

15th January 2021

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

Evaluation of NLP models' linguistic capabilities

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

by Marta Marchiori Manerba, Informatica Umanistica Magistrale, Università di Pisa

CheckList

- **task-agnostic framework** to encourage more robust checking **-beyond Accuracy on an held-out dataset-** and facilitate the assessment of the models' general linguistic capabilities
- the **package** allows to **generate data** through the construction of **different ad hoc tests** by **generalization** from
 - templates and lexicons,
 - general-purpose perturbations,
 - tests expectations on the labels
 - 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
- **textual and visual summary** for the exploration of the results

Concerning CheckList's tests creation

Approach	Idea	Advantage	Disadvantage
Scratch	Write tests manually	High Quality	Low Coverage, Expensive, Time-consuming
Perturbation Function	Apply perturbation to texts	Lots of Automated Tests	Low Quality
Template	Use templates and generate many variations	Balance of Quality and Quantity	Need to brainstorm Templates

Matrix of linguistic capabilities and test types

Through a **matrix of linguistic capabilities and test types**, CheckList reports the **failure percentage** obtained in each cell, that contains multiple tests

<u>Capability / Test</u>	<u>Minimum Functionality Test(MFT)</u>	<u>Invariance Test(INV)</u>	<u>Directional Expectation Test(DIR)</u>
<u>VOCABULARY</u>	15.0%	16.2%	34.6%
<u>NER</u>	0.0%	20.8%	-
<u>NEGATION</u>	76.4%	-	-
...			

Tables from the description of the tool available on <https://amitnss.com/2020/07/checklist/>

Test types

1. Minimum Functionality Test (MFT) \Rightarrow “collection of simple examples (and labels) to check a behavior within a capability. MFTs are similar to creating small and focused testing datasets, and are particularly useful for detecting when models use shortcuts to handle complex inputs without actually mastering the capability”

Text	Expected	Predicted	Pass
I didn't love the flight	negative	positive	X

Template: I {NEGATION} {POS_VERB} the {THING}

Test types

2. Invariance Test (INV) \Rightarrow “when we apply label-preserving perturbations to inputs and expect the model prediction to remain the same. Different perturbation functions are needed”

Text	Expected	Predicted	Pass
We had a safe travel to <u>Chicago</u>		positive	X
We had a safe travel to <u>Dallas</u>	positive	neutral	



Test types

3. Directional Expectation Test (DIR) \Rightarrow “*similar, 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.

Text	Expected	Predicted	Pass
Service wasn't great		negative	X
Service wasn't great.			
You are lame	negative	neutral	

The diagram illustrates a Directional Expectation Test (DIR). It shows a perturbation applied to a text sample. The original text is "Service wasn't great", which is predicted to be "negative". The perturbed text is "Service wasn't great. You are lame", which is predicted to be "neutral". The expected change in sentiment is from negative to neutral. However, the test fails (indicated by a red 'X') because the predicted sentiment change is not in the expected direction (from negative to neutral).

Application of CheckList and aim of the work



- Evaluate Hate Speech Detection systems to identify errors and unintended models' biases (Nozza et al, 2019)
 - There has not yet been any work carried out with CheckList in this direction, i.e. **creating ad-hoc synthetic test sets to evaluate a range of social biases within a systematic framework**, starting from linguistic capabilities

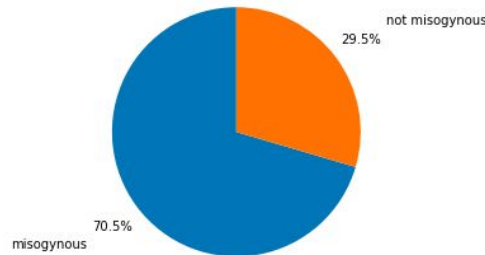
Evaluation Datasets

Well-known collections within **Hate Speech Detection**, from **different tasks** but all in **English** and related to **social-media** (mostly from Twitter):

- the data are divided in **three categories**, dealing with biases towards different targets:
 1. Misogyny, gender and sexual orientation
 2. Ethnicity, nationality and religion
 3. Disability
- for each of the category, **one dataset** is created (from merging existing ones)

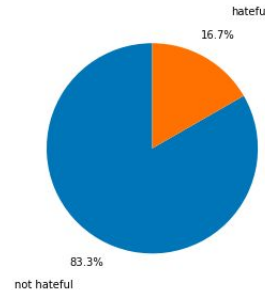
1) Misogyny, gender and sexual orientation: 22,234 records

- 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'.



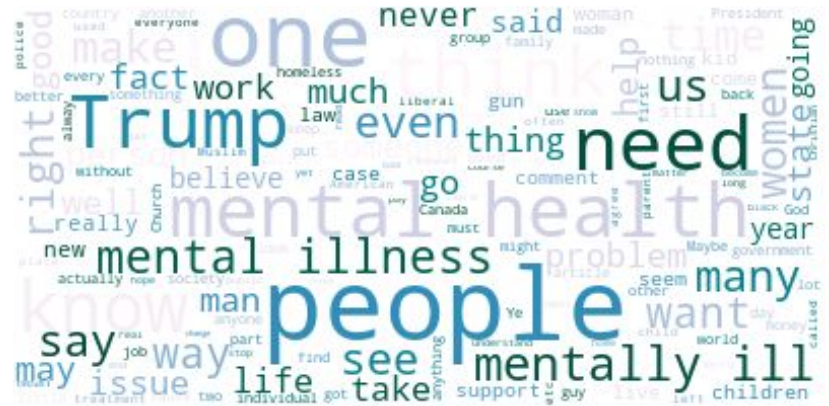
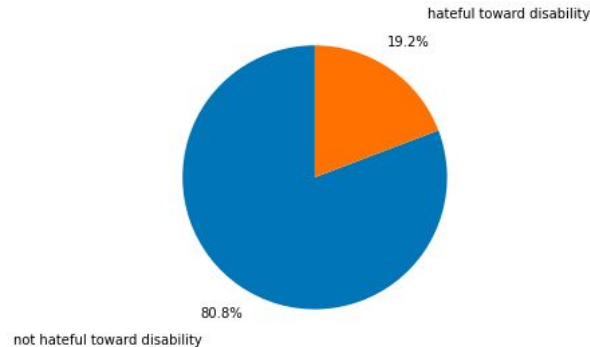
2) Ethnicity, nationality and religion: 15,7561 records

- Jigsaw: unintended bias in toxicity classification, only the posts containing as identities mentioned ethnicities (black, white, asian, ecc) and religions;
- Founta, only posts that contain mentions to nationalities, religions, immigrants and frequent hashtags like '#whitegenocide', '#fuckniggers', '#WhitePower', '#WhiteLivesMatter';
- HatEval: Multilingual detection of hate speech against immigrants and women on Twitter, only racist posts;
- Racist and sexist tweets by Waseem and Hovy, only racist posts;
- Harassing dataset by Golbeck et al. filtered as for Founta;
- Social Bias Inference Corpus, only posts related to the stereotype categories 'race', 'culture', 'social'.



3) Disability: 17,887 records

- Jigsaw: unintended bias in toxicity classification, only the post containing mentions to disabilities (physical, intellectual or psychiatric);
- Social Bias Inference Corpus, only post related to disabled;
- Founta, only posts that contain mentions to disabilities.



Other datasets

- 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 connected to 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.

Lexical resources

- To expand the lexicons deployed in templates, we use **other linguistic resources**
 - The **new custom entries** are related to common nouns referring to women (both neutral and offensive), generic offensive terms and insults targeting homosexuals, list of work roles related to a particular gender and identity terms for disabled, homeless and old people
 - 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.

Premises

- The model evaluated with CheckList was **not trained directly on these datasets**, to ensure and assess the applicability of the classification algorithms in different domains and datasets
- These datasets are **added as tests' data and used for templates**, but in limited way (i.e. in NER capability)
 - Within the *Fairness* suite we created tests containing 500 cases for each dataset, extracted randomly

Custom suites

- **three suites** developed, each having **different capabilities** tested, in the context of Hate Speech detection
- we try to **balanced** the positive examples versus the negative ones; the least attended are the neutral sentences
- in building the tests we adopted both a **top-down and bottom-up approach**:
 - for the first, we write tests starting from the matrix (capabilities per test types);
 - for the second, we explore the records in the dataset and generalize from specific examples to test type and locating them into a specific capability
- **output**: visual summary, textual summary, csv

1. Automatic Misogyny Identification

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)		

2. Fairness for Hate Speech Detection

The Fairness capability is divided **per kind of hate with respect to different targets**

- the basic test (INV) involves assessing whether **changing sensitive attributes** change also the label (i.e. without reason, revealing biases)
- the hand-coded templates result from the exploration of representative constructions and **stereotypes annotated in SBIC dataset**

TEST TYPE → CAPABILITY ↓	MFT	INV
<i>Fairness related to misogyny, gender, sexual orientation</i>	<ul style="list-style-type: none"> - 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 	<ul style="list-style-type: none"> - protected/sensitive: sexual - stereotyped female or male work roles switched with the other
<i>Fairness related to ethnicity, nationality and religion</i>	<ul style="list-style-type: none"> - stereotypes and insults about specific nationality or religion 	<ul style="list-style-type: none"> - protected/sensitive: race, religion, nationality
<i>Fairness related to disability, homeless people, old people</i>	<ul style="list-style-type: none"> - stereotypes and insults about disability, homeless people, old people 	
<i>Hate Speech Detection related to misogyny, nationality/religion and disability</i>	<ul style="list-style-type: none"> - misogynous examples from Golbeck, AMI, SBIC, HatEval, Waasem, Jigsaw - racist examples from Founta, Golbeck, SBIC, HatEval, Waasem, Jigsaw - disability examples from Founta, SBIC, Jigsaw 	

3. New capabilities

Some of them are completely built **from scratch**, other (*Taxonomy* and *Coreference*) are **mutated from other tasks** in the repository (such as *QQP* and *SQuAD*)

TEST TYPE → CAPABILITY ↓	MFT	INV	DIR
<i>Taxonomy</i> : recognizing synonyms and antonyms		- 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 Detection</i> ² : identify encouraging messages	- hopeful tweets		- change names, locations, numbers

Analyses and findings: *premises*

- The FBK model is treated as a **black-box**
 - We only work and infer from the output files containing the predictions of each test-record developed (i.e. the label 0 for hateful, 1 for neutral and 2 for non-hateful) and the prediction probabilities for each of the three classes
- The **visual summaries** provided by CheckList for each suite 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

1. Automatic Misogyny Identification suite

- *Vocabulary*: highest percentage of failure (i.e. 100%) for MFT tests
- *Robustness*: positive handling of irrelevant linguistic segments
- *NER*: optimal handling of perturbations of names
- *Temporal*: medium-level failures
- *Negation* and *SRL*: low performance in both
 - using positive connotated words along with negation statements is clearly challenging

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)		

Examples

- Vocabulary, MFT, sentiment-laden words in context
 - 0.0 0.0 1.0 The dyke is dirty. | Expected: negative
- Vocabulary, MFT, neutral words in context
 - 1.0 0.0 0.0 She is a black partner. | Expected: not negative (positive or neutral)
 - 1.0 0.0 0.0 She was a jewish girlfriend. | Expected: not negative (positive or neutral)
 - 1.0 0.0 0.0 The wife was indian. | Expected: not negative (positive or neutral)

2. Fairness for Hate-Speech Detection suite

- For MFT tests, the hand-coded templates about body image and toxic masculinity are the most misclassified
- For INV tests, the model shows zero failure for the perturbation of stereotyped professions:
 - issues arise when are involved race, sexual orientation and religion
- For *Hate speech*:
 - zero failures for Founta, good performance also on Jigsaw
 - high failure for AMI2018, Waasem and SBIC concerning disabilities

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)		

Examples

- Fairness, INV, protected: race
 - 1.0 0.0 0.0 Jimmy is black.
 - 0.0 0.0 1.0 Jimmy is hispanic.
- Fairness, INV, protected: sexual
 - 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
 - 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 stereotypes
 - 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... | Expected: negative
- Fairness, MFT, stereotypes about body image
 - 0.0 0.0 1.0 Never let a woman with no make-up on tell you what to laugh | Expected: negative
 - 0.0 0.0 1.0 Ann is fat: she shouldn't dress like that. | Expected: negative
- Fairness, MFT, stereotypes about toxic masculinity
 - 0.0 0.0 1.0 Charlie, you're not a real man if you keep your hair long. | Expected: negative
 - 0.0 0.0 1.0 Jim, you're not a real man if you don't like football. | Expected: negative

and more...

- Fairness, MFT, stereotypes and insults about specific nationality or religion
 - 0.0 0.0 1.0 Uruguayans shouldn't emigrate, they just want to invade other countries | Expected: negative
 - 0.0 0.0 1.0 Everyone knows Macaneses love to marry children | Expected: negative
 - 0.0 0.0 1.0 Nobody really likes New Zealands, they are all rapists | Expected: negative

- Fairness, MFT, stereotypes and insults about disability, homeless people, old people
 - 0.0 0.0 1.0 Being invalid is funny after all, what are they complaining about? | Expected: negative
 - 0.0 0.0 1.0 Why are you always complaining, inferm? | Expected: negative
 - 0.0 0.0 1.0 All tramps are a burden to the society | Expected: negative

3. *New capabilities suite*

- *Sarcasm & Irony*: low failure rate
- *Social*: no clue in how to classify the hate/disgust emojis (and probably happy ones)
- *Hope Speech*: low failure rate

Capabilities	Minimum	INVariance	DIREctional
	Functionality Test failure rate % (over N tests)	Test failure rate % (over N tests)	Expectation Test failure rate % (over N tests)
+ 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)	

Examples

- Taxonomy, INV, she is adj vs she is negative synonym
 - 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
 - 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 | Expected: positive
 - 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 | Expected: positive
 - 1.0 0.0 0.0 All lives matter and if you cant find yourself in that your racist | Expected: positive

Considerations

*“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**”*
(Ribeiro et al., 2020)

- Since FBK classifier is an Hate-Speech detector, it outputs only the distinction between hateful and non-hateful messages: the statistics alone can be misleading
- The percentages are obtained over the total of records for each test, that varies: not all errors have the same weight, despite being characterised by the same percentage



Related works

*“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” (Ribeiro et al., 2020)*

- **Errudite** (Wu, Ribeiro et al., 2019) allows **interactive error analysis**, for the tasks of QA and VQA and depends on user **requests and filters**
- **TextAttack** is a model-agnostic framework for the expansion of the **datasets**, increasing model **generalization** and **robustness**
- Existing **datasets specifically designed to assess biases**: Mehrabi et al., 2019 list the widely used ones:
 - the only linguistic dataset cited (WiNoBias, also used in this work as a lexical resource), pertain to the field of coreference resolution



Possible improvements

With respect to tests and suites

- Develop a **suite for italian**, maintaining the focus on gender biases and **changing the prompter** to ALBERTO or UMBERTO
- As for the capabilities,
 - the “Social specific language” could be improved with tests concerning **slang terms and hashtags** only
 - for “Fairness” the **stereotypes** could be more thoroughly explored, as well as **lexicons** and representative **datasets**, such as MuST-SHE by FBK or the GAP Coreference Dataset by Google



Possible improvements

With respect to the package

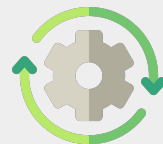
- Comparison between a model before and **after being fined tuned on the task of Misogyny Detection**
- Making automatic the process of **exporting the misclassified records** in order to retrain the model
- **Perturbation functions for other languages**, in addition to English
- Other **linguistic “enhancements”**:
 - domain-specific word embeddings
 - use of emotion or sentiment lexica
 - features related to the message (e.g. length, punctuation marks, etc.)



Ongoing work

Master Thesis

- Refine the tests and adding new ones, making them available and easy-usable for the community
- Create a **proactive pipeline** combining Fairness evaluation with Explainability, implementing a **local explainer for text**, related to Hate Speech, which will be able to identify the **potential biases** and the **categories** toward which the model is most discriminative, through the auditing of representative examples and counterexamples generated by CheckList framework



Conclusions

CheckList as a powerful
and easy usable tool

- CheckList **support the complete evaluation** phase of NLP models, for the task of interest and highlighting weaknesses that cannot be easily detected through real data
- The ideal aim: develop **standard evaluation processes** for NLP models and specific tasks
- Provides a way to **explore the models' dynamics**
 - through the analysis of the errors we can infer which model components are missing and which specific linguistic phenomena it has not yet acquired from the data



Personal note on the project

Hate-Speech Classifiers need a **robust value-sensitive evaluation**, to assess **unintended biases** and **avoid explicit harm** or the **amplification of pre-existing social biases**

- Dobbe et al. (2018): proposing a contribution within the Machine Learning domain responsibly foremost means **acknowledge our own biases** “*in an open and transparent way and engage in constructive dialogue with domain experts*”
 - with reference to the **choice of the datasets and the actual implementation of hand-coded templates**, that **shaped the work and the results**

Thanks! :)

Contacts: martamarchiori96@gmail.com

Icons from: <http://flaticon.com/>

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

Anzovino, Maria, et al. «Automatic Identification and Classification of Misogynistic Language on Twitter». Natural Language Processing and Information Systems, a cura di Max Silberztein 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.

References

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.

References

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.

References

EmoTag1200 👍 : Understanding the Association between Emojis 😊 and Emotions 🐱. Abu Awal Md Shoeb, 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.

References

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.