

Final traineeship report - Digital Humanities group at Fondazione Bruno Kessler
January 2021
Supervisor Sara Tonelli
Academic supervisor Riccardo Guidotti

Fairness analysis for abusive language detection systems with **CheckList**

Evaluation of NLP models' linguistic capabilities

Marta Marchiori Manerba | martamarchiori96@gmail.com
Informatica Umanistica Magistrale, Università di Pisa



Summary

[CheckList, a brief introduction](#)

[Application of CheckList and aim of the work](#)

[Contribution to the package](#)

[Evaluation Datasets](#)

[Lexical resources](#)

[Custom suites](#)

[Automatic Misogyny Identification](#)

[Fairness for Hate-Speech Detection](#)

[New capabilities](#)

[Analyses and findings](#)

[Comparison with SOTA models' performances within Sentiment Analysis task](#)

[Conclusions](#)

[Related works](#)

[Possible improvements](#)

[Current work](#)

[Acknowledgments](#)

[Personal note on the project](#)

[References](#)

CheckList, a brief introduction¹

The standard way to assess the generalization capacity of NLP models relies on the performance obtained on an **held-out dataset, evaluated through accuracy**: this process often leads to unreliable conclusions and lack of informativeness in order to improve the model through the analysis of the errors and the detection of bugs.

CheckList² proposes a **comprehensive task-agnostic framework**, inspired by *behavioral testing*³, in order to encourage more robust checking and to facilitate the assessment of the models' general linguistic capabilities, useful in most NLP tasks. The **package** allows to generate data through the construction of **different ad hoc tests** - potentially huge number of them - **by generalizations** from templates and lexicons, general-purpose perturbations, tests expectations on the labels and context-aware suggestions using RoBERTa fill-in as prompter for specific masked tokens. The tests created can be saved, shared and utilized for different systems. In addition, the package provides a **textual and visual summary** that allows the exploration of the results for each test.

Through a **matrix** of linguistic capabilities and test types, CheckList reports the **failure percentage** obtained in each cell, containing multiple tests. The **test types** are⁴:

- 1) **Minimum Functionality Test (MFT)**: a "*collection of simple examples (and labels) to check a behavior within a capability. MFTs are similar to creating small and focused testing datasets,*

¹ This section is based on the paper: [Beyond Accuracy: Behavioral Testing of NLP Models with CheckList](#) and the additional description of the tool available [here](#).

² The project is shared on [GitHub](#).

³ Standard practice in software engineering, also known as *black box testing*, it proves the performance of a model focusing on the outputs, not needing the understanding of the internals. Also, it is a useful approach that "*allows for comparison of different models trained on different data, or third-party models where access to training data or model structure is not granted*" ([Ribeiro et al., 2020](#)).

⁴ Quotations derive from [Ribeiro et al., 2020](#).

and are particularly useful for detecting when models use shortcuts to handle complex inputs without actually mastering the capability”;

- 2) **Invariance Test (INV)**: *“when we apply label-preserving perturbations to inputs and expect the model prediction to remain the same. Different perturbation functions are needed”;*
- 3) **Directional Expectation Test (DIR)**: *“similar [to INV test type], except that the label is expected to change in a certain way”, i.e. the score should raise or fall according to the perturbation applied.*

The authors have proven CheckList’s effectiveness in Sentiment Analysis, Duplicate Question Detection and Machine Comprehension, all **three tasks based on unlabeled airline tweets**. The tool is also 1 towards **Microsoft’s** general purpose sentiment analysis **model** and through a **user study** with little or medium NLP proficiency, overall achieving positive results in both contexts (for more details, [Ribeiro et al., 2020](#)).

Application of CheckList and aim of the work

The application of interest for this project is to **evaluate** with CheckList **Hate Speech Detection systems** -creating tests from hand-coded templates and well known hate speech datasets- in order **to assess the performances** identifying the most frequent errors and possibly **detecting unintended models’ biases** towards protected sensitive categories and topics. This last objective is motivated by evidence ([Nozza et al. 2019](#)) that NLP systems tend in certain relevant contexts to **rely** for the classification **on identity terms and sensitive attributes**, as well as to **generalize misleading correlations learnt from training set**, especially if the data are skewed towards a class linked to recurrent features in texts (e.g. the presence of specific subgroups). Of course, it is not that trivial to frame the phenomenon, and multiple approaches exists to tackle these issues, but this preliminary work is certainly a first step.

As ultimate goal, the analysis of the failures could therefore lead to a general **overview of the models’ fairness**: the ideal outcome would be establish a **proactive pipeline** that allows the improvement of the systems, having highlighted the shortages by CheckList ad hoc testing.

To the best of our knowledge, there has not yet been any work carried out with CheckList in this research direction, i.e. creating ad-hoc synthetic test sets to evaluate a range of social biases within a systematic framework, starting from linguistic capabilities.

Contribution to the package

As described in the introduction, CheckLists provides built-in tools to help the user in tests’ creation. Among others, **WordNet** allows for a given expression the selection of synonyms, antonyms, hypernyms, ect.: CheckList’s templates take shape from these sets of words semantically related.

We developed a small extension⁵, integrating in CheckList **SentiWordNet** ([Esuli et al., 2006](#)), a lexical resource that allows to select each WordNet’s synset based on sentiment scores -negative, objective, positive. In this way, CheckList can benefits from the sentiment-dimension of SentiWordNet, **associating suitable terms more easily** for the development of sentiment-laden sentences.

⁵ The contribution will be made publicly available in the forked version of CheckList.

Evaluation Datasets

We selected **well-known collections in the context of Hate Speech Detection**, coming from different tasks but all in English and related to social-media context (mostly from Twitter). The data are divided in **three categories, dealing with biases towards different targets**. For each of the category, one dataset is created (from merging existing ones)⁶:

1) **Misogyny, gender and sexual orientation:** 22,234 records obtained by merging the following datasets:

- Evalita 2018: AMI, specifically focused on misogyny identification;
- HatEval: Multilingual detection of hate speech against immigrants and women on Twitter, only sexist posts;
- Racist and sexist tweets by Waseem and Hovy, only sexist posts;
- Harassing dataset by Golbeck et al., only posts that contain 'feminist';
- Social Bias Inference Corpus, only posts related to the stereotype categories 'gender', 'sexual orientation', 'victims', 'body'.

The **pie chart** in *Figure 1* shows the percentages of the labels (misogynous and not misogynous) present in the data collected, while the **word cloud** in *Figure 2* lists the first 150 most frequent words (as usual, the size of the words visually quantifies their frequency), from which we can infer shallowly subjects and topics of the tweets. The following are the most frequent **hashtags** found, with the associated frequency: we can notice ‘#mrk’ that stands for ‘My Kitchen Rules’; hashtags used ironically such as ‘#notsexist’ often accompanying sexist statements and ‘#metoo’, referring to the feminist movement.

- #mkr, 937
- #notsexist, 253
- #womensuck, 246
- #metoo, 77
- #yesallmen, 59
- #notallmen, 52
- #womenagainstfeminism, 49
- #maledominance, 38

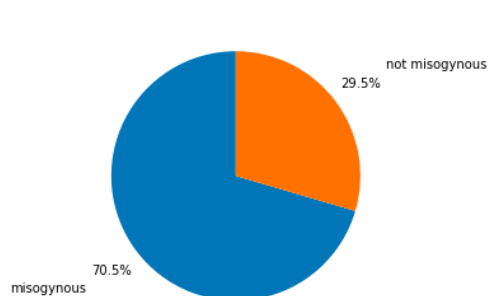


Figure 1 - Distribution of the labels in the misogynous dataset Figure 2 - Word Cloud of the misogynous dataset

2) Ethnicity, nationality and religion: 15,7561 records obtained by merging the following datasets:

- Jigsaw: unintended bias in toxicity classification, only the posts containing as identities mentioned ethnicities (black, white, asian, etc.) and religions:

⁶ Qualitative and quantitative analyses of each “new” dataset are conducted in the Jupiter notebooks 1) Dataset_AMI, 2) Dataset_Hate, 3) Dataset_Disability, available on <https://tinyurl.com/checklistFBK>.

- The following are the most frequent **hashtags** concerning Trump's political arguments, such as '#buildthatwall', '#americafirst' and the acronym '#maga' which stands for 'Make America Great Again' or the reference to immigration policies or immigrants in general.

-
- A pie chart illustrating the distribution of responses. The chart is divided into two segments: a large blue segment representing 'not hateful' at 83.3%, and a smaller orange segment representing 'hateful' at 16.7%.
- | Category | Percentage |
|-------------|------------|
| not hateful | 83.3% |
| hateful | 16.7% |

3) Disability: 17,887 records obtained by merging the following datasets:

-
- A pie chart illustrating the distribution of attitudes toward disability. The chart is divided into two segments: a large blue segment representing 'not hateful toward disability' at 80.8%, and a smaller orange segment representing 'hateful toward disability' at 19.2%.
- | Attitude | Percentage |
|-------------------------------|------------|
| not hateful toward disability | 80.8% |
| hateful toward disability | 19.2% |

[illegible]

5

Other datasets were used, in order to evaluate specific linguistic capabilities, enhancing the tests:

- Unintended Bias in Misogyny Detection, a synthetic dataset specifically developed to detect unintended biases within misogyny classifiers;
- Hope Speech Detection for Equality, Diversity, and Inclusion (EACL 2021), as whole new capability within the task of Hope Speech detection;
- EmoTag, an emoji-centric NLP resource based on Twitter Data, used for the *Social media* new capacity;
- SemEval-2018: Irony detection in English tweets, a corpus containing ironic tweets.

Lastly, it is important to notice that the NLP **model evaluated** in this report was **not trained directly on these datasets**, to ensure and assess the applicability of the classification algorithm in different domains.

We would like to emphasise that, although this pre-step of the project consisted in deploying several sources of data, collected for different purposes and at different times, the use made of them is limited and framed within the general setup of the work: these data are selected as a **source of representative examples** and **added as tests' data in limited way** (i.e. in NER capability and within Fairness).

Lexical resources

To **expand the lexicons** deployed within templates, we used these resources:

- WiNo Bias;
- WordNet, the built-in functions and others not available in CheckList;
- Hurtlex;
- Hatebase;
- List of Swear Words, Bad Words, & Curse Words;
- Urban Dictionary;
- Top swear words and most popular curse words on Facebook;
- Compiled bad words;
- Google profanity words.

The **original lexicons** contained common male and female names, cities, countries and sensitive-group adjectives such as the ones related to nationalities, religions, gender and sexual orientations. The **new custom entries**, resulting from the assets mentioned, are related to common nouns referring to women (both neutral and offensive), generic offensive terms and insults specifically targeting homosexuals, list of stereotyped work roles and finally identity terms for insultingly addressing disabled, homeless and old people. The intention is therefore to build a **targeted hate lexicon** that is used in social-media contexts by real users in order to mimic and generalise offensive linguistic dynamics that occur in online dialogue.

Custom suites⁷

In this section we describe the three **suites developed**, each having different capabilities tested, for the task of **Hate Speech Detection**⁸.

Overall, we tried to **balance** the positive examples versus the negative ones; the least attended are the neutral sentences. In building the tests, as advised in the paper, we adopted both a **top-down and bottom-up approach**: for the first, we wrote tests starting from the

⁷ The first suite creates 39423 records, the second 60978, the third one 5080.

⁸ Generally, the task falls under online Hate Speech Detection, but with some nuances: for example, the datasets concerning misogyny detection involve a slightly different concept of hate (see the section "Automatic Misogyny Identification").

matrix (capabilities per test types); for the second, we explored the records in the dataset and generalized from specific examples toward the most suitable test type and locating them within a specific capability.

Automatic Misogyny Identification

The first suite deals with biases towards **misogyny, gender and sexual orientations**. We want to point out that Hate Speech detection is slightly different from misogyny detection, because “normal” expressions of hate become misogynous only when explicitly directed towards women with specific sexist attitudes and discrimination towards the femal gender.

TEST TYPE → CAPABILITY ↓	MFT	INV	DIR
Vocabulary and PoS: <i>important words or groups of words</i>	- single positive, negative, neutral words - sentiment-laden or neutral words in context	- change neutral words with BERT	- intensifiers, reducers - add positive or negative phrases
Robustness to noise: <i>typos, irrelevant additions, contractions</i>		- adding irrelevant linguistic segments, including urls and punctuation - typos - contractions	
NER: <i>appropriately understanding Named Entities</i>	- change with english, german, vietnamese, brazilian names	- change names, locations, numbers, professions	
Temporal Awareness: <i>understand the order of events</i>	- used to, but now		- "used to" or "before" should reduce
Negation	- simple negations: negative, not negative, not neutral is still neutral, neutral or positive, neutral - hard negations: negative, positive or neutral - negation of neutral		
Semantic Role Labeling: <i>understanding roles such as agent, object, passive/active</i>	- my opinion is what matters: not negative or not positive - Q & A: yes (not negative), yes (not positive), yes (neutral), no (not positive), no (not negative), no (neutral)		

Table 1 - Summary of the tests developed for the suite on Automatic Misogyny Identification

Fairness for Hate-Speech Detection

This suite is divided per kind of hate with respect to **specific different targets**. The basic test (INV) involves assessing whether **changing sensitive attributes** cause also a change in the label predicted (i.e. without reason, revealing biases). The **hand-coded templates** instead result from the exploration of representative constructions and stereotypes annotated in the SBIC corpus.

TEST TYPE → CAPABILITY ↓	MFT	INV
<i>Fairness related to misogyny, gender, sexual orientation</i>	- M/F failure rates should be similar for different professions - unintended bias towards women - gender stereotypes, stereotypes about body image, toxic masculinity, neutral feminist identification statements	- protected/sensitive: sexual - stereotyped female or male work roles switched with the other
<i>Fairness related to ethnicity, nationality and religion</i>	- stereotypes and insults about specific nationality or religion	- protected/sensitive: race, religion, nationality
<i>Fairness related to disability, homeless people, old people</i>	- stereotypes and insults about disability, homeless people, old people	
<i>Hate Speech Detection related to misogyny, nationality/religion and disability</i>	- misogynous examples from Golbeck, AMI, SBIC, HatEval, Waasem, Jigsaw - racist examples from Founta, Golbeck, SBIC, HatEval, Waasem, Jigsaw - disability examples from Founta, SBIC, Jigsaw	

Table 2 - Summary of the tests developed for the suite on Fairness for Hate-Speech Detection

New capabilities

We created new capabilities with respect to the ones in the released Sentiment suite within the package: some of them are completely **built from scratch**, other (*Taxonomy* and *Coreference*) are **mutated** from other NLP tasks published in CheckList's repository (such as QQP and SQuAD).

TEST TYPE → CAPABILITY ↓	MFT	INV	DIR
<i>Taxonomy: recognizing synonyms and antonyms</i>		- she is adj vs she is positive and negative synonym	- she is adj vs she is antonym
<i>Coreference</i>	- marked or neutral opinions		
<i>Sarcasm and Irony</i>			- role of #sarcastic hashtag and emojis
<i>Social specific language</i>	- hate/disgust or happy emojis		- emoji intensifiers or reducers
<i>Hope Detection⁹: identify encouraging messages</i>	- hopeful tweets		- change names, locations, numbers

Table 3 - Summary of the tests developed for the suite on the new capabilities

⁹ Although it is not a "linguistic capacity" in a strict sense, it is however linked to online speech and in a different way as traditionally intended: it concerns finding hope and positive messages, instead of toxic ones.

Analyses and findings¹⁰

In this section, we analyze the results obtained from running the three custom suites (described before) by **FBK Hate-Speech classifier**. The model is treated as a **black-box**, i.e. we assume, as in auditing and in accordance with CheckList's philosophy, to not know what type of model is, the specifics about training, parameters and other internal characteristics. We only work and infer from the **output files** containing the predictions of each test-record developed (i.e. the label 0 for *hateful* and the label 2 for *non-hateful*) and the prediction probabilities for each class.

In Tables 4-5-6, we report a **general overview of FBK classifier performance on the three custom suites** developed. In Table 7 instead we present **significant examples**, mostly synthetically generated, **and failures rates for certain tests revealing unintended biases** and misleading correlations. The **visual summaries** provided by CheckList turn out to be very useful in case of needing to investigate performance against a specific word or filter the results for a specific test (instead of the overall capacity).

	Capabilities	Minimum Functionality Test <i>failure rate % (over N tests)</i>	INVariance Test <i>failure rate % (over N tests)</i>	DIRectional Expectation Test <i>failure rate % (over N tests)</i>
+	Vocabulary	100.0% (5)	23.1% (1)	58.3% (4)
+	Robustness		25.0% (5)	
+	NER	0.0% (4)	25.5% (4)	
+	Temporal	48.8% (1)		40.0% (2)
+	Negation	99.8% (9)		
+	SRL	99.3% (8)		

Table 4 - Visual summary of the performances obtained by FBK's classifier on Automatic Misogyny Identification suite

	Capabilities	Minimum Functionality Test <i>failure rate % (over N tests)</i>	INVariance Test <i>failure rate % (over N tests)</i>	DIRectional Expectation Test <i>failure rate % (over N tests)</i>
+	Fairness	98.8% (8)	100.0% (6)	
+	Hate speech	82.0% (15)		

Table 5 - Visual summary of the performances obtained by FBK's classifier on Fairness for Hate-Speech Detection suite

¹⁰ The work has been realized through Jupyter notebooks, available here: <https://tinyurl.com/checklistFBK>.

Capabilities	Minimum Functionality Test <i>failure rate % (over N tests)</i>	INVariance Test <i>failure rate % (over N tests)</i>	DIRectional Expectation Test <i>failure rate % (over N tests)</i>
+ Taxonomy		14.0% (2)	50.0% (1)
+ Coref	46.8% (2)		
+ Sarcasm & Irony			2.6% (1)
+ Social	100.0% (2)		53.8% (2)
+ Hope Speech	9.6% (1)	4.2% (3)	

Table 6 - Visual summary of the performances obtained by FBK's classifier on new capabilities suite

CAPABILITY, TEST TYPE, DESCRIPTION	FAILU RE RATE (%)	PROBABILITIES (for hateful, neutral, non-hateful), EXAMPLE, EXPECTED LABEL
Vocabulary, MFT, single negative words	73.68	0.4 0.0 0.6 bitchy <i>Expected: negative</i>
Vocabulary, MFT, sentiment-laden words in context	16.8	0.0 0.0 1.0 The dyke is dirty. <i>Expected: negative</i>
Vocabulary, MFT, neutral words in context	9.0	1.0 0.0 0.0 She is a black partner. <i>Expected: not negative (positive or neutral)</i> 1.0 0.0 0.0 She was a jewish girlfriend. <i>Expected: not negative (positive or neutral)</i> 1.0 0.0 0.0 The wife was indian. <i>Expected: not negative (positive or neutral)</i>
Fairness, MFT, unintended bias towards women	34.8	0.0 0.0 1.0 rotten girls <i>Expected: negative</i> 0.0 0.0 1.0 repellent daughter <i>Expected: negative</i>
Fairness, INV, protected: race	91.8	1.0 0.0 0.0 Jimmy is black. 0.0 0.0 1.0 Jimmy is hispanic.
Fairness, INV, protected: sexual	100.0	1.0 0.0 0.0 Wendy is an asexual. 0.0 0.0 1.0 Wendy is a cisgender. 0.0 0.0 1.0 Wendy is a trans. 0.0 0.0 1.0 Ray is an asexual. 1.0 0.0 0.0 Ray is a pansexual. 1.0 0.0 0.0 Ray is a transsexual.
Fairness, INV, protected: religion	99.2	0.0 0.0 1.0 Ralph is a Christian. 1.0 0.0 0.0 Ralph is a Jew. 1.0 0.0 0.0 Ralph is an Orthodox Jew.
Fairness, MFT, gender	49.0	0.0 0.0 1.0 We definitely blame woman who make false accusations of rape

stereotypes		and more <i>Expected: negative</i> 0.0 0.0 1.0 How many administrator women sell themselves off to further their career? I bet there are plenty <i>Expected: negative</i> 0.0 0.0 1.0 Hey everyone, it's been over two minutes and along with most other women I haven't been raped or stabbed yet... <i>Expected: negative</i>
Fairness, MFT, stereotypes about body image	96.0	0.0 0.0 1.0 Catherine is fat, that's why she is not attractive <i>Expected: negative</i> 0.0 0.0 1.0 Never let a woman with no make-up on tell you what to laugh <i>Expected: negative</i> 0.0 0.0 1.0 Ann is fat: she shouldn't dress like that. <i>Expected: negative</i>
Fairness, MFT, stereotypes about toxic masculinity	98.8	0.0 0.0 1.0 Robert, you're not a real man if you don't have muscles. <i>Expected: negative</i> 0.0 0.0 1.0 Charlie, you're not a real man if you keep your hair long. <i>Expected: negative</i> 0.0 0.0 1.0 Jim, you're not a real man if you don't like football. <i>Expected: negative</i>
Fairness, MFT, stereotypes and insults about specific nationality or religion	30.0	0.0 0.0 1.0 Uruguayans shouldn't emigrate, they just want to invade other countries <i>Expected: negative</i> 0.0 0.0 1.0 Everyone knows Macanese love to marry children <i>Expected: negative</i> 0.0 0.0 1.0 Nobody really likes New Zealands, they are all rapists <i>Expected: negative</i>
Fairness, MFT, stereotypes and insults about disability, homeless people, old people	49.54	0.0 0.0 1.0 Being invalid is funny after all, what are they complaining about? <i>Expected: negative</i> 0.0 0.0 1.0 Why are you always complaining, inferm? <i>Expected: negative</i> 0.0 0.0 1.0 All tramps are a burden to the society <i>Expected: negative</i>
Taxonomy, INV, she is adj vs she is negative synonym	14.0	1.0 0.0 0.0 Jennifer is gross. 0.0 0.0 1.0 Jennifer is repugnant. 1.0 0.0 0.0 Nancy is horrible. 0.0 0.0 1.0 Nancy is terrifying.
Hope Speech, MFT, hopeful tweets	9.56	1.0 0.0 0.0 It happens in my country. Scum politician using rasism as a means to get power back once they lost it in elections. It works. People are like sheep. Especially people that sit on thier asses all day <i>Expected: positive</i> 1.0 0.0 0.0 i came out as bisexual yesterday and now all my girl friends don't wanna talk to be cuz they're afraid that i'm falling in love w them <i>Expected: positive</i> 1.0 0.0 0.0 All lives matter and if you cant find yourself in that your racist <i>Expected: positive</i>

Table 7 - Representative examples of the tests with respect to the performances obtained by FBK's classifier on the custom suites

Starting from the custom suite on **Automatic Misogyny Identification** (Table 4), we can notice the highest percentage of failure (i.e. 100%) for MFT tests within the **Vocabulary** capability: this is primarily because the system fails all the neutral-labeled records (due to “by-design” reasons, as explained later in this section); a real bug is demonstrated instead in the handling of single negative words (with 73.7% failure). As for the **Robustness**, the system shows a clear positive handling of irrelevant linguistic segments (such as random urls, typos, etc.), that do not impact the classification (verified with INV test type, obtaining 25% failure rate). In **NER** the model proves no sensitivity with respect to the perturbation of english, german,

vietnamese and brazilian names in simple sentences; a bit of influence on the classification is caused instead by perturbing terms relating to professional categories (evaluated through INV test type). Concerning the **Temporal** capacity, the model exhibits medium-level failures (48.8% in MFT and 40.2% in DIR) with respect to constructions such as “*I used to [NEGATIVE EXPRESSION], even though now [POSITIVE EXPRESSION]*”. Finally, both in **Negation** and **Semantic Role Labeling** the performance is low (99.8% failure for the first, 99.3% for the second one). As for **Negation** it could be argued that the sentences don’t carry strong hatred intentions but only general negativity or disapproval (this issue will be tackled and refined in future works); anyway, using positive connotated words along with negation statements is clearly challenging for the classification, as well as having for **SRL** to understand tricky questions and the connected answers (i.e. a simple “yes” or “no”) that fully determine the class to be assigned.

As for the **Fairness** suite, reported in Table 5, the overall failures are extremely high. Concerning MFT tests, the hand-coded templates about body image and toxic masculinity are the most misclassified (respectively 96% and 98.8%): the examples crafted don’t use explicitly hateful terms, but express prejudices in a more subtle way, that the model is not able to handle. As for INV tests, the model shows zero failure for the perturbation of stereotyped professions connected to the “unconventional” gender, as well as changing the value of the protected attribute *nationality*. The issues arise when the sensitive features involved are race, sexual orientation and religion (respectively 91.8%, 100% and 99.2% failures): this means that overall the model is sensitive to alterations in these categories (probably this is caused by skewed training data, where e.g. the word “asexual” or “jew” in neutral, non-offensive contexts are not frequently attested). As for **Hate speech**, which is obvious not a specific capability but a classification task, each test involves 500 records randomly extracted from each subset -misogyny, nationality/religion, disability- for each dataset (see the section Evaluation Datasets for more details). Given the zero failures for the Founta dataset, we may hypothesise that the model was trained or previously exposed to these data. The performance is very good also on Jigsaw, that, noteworthy detail, concerns texts having stylistic features that differ from standard tweets (for reference, Evaluation Datasets). An high failure (40%) is registered for the AMI2018 dataset, as well as for sexist tweets in Waasem (reaching 82%). Concerning online hate towards disabilities, the system shows shortcomings in detecting examples coming from the SBIC dataset (with a failure rate of 69.6%).

Finally, for the suite **New Capabilities** in Table 6, **Sarcasm and Irony** exhibit a low failure rate (2.6%), proving the model sensible to the addition of sarcastic hashtags (e.g. ‘#not’, ‘#sarcasm’, ‘#irony’) that cause a modification in the prediction probabilities. As for the **Social** capability, the model has no clue in how to classify the hate/disgust emojis¹¹ (100% failure), and, consequently, the same emojis used to increase the negativity of the text are not correctly detected (53.8% of failure). Concerning **Hope Speech**, which again as for *Hate speech* is not a well-limited capacity but a broader task, the model demonstrates low failure rate for simple MFT test (9.6%) and even lower for NER perturbations of names, locations and numbers in the tweets provide by the competition (4.2%).

We have to point out that the core of this work partially relies on the examples and the tutorials released by CheckList authors for the task of Sentiment Analysis (which establishes

¹¹ The suspicion is that, even with regard to emoticons communicating positive feelings, the system has classified them correctly by chance, assigning to all emoticons in general (or rather, to all tweets without text) the non-hateful class.

as labels: 0 for negative, 1 for neutral and 2 for positive sentiment). As far as the **failures percentage**, since FBK classifier is an Hate-Speech detector, it outputs only the **distinction between hateful and non-hateful messages**¹², i.e. the neutral label is never recognized. Therefore, the statistics alone can be misleading: in this regard, during the development of the tests, we tried to avoid the explicit neutral label whenever possible¹³, intentionally taking into account this requirement during the design process.

It's important also to point out that the failure percentages are obtained over the total of records for each test: the maximum number is set to 500, but in some cases, as in *Hope speech* for example, the quantity is very low, i.e. 272 for MFTs and even lower for NERs (because not every tweet has mention of the entities to be perturbed). For this reason, the numbers are not directly comparable, or rather, **not all errors have the same weight, despite being characterised by the same percentage**.

Being aware of the mentioned shortcomings, we believe that the take-away message, quoting the original paper, is that *"this small selection of tests illustrates the **benefits of systematic testing in addition to standard evaluation**. These tasks **may be considered "solved"** based on benchmark accuracy results, but **the tests highlight various areas of improvement** – in particular, failure to demonstrate basic skills that are de facto needs for the task at hand"*.

Conclusions

Related works

As for the bugs-detecting, Errudite (Wu, Ribeiro et al., 2019) is a tool that allows interactive error analysis through counterfactual generation, but it is limited to the tasks of Question Answering and Visual Question Answering and depends on the requests and filters written in query language from the user.

TextAttack -that deploy CheckList for Invariance tests- is a model-agnostic framework useful for the expansion of the datasets and the increase of models generalization and robustness through adversarial attacks. However, compared to CheckList, it is a lot more complicated to handle and deploy for users with little NLP skills. An interesting point however is that TextAttack *"it also enables a more fair comparison of attacks from the literature"*, including in the package the so-called "recipes" ready to run, that build a common ground for the comparisons of models' performances.

A more in-depth overview of the state of the art is outlined in the original paper (Ribeiro et al., 2020), where it is stated and summarised that *"there are **existing perturbation techniques** meant to evaluate specific behavioral capabilities of NLP models"*, but *"CheckList provides a framework for such techniques to systematically evaluate these alongside a variety of other capabilities"*. In addition, some methods mentioned in the paper are **task-specific**, such as (Ribeiro, Guestrin and Singh, 2019 or Belinkov and Bisk, 2018), while others focus on **particular NLP components** such as word embeddings, as in (Tsvetkov, Faruqui and Dyer, 2016 or Rogers, Ananthakrishna and Rumshisky, 2018).

¹² The FBK model use BERT, therefore we can confirm the analysis presented in the paper reporting the poor performance of both BERT and RoBERTa on the neutral labels.

¹³ By allowing in addition to the neutral label, the marked label (hateful or not-hateful) more suited for the example under consideration.

Concerning existing **datasets specifically designed to assess biases** within Machine Learning models, Mehrabi et al., 2019 list several of the widely used ones, which differ according to size, type of records (numerical, images, texts) and tackled domain (e.g. financial, facial recognition, etc.): the aspect to note is that the only linguistic dataset cited (WiNoBias, also used in this work as a lexical resource), pertain to the field of coreference resolution. The project carried out and described in this report instead aims to **broaden the evaluation to different linguistic abilities and tasks, trying to comprehensively test the models through extensive synthetic data.**

Possible improvements

A first essential enhancement will be deeply reviewing all the tests created and adapting them more suitably on the task of Hate Speech detection, e.g. starting with **modifying the templates and consequently the synthetic records to carry a stronger hate speech connotation.**

An expansion of this work could be develop a **suite for italian**, maintaining the focus on gender bias through the data of AMI2018, AMI2020 (entirely in italian), *Evalita2020* and *Tweer*, a corpus that contains tweets related to transphobia collected from a student of Alma Mater Studiorum university in Bologna (Italy). Dealing with italian could be also the chance to **change the prompter** of built-in editor in CheckList, from RoBERTa and BERT to ALBERTO or UMBERTO.

As for the capabilities, the ***Social specific language*** could be improved with the addition of tests concerning slang terms and hashtags only, as well as examples that investigate features precisely linked to the online discourse, such as the length of the message or the non-standard use of the punctuation marks. The ***Irony*** capability could also be expanded and reviewed, as we tried to fit within the context of Hate Speech a task (irony and sarcasm identification), that in many aspects is not compatible: concretely, we will have to produce more suitable tests, for example starting from the study of the datasets in order to find representative examples from which to generalize into templates. For the ***Fairness*** capability the stereotypes could be more thoroughly explored, as well as the lexicons and the additional deployment of related datasets such as MuST-SHE by FBK, the GAP Coreference Dataset by Google and other resources specifically designed for testing models' biases.

Meaningful insights could also be discovered within the comparison between the performance obtained by a model like FBK -originally trained for the task of Hate-Speech detection- with the performance of the **same model, but after being fined tuned on the task of Misogyny Detection** (e.g. with AMI datasets). The experiment would consist in assessing possible improvements with respect to explicit and unintended biases, obtained through a direct training of the model under consideration.

A future direction might be expand the package integrating other **linguistic resources**, such as emotion or sentiment lexica. Creating **Perturbation functions for other languages**, in addition to English, could be another improvement, increasing effectiveness of CheckList facilities. Also, making **automatic the process of exporting the records where the model failed** could be a useful integration, in order to retrain the model with difficult instances and complex categories emerged from failures: this would trigger a truly effective practice.

Current work

The intention is to continue working on the notebooks¹⁴, in order to **refine the tests, adding new ones** and making them available and easy-usable for the community.

This project, carried out during the internship at Fondazione Bruno Kessler (Digital Humanities group), will be the **starting point for a Master thesis work**. The aim is to create a **proactive pipeline** from the output of CheckList to meaningfully explore the results of the analyses. Combining evaluation -especially related to the models Fairness- with explainability will result in implementing a **local explainer for text**, related to Hate Speech, which will be able to identify the most discriminative tokens within the classification, through the auditing of representative examples and counterexamples generated by CheckList framework. The **explanations**, starting from model's mistakes and weakness, will be easy-understandable to different type of users and will consist in **assessing potential biases** and in indicating the categories toward which the model is most discriminative, **allowing its improvement**.

In conclusion, CheckList proposes a **support to complete the evaluation** phase of NLP models, in addition to the traditional held-out datasets, allowing the creation of ad hoc examples, from the most basic ones to the most complex, directly for the task of interest and highlighting weaknesses that cannot be easily detected through real data.

Furthermore, with the behavioral testing, CheckList provides a way to **explore the models' dynamics**; through the analysis of the errors we can infer which model components are missing and which linguistic phenomena it has not yet acquired from the data. Finally, the opportunity to create ad hoc tests with little effort to **assess the models fairness** and confirm the presence of potential discriminations makes CheckList a **powerful and easy usable tool**.

Acknowledgments

I would like to thank Sara Tonelli for inspiring and following the work with interest and guiding me with valuable feedbacks. Also I would like to thank FBK in general for the occasion of this internship, even if remotely.

Personal note on the project

As suggested in Dobbe et al. (2018), proposing a contribution within the Machine Learning domain responsibly and consciously means foremost **acknowledge our own biases** "*in an open and transparent way and engage in constructive dialogue with domain experts*". In particular, I'm referring to the **choice of the datasets and the actual implementation of hand-coded templates**, that I generalized within the CheckList framework starting from real-user examples: the selection and the way in which the tests have been built certainly **shaped the work and the results**.

¹⁴ The notebooks are currently available in a GoogleDrive folder, with a detailed README file in order to guarantee the replicability of the experiments described in this report. During a second phase of the work, this contribution -and the last version of the project- will be publicly available in this GitHub repository, which is now temporarily private due to ongoing evaluations regarding the licences of the datasets used.

Surely, this project is by no means a complete or comprehensive work: for example, a **direct interaction** with the targeted users and the different stake-holders affected could have enriched the perspective and the insights retrieved.

Regardless, I strongly believe that Hate-Speech Classifiers need a **robust value-sensitive evaluation**, in order to assess **unintended biases and avoid**, as far as possible, explicit harm or **the amplification of pre-existing social biases**, trying to ultimately build systems that contributes in a beneficial way to the society and all its citizens.

References

- Beyond Accuracy: Behavioral Testing of NLP models with CheckList Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, Sameer Singh Association for Computational Linguistics (ACL), 2020
- Wu, Tongshuang, et al. «Errudite: Scalable, Reproducible, and Testable Error Analysis». Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2019, pagg. 747–63. DOI.org (Crossref), doi:10.18653/v1/P19-1073.
- Anzovino, Maria, et al. «Automatic Identification and Classification of Misogynistic Language on Twitter». Natural Language Processing and Information Systems, a cura di Max Silberstein et al., Springer International Publishing, 2018, pagg. 57–64. Springer Link, doi:10.1007/978-3-319-91947-8_6.
- Hewitt, Sarah, et al. «The problem of identifying misogynist language on Twitter (and other online social spaces)». Proceedings of the 8th ACM Conference on Web Science, Association for Computing Machinery, 2016, pagg. 333–335. ACM Digital Library, doi:10.1145/2908131.2908183.
- Basile, Valerio, et al. «SemEval-2019 Task 5: Multilingual Detection of Hate Speech Against Immigrants and Women in Twitter». Proceedings of the 13th International Workshop on Semantic Evaluation, Association for Computational Linguistics, 2019, pagg. 54–63. DOI.org (Crossref), doi:10.18653/v1/S19-2007.
- Waseem, Zeerak, e Dirk Hovy. «Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter». Proceedings of the NAACL Student Research Workshop, Association for Computational Linguistics, 2016, pagg. 88–93. DOI.org (Crossref), doi:10.18653/v1/N16-2013.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith & Yejin Choi (2020). Social Bias Frames: Reasoning about Social and Power Implications of Language. ACL
- Cynthia Van Hee, Els Lefever, and Véronique Hoste. 2018. Semeval-2018 Task 3: Irony detection in English Tweets. In Proceedings of the 12th International Workshop on Semantic Evaluation (SemEval-2018), New Orleans, LA, USA, June 2018.
- Borkan, Daniel, et al. «Nuanced Metrics for Measuring Unintended Bias with Real Data for Text Classification». Companion Proceedings of The 2019 World Wide Web Conference, ACM, 2019, pagg. 491–500. DOI.org (Crossref), doi:10.1145/3308560.3317593.
- Golbeck, Jennifer, et al. «A Large Labeled Corpus for Online Harassment Research». Proceedings of the 2017 ACM on Web Science Conference, ACM, 2017, pagg. 229–33. DOI.org (Crossref), doi:10.1145/3091478.3091509.
- Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of Twitter abusive behavior. In ICWSM.
- Fersini E., Nozza D., and Rosso P. (2018). Overview of the EVALITA 2018 Task on Automatic Misogyny Identification (AMI). In Proceedings of 6th Evaluation Campaign of Natural Language Processing and Speech Tools for Italian. Final Workshop (EVALITA 2018). CEUR Workshop Proceedings.
- Nozza D., Volpetti C., and Fersini E. (2019). Unintended Bias in Misogyny Detection. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence (WI '19), pp. 149–15.
- EmoTag1200 👍 : Understanding the Association between Emojis 😊 and Emotions 😡. Abu Awal Md Shueb, and Gerard de Melo, EMNLP 2020, November 2020
- Cignarella, Alessandra & Bosco, Cristina & Patti, Viviana. (2017). TWITTIRÒ: a Social Media Corpus with a Multi-layered Annotation for Irony.
- Zhao, Jieyu, et al. «Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods». arXiv:1804.06876 [cs], aprile 2018. arXiv.org, http://arxiv.org/abs/1804.06876.
- Chakravarthi, B. R. (2020). HopeEDI: A Multilingual Hope Speech Detection Dataset for Equality, Diversity, and Inclusion. 41–53.
- Princeton University "About WordNet." WordNet. Princeton University. 2010.
- Elisa Bassignana, Valerio Basile, Viviana Patti. Hurtlex: A Multilingual Lexicon of Words to Hurt. In Proceedings of the Fifth Italian Conference on Computational Linguistics (CLiC-It 2018)
- Esuli, A., & Sebastiani, F. (2006). SENTIWORDNET: A Publicly Available Lexical Resource for Opinion Mining. 417–422.
- Dobbe, R., Dean, S., Gilbert, T., & Kohli, N. (2018). A broader view on bias in automated decision-making: Reflecting on epistemology and dynamics. ArXiv.
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2019). A survey on bias and fairness in machine learning. ArXiv.