unigrams and bigrams for words and part-of-speech tags);
- C2: a combination of eight SVMs (same features as above), each considering a single unfairness category as the positive class, whereby a sentence is predicted as potentially unfair if at least one of the SVMs predicts it as such;
- C3: a single SVM exploiting TK for sentence representation;
- C4: a CNN trained from plain word sequences;
- C5: an LSTM trained from plain word sequences;
- C6: an SVM-HMM performing collective classification of sentences in a document (same features as C1);
- C7: a combination of eight SVM-HMMs, each performing collective classification of sentences in a document on a single unfairness category as the positive class (same setting as C2);
- C8: an ensemble method, that combines the output of C1, C2, C3, C6 and C7 with a voting procedure (sentence predictive as positive if at least 3 systems out of 5 classify it as such).

As a reference for the complexity of the task, we also report the performance of the following baselines: a *random* classifier, which predicts potentially unfair clauses at random,⁴ and an *always positive* baseline, which classifies every sentence as potentially unfair. For all the classifiers, the validation set was used to select the best hyper-parameters. For all SVMs we used a linear kernel, thus optimizing the C parameter only. For SVM-HMM we used an order of dependencies equal to 2 and 1 for transitions and emissions, respectively; different from SVMs, we also used trigrams besides unigrams and bigrams, as they slightly increased performance. For CNNs, we considered one layer with 64 filters of size equal to 3, followed by two fully connected layers with 32 and 16 neurons, respectively. We applied dropout equal to 0.5, batch size equal to 16. An embedding of size 64 was learned after the input layer. For LSTMs, we considered a 2-layer network with 64 and 32 cells, respectively, with 0.25 dropout and mini-batch size equal to 16. An embedding of size 32 was learned after the input layer. Both for CNNs and LSTMs, no improvement was observed if using pre-trained word embeddings.

Table 3 shows the results achieved by each of these variants. If we exclude the ensemble approach, the best classifier in terms of F_1 results to be C2, that is the system combining one different SVM trained for each unfairness category, with a precision above 80%, and a recall of 78%. The structured SVMs exploiting the sequentiality of the sentences achieve slightly lower results, yet very interestingly the results of the sentence-wise and document-wise approaches are different across different documents. Moreover, the worse performance associated with TK suggests that the syntactic structure of the sentence is less informative than the lexical information captured by n -grams. This makes the task of detecting unfair clauses different from other text retrieval problems in the legal domain, such as, for example, the detection of claims and argu-

⁴ Sampling takes into account the class distribution in the training set.

Table 3 Results on leave-one-document-out procedure.

Classifier	Method	P	R	F_1
C1	SVM – Single Model	0.729	0.830	0.769
C2	SVM – Combined Model	0.806	0.779	0.784
C3	Tree Kernels	0.777	0.718	0.739
C4	Convolutional Neural Networks	0.729	0.739	0.722
C5	Long Short-Term Memory Networks	0.696	0.723	0.698
C6	SVM-HMM – Single Model	0.759	0.778	0.758
C7	SVM-HMM – Combined Model	0.848	0.720	0.772
C8	Ensemble (C1+C2+C3+C6+C7)	0.828	0.798	0.806
	Random Baseline	0.125	0.125	0.125
	Always Positive Baseline	0.123	1.000	0.217

ments (Lippi and Torroni 2016a). As for CNNs and LSTMs, the slightly worse performance with respect to the other approaches could also be ascribed to the limited size of the training set.

All these observations led us to the implementation of an ensemble method (C8), combining the five best performing approaches. This system achieves an F_1 of around 81%, thus beating all the competitors. Such a result is particularly interesting, because it confirms that the different systems capture complementary information for the detection of potentially unfair clauses. The ensemble method correctly detects over 75% of the potentially unfair clauses of all the categories, from 76.6% of Unilateral Change up to 89.7% for Jurisdiction.

In order to gain a better understanding of which are the n -grams that contribute the most to the discrimination between fair and potentially unfair clauses, we computed the frequencies of 2-grams in both positive and negative support vectors of classifier C2, and we looked for those with the largest discrepancy in appearing in the positive class rather than in the negative one. These were some of the most significant 2-grams, according to such ranking: *for any, the right, these terms, any time, at any, right to, reserves the, we may, liable for, terminate your, sole discretion, the services*. This analysis confirms that the discriminative lexicon is quite general and widespread both across the different unfairness categories and the different types of services we considered.

The second task we considered is unfairness categorization, for which we employed eight SVM classifiers, each trained to discriminate between potentially unfair clauses of one category with respect to all the other categories. Note that this task differs from that addressed by the previously introduced classifiers, since in this case the classifiers is trained on potentially unfair clauses only. In Table 4 we report the precision, recall, and F_1 of such classifiers, one for each separate tag category, micro-averaged on the whole dataset. The results show that discriminating amongst the different categories is a simpler task, since the F_1 is larger than 74% for all tags, and larger than 93% for four tags (jurisdiction, choice of law, limitation of liability, contract by using).

Table 4 Micro-averaged precision, recall and F_1 of abusive clauses for each tag category.

Tag	Precision	Recall	F_1
Arbitration	0.832	0.814	0.823
Unilateral change	0.832	0.814	0.823
Content removal	0.713	0.780	0.745
Jurisdiction	1.000	0.941	0.970
Choice of law	0.984	0.886	0.932
Limitation of liability	0.961	0.905	0.932
Unilateral termination	0.786	0.932	0.853
Contract by using	0.949	0.957	0.953

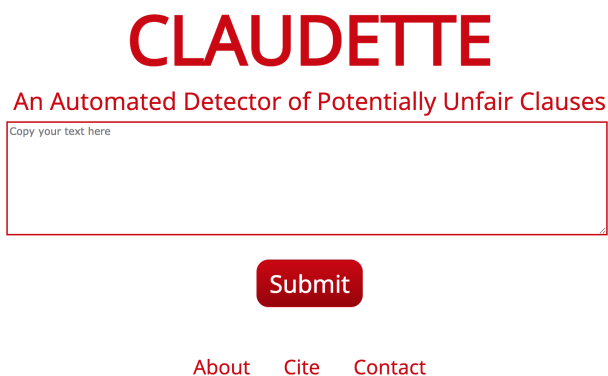


Fig. 2 The interface of the CLAUDETTE web server, consisting of a box where a user can copy-paste the text of a terms of service.

6 The CLAUDETTE Web Server

The proposed approach was implemented and developed within a web server, reachable at the address `http://155.185.228.137/claurette/`, so as to produce a prototype system that users can easily access and test.

As shown in Figure 2, the interface is easy to use. A user only needs to paste the text to be analyzed and push a button. The system will then produce an output file that highlights the sentences predicted to contain a potentially unfair clause. The output will also indicate the predicated category the unfair clause belongs to, as illustrated in Figure 3. The output of the system can be obtained in several formats including HTML, XML, JSON, and plain text.

For this online service, for the detection stage we implemented only one system (namely, classifier C2) rather than the ensemble method, because it resulted to be a much more efficient solution, despite producing a slightly lower performance accuracy.

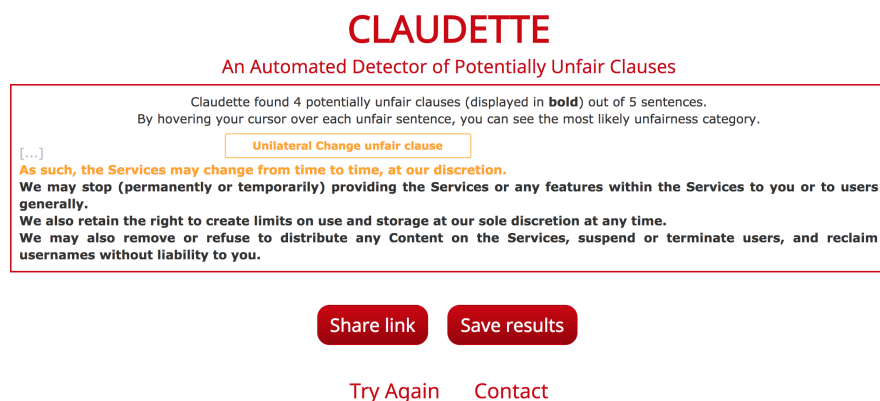


Fig. 3 Results of a query to the CLAUDETTE web server. Hovering over a detected clause with the pointer provides an indication of the type of potentially unfair clause. In this example the detected clauses are predicted to be of types *unilateral change*, *unilateral termination*, and *content removal*, and the cursor was left hovering over the first potentially unfair clause.

7 Related Work

The use of artificial intelligence, machine learning and natural language processing techniques in the analysis and classification of legal documents is gaining a growing interest (Ashley 2017). Among others, Moens et al (2007) proposed a pipeline of steps for the extraction of arguments from legal documents, exploiting supervised classifiers and context-free grammars, whereas Biagioli et al (2005) proposed to employ multi-class SVM for the identification of significant text portions in normative texts. Recent approaches have focused on the detection of claims (Lippi et al 2018) and of cited facts and principles in legal judgments (Shulayeva et al 2017), as well as on the prediction of judicial decisions (Aletras et al 2016). A case study regarding the construction of legal arguments in the legal determinations of vaccine/injury compensation compliance using natural language tools was given by Ashley and Walker (2013). Finally, privacy policies represent another strictly related application where machine learning approaches have proved effective, as discussed by Fabian et al (2017) and references therein, as well as Harkous et al (2018). Besides the work by Lippi et al (2017) which we discussed in the introduction, we are not aware of other recent text processing applications in consumer law.

8 Conclusions

Our study investigates the use of machine learning and natural language methods for the automated detection of potentially unfair clauses in online contracts. We addressed two tasks: clause detection and clause type classification. For clause detection, our results are very encouraging: using a relatively small

training set we could automatically detect over 80% clauses, with an 80% precision. The categorization task turned out to be simpler. Given that most unfair clauses are currently hidden within long and hardly readable ToS, the recall and precision offered by our approach may already be significant enough to enable useful applications.

It is interesting to notice the comparatively better performance of the BoW approach with respect to other more sophisticated approaches. That is in agreement with the surveyed literature, where classic lexical approaches such as BoW still represent a crucial ingredient of automated systems. It is also worth remarking that an ensemble method produced the best performance, thus indicating that different machine learning approaches are capable of capturing different characteristics of potentially unfair clauses.

This study was motivated by a long-term goal such as the pursuit of effective consumer protection by way of tools that support consumers and their organizations in detecting unfair contractual clauses. We plan to extend our analysis to other machine learning methods that could contribute to such tools. In particular, we are studying ways to exploit contextual information, since it was pointed out that the fairness of clauses might very well depend on the context. For example, a potentially unfair jurisdiction clause might actually be fair according to EU regulation if is followed by a paragraph stipulating relevant exceptions according to the user's country of residence.

As a further, challenging line of research, we are planning to apply similar methodologies also to privacy policies: an important area of consumer protection that has recently gained media focus due to its enormous implications not only for individuals but also for society at large.

References

- Aletras N, Tsarapatsanis D, Preoiuc-Pietro D, Lamos V (2016) Predicting judicial decisions of the European Court of Human Rights: a natural language processing perspective. *PeerJ Computer Science* 2:e93
- Ashley K (2017) *Artificial Intelligence and Legal Analytics: New Tools for Law Practice in the Digital Age*. Cambridge University Press
- Ashley KD, Walker VR (2013) Toward constructing evidence-based legal arguments using legal decision documents and machine learning. In: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Law*, ACM, pp 176–180
- Bakos Y, Marotta-Wurgler F, Trossen DR (2014) Does anyone read the fine print? consumer attention to standard-form contracts. *The Journal of Legal Studies* 43(1):1–35
- Biagioli C, Francesconi E, Passerini A, Montemagni S, Soria C (2005) Automatic semantics extraction in law documents. In: *Proceedings of ICAIL*, ACM, pp 133–140
- Collins M, Duffy N (2002) New ranking algorithms for parsing and tagging: Kernels over discrete structures, and the voted perceptron. In: *Proceedings of the 40th Annual Meeting of the ACL*, ACL, pp 263–270
- Department of Commerce (2010) *Commercial data privacy and innovation in the Internet economy: A dynamic policy framework*. Tech. rep., Department of Commerce Internet Policy Task Force, URL https://www.ntia.doc.gov/files/ntia/publications/iptf_privacy_greenpaper_12162010.pdf
- Fabian B, Ermakova T, Lentz T (2017) Large-scale readability analysis of privacy policies. In: *Proceedings of the International Conference on Web Intelligence*, ACM, pp 18–25

- Graves A, Schmidhuber J (2005) Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks* 18(5):602–610
- Harkous H, Fawaz K, Lebrete R, Schaub F, Shin KG, Aberer K (2018) Polisis: Automated analysis and presentation of privacy policies using deep learning. arXiv preprint arXiv:180202561
- Joachims T (1998) Text categorization with support vector machines: Learning with many relevant features. *ECML 98* pp 137–142
- Kim Y (2014) Convolutional neural networks for sentence classification. In: Moschitti A, Pang B, Daelemans W (eds) *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, Doha, Qatar, ACL*, pp 1746–1751
- Leopold E, Kindermann J (2002) Text categorization with support vector machines. How to represent texts in input space? *Machine Learning* 46(1-3):423–444
- Lippi M, Torroni P (2016a) Argumentation mining: State of the art and emerging trends. *ACM Trans Internet Technol* 16(2):10:1–10:25
- Lippi M, Torroni P (2016b) Margot: A web server for argumentation mining. *Expert Systems with Applications* 65(C):292–303, DOI 10.1016/j.eswa.2016.08.050
- Lippi M, Palka P, Contissa G, Lagioia F, Micklitz H, Panagis Y, Sartor G, Torroni P (2017) Automated detection of unfair clauses in online consumer contracts. In: Wyner AZ, Casini G (eds) *Legal Knowledge and Information Systems - JURIX 2017: The Thirtieth Annual Conference, Luxembourg, 13-15 December 2017*, IOS Press, *Frontiers in Artificial Intelligence and Applications*, vol 302, pp 145–154
- Lippi M, Lagioia F, Contissa G, Sartor G, Torroni P (2018) Claim detection in judgments of the EU Court of Justice. In: *Artificial Intelligence and the Complexity of Legal Systems, VI International Workshop (AICOL)*, selected revised papers. *Lecture Notes in Artificial Intelligence*, Springer, forthcoming
- Loos M, Luzak J (2016) Wanted: a bigger stick. On unfair terms in consumer contracts with online service providers. *Journal of Consumer Policy* 39(1):63–90
- McDonald A, Cranor L (2008) The cost of reading privacy policies. *I/S: A Journal of Law and Policy for the Information Society* 4(3), URL http://moritzlaw.osu.edu/students/groups/is/files/2012/02/Cranor_Formatted_Final.pdf, issue: *2008 Privacy Year in Review*
- Micklitz HW, Reich N (2014) The court and sleeping beauty: The revival of the unfair contract terms directive (uctd). *Common Market Law Review* 51(3):771–808
- Micklitz HW, Palka P, Panagis Y (2017) The empire strikes back: Digital control of unfair terms of online services. *Journal of Consumer Policy* pp 1–22
- Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient estimation of word representations in vector space. arXiv preprint arXiv:13013781
- Moens MF, Boiy E, Palau RM, Reed C (2007) Automatic detection of arguments in legal texts. In: *Proceedings of the 11th international conference on Artificial intelligence and law*, ACM, pp 225–230
- Moschitti A (2006) Efficient convolution kernels for dependency and constituent syntactic trees. In: Frnkranz J, Scheffer T, Spiliopoulou M (eds) *ECML 2006*, Springer Berlin Heidelberg
- Nebbia P (2007) *Unfair contract terms in European law: A study in comparative and EC law*. Bloomsbury Publishing
- Obar JA, Oeldorf-Hirsch A (2016) The biggest lie on the internet: Ignoring the privacy policies and terms of service policies of social networking services. In: *TPRC 44: The 44th Research Conference on Communication, Information and Internet Policy*
- Reich N, Micklitz HW, Rott P, Tonner K (2014) *European consumer law*. Intersentia
- Schulte-Nölke H, Twigg-Flesner C, Ebers M (2008) *EC consumer law compendium: The consumer acquis and its transposition in the member states*. Walter de Gruyter
- Sebastiani F (2002) Machine learning in automated text categorization. *ACM Comput Surv* 34(1):1–47
- Shulayeva O, Siddharthan A, Wyner A (2017) Recognizing cited facts and principles in legal judgements. *Artificial Intelligence and Law* 25(1):107–126
- Tsochantaridis I, Joachims T, Hofmann T, Altun Y (2005) Large margin methods for structured and interdependent output variables. *Journal of Machine Learning Research* 6:1453–1484