

Building a smart rejector for hate speech detection

Philippe Lammerts

4563182
24-05-2022

Supervised by:
Dr. J. Yang
Dr. Y-C. Hsu
P. Lippmann



Contents

- Context
- Problem
- Rejection
- Research questions
- Part 1: rejection metric
- Part 2: costs of predictions
- Part 3: unknown (un)knowns
- Planning

Context

- The amount of hateful content online is a growing concern
- It harms people and can even lead to acts of violence [1-5]
- Tackling hate speech is in the interest of all:
 - Governments
 - Social media companies
 - Everybody



[1] Balayn, A., Yang, J., Szlavik, Z., & Bozzon, A. (2021). Automatic Identification of Harmful, Aggressive, Abusive, and Offensive Language on the Web: A Survey of Technical Biases Informed by Psychology Literature. *ACM Transactions on Social Computing (TSC)*, 4(3), 1-56.

[2] Müller, K., & Schwarz, C. (2021). Fanning the flames of hate: Social media and hate crime. *Journal of the European Economic Association*, 19(4), 2131-2167.

[3] <https://www.nytimes.com/2022/02/08/world/asia/india-hate-speech-muslims.html>.

[4] <https://www.nytimes.com/2018/10/15/technology/myanmar-facebook-genocide.html>

[5] <https://www.coe.int/en/web/european-commission-against-racism-and-intolerance/hate-speech-and-violence>

Context

Two ways to detect hate speech [1]:

1. **Reactive moderation**
 - a. Flagging/reporting
2. **Proactive moderation**
 - a. Manually
 - b. Automatically

Problem

Manual proactive moderation

Human moderators

- + Most reliable
- Infeasible

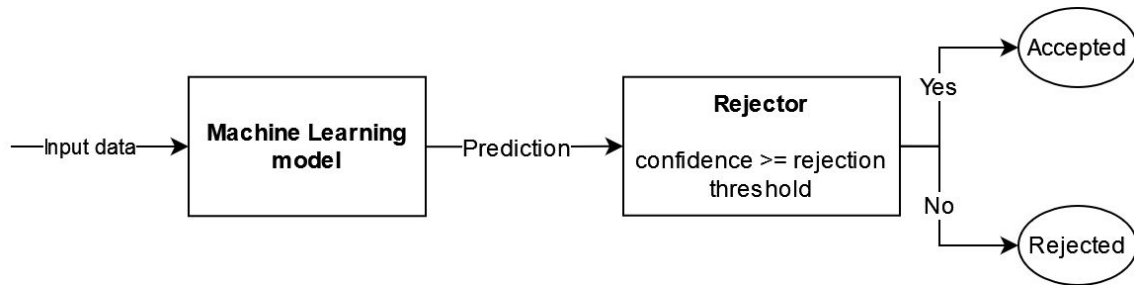
Automated proactive moderation

Machine Learning algorithms

- + Fast
- Can be unreliable [1, 2]:
 - Performs poor on deployment data
 - 69% F1-score drop when using different test datasets [2]

Rejection

“The goal of a machine learning model's reject option is to abstain from making a prediction when a model receives a test example where the risk of making a misprediction is too large.” [1]



Research questions

RQ: How can we maximize the utility of Machine Learning models in hate speech detection using a reject option?

SRQ1 How can we determine when the Machine Learning model is not confident enough?

- **SRQ1.1** How can we measure the utility of Machine Learning models with a reject option?
- **SRQ1.2** How can we determine the relative costs of rejections and True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) predictions?

SRQ2 Can unknown (un)known detection further improve the reject option?

Part 1: rejection metric

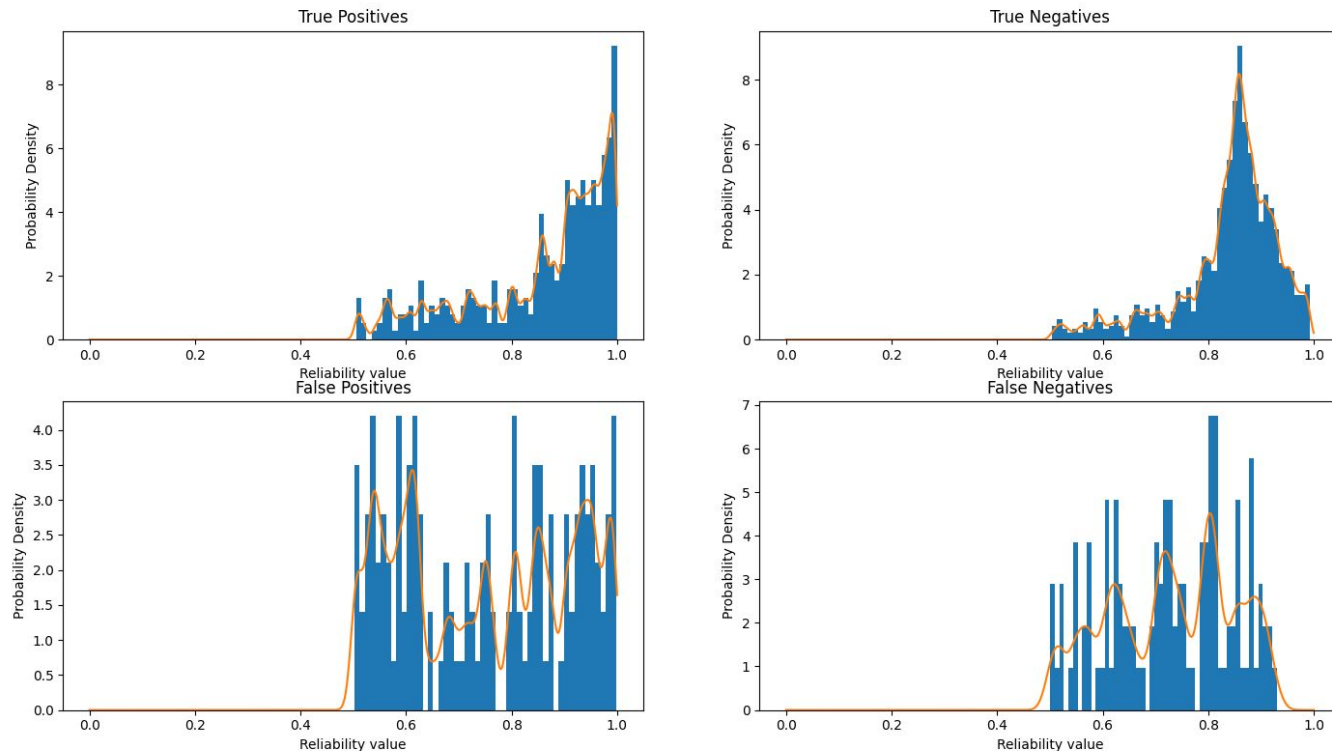
SRQ1.1 How can we measure the utility of Machine Learning models with a reject option?

- Metric that measures the effectiveness of the reject option
- Advanced version of the metric in [1]
- Calculation based on:
 - Gain of True Positive
 - Gain of True Negative
 - Cost of False Positive
 - Cost of False Negative
 - Cost of rejection
 - List of predictions with confidence values

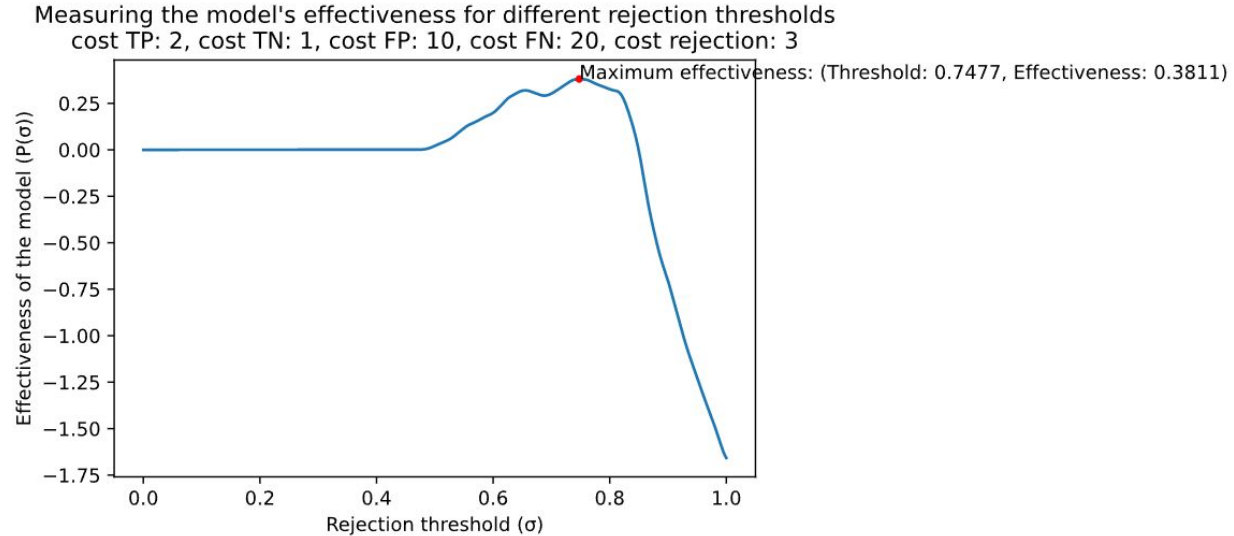
Part 1: rejection metric - example

Probability Density Functions for the sets of TP, TN, FP, and FN

The orange line is the estimated PDF that is derived using Kernel Density Estimation by fitting it with the original data. The blue histogram is the probability density of the original data



Part 1: rejection metric - example



Part 2: costs of predictions

SRQ1.2 How can we determine the relative costs of rejections and True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) predictions?

- Objective cost analysis in hate speech is infeasible [1]
- We focus on subjective cost analysis
- Magnitude Estimation scale
 - Unbounded rating scale
 - Provides ratio data
 - Used in many different domains [2-6]

[1] Sunstein, C. R. (2018). Does the Clear and Present Danger Test Survive Cost-Benefit Analysis?. *Cornell L. Rev.*, 104, 1775.

[2] Maddalena, E., Mizzaro, S., Scholer, F., & Turpin, A. (2017). On crowdsourcing relevance magnitudes for information retrieval evaluation. *ACM Transactions on Information Systems (TOIS)*, 35(3), 1-32.

[3] Lodge, M., & Tursky, B. (1979). Comparisons between category and magnitude scaling of political opinion employing SRC/CPS items. *American Political Science Review*, 73(1), 50-66.

[4] Lodge, M., Tanenhaus, J., Cross, D., Tursky, B., Foley, M. A., & Foley, H. (1976). The calibration and cross-modal validation of ratio scales of political opinion in survey research. *Social Science Research*, 5(4), 325-347.

[5] McGee, M. (2004, April). Master usability scaling: magnitude estimation and master scaling applied to usability measurement. In *Proceedings of the SIGCHI conference on Human factors in computing systems* (pp. 335-342).

[6] Bard, E. G., Robertson, D., & Sorace, A. (1996). Magnitude estimation of linguistic acceptability. *Language*, 32-68.

Part 2: costs of predictions

Experiment

Present TP, TN, FP, FN, and rejection scenarios to a group of subjects:

- Show (non)hateful tweet
- Show decision of platform
- Ask subjects whether they (dis)agree with the decision using the ME scale

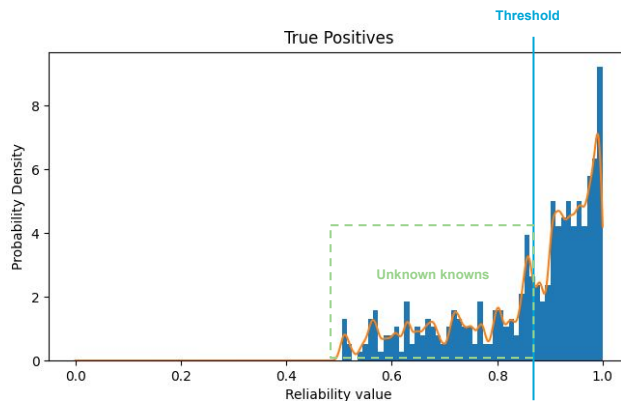
Progress

- Pre-registration report is ready
 - Experimental setup
 - Procedure
 - Analysis
- Waiting for approval of Human Research Ethics

Part 3: unknown (un)knowns

SRQ2 Can unknown (un)known detection further improve the reject option?

- A single optimal rejection threshold is not enough
- Different methods for unknown (un)known detection [1-3]



Planning

- Currently finishing part 2
- Combine findings of parts 1 + 2 and create a paper submission for:
 - *The 10th AAAI Conference on Human Computation and Crowdsourcing (HCOMP 2022)*
 - Deadline June 24, 2022
- Work on part 3 after the deadline
- Finish around the end of Q1 2022-2023