# Machine Learning Model for Classifying Tweet's Content

A model which classifies the content of text (tweets specifically) in 3 classes:

Neutral Speech, Offensive Speech and Hate Speech; using a dataset with

Created by: Kalovan Dragiychev
already annotated tweets for training and testing purposes, 2024 10:44 AM

Changed on: April 25, 2024 11:09 AM

Context of use: Education Level of education: Bachelor

Machine Learning Model for Classifying Tweet's Content

#### Impact on society

What impact is expected from your technology?

This category is only partial filled.

#### What is exactly the problem? Is it really a problem? Are you sure?

The problem this technology addresses is the prevalence of hate speech and offensive language on social media, which undermines safety and inclusivity online. This "pain" is felt by social media users exposed to harmful content, causing emotional distress and discouraging engagement. Platform administrators also struggle to moderate content effectively without infringing on free speech. This technology aims to ease these challenges by enhancing the ability to detect and manage such content efficiently. Solving this problem improves digital communication safety and fosters a respectful, inclusive online community, making it a worthwhile endeavor.

Are you sure that this technology is solving the RIGHT problem? This question has not been answered yet.

How is this technology going to solve the problem? This question has not been answered yet.

What negative effects do you expect from this technology? This question has not been answered yet.

In what way is this technology contributing to a world you want to live in?

This question has not been answered yet.

Now that you have thought hard about the impact of this technology on society (by filling out the questions above), what improvements would you like to make to the technology? List them below. This question has not been answered yet.

Machine Learning Model for Classifying Tweet's Content

#### Hateful and criminal actors

What can bad actors do with your technology?

This category is only partial filled.

In which way can the technology be used to break the law or avoid the consequences of breaking the law?

Hate speech detection technology, while designed to improve social media safety, could potentially be misused in ways that infringe on privacy or freedom of speech. For instance, governments or other entities might exploit the technology to surveil and suppress dissenting voices under the guise of moderating hate speech. This could lead to censorship of political opposition or marginalized groups, impacting their right to express dissent. Additionally, the technology could be manipulated to wrongfully flag content, causing harassment or reputational damage to individuals by falsely labeling their speech as offensive or hateful. Moreover, there is a risk that the misuse of this technology could lead to biased enforcement of rules, disproportionately impacting certain groups based on flawed data or biased algorithms.

Can fakers, thieves or scammers abuse the technology? This question has not been answered yet.

Can the technology be used against certain (ethnic) groups or (social) classes?

This question has not been answered yet.

In which way can bad actors use this technology to pit certain groups against each other? These groups can be, but are not constrained to, ethnic, social, political or religious groups.

This question has not been answered yet.

How could bad actors use this technology to subvert or attack the truth?

This question has not been answered yet.

Now that you have thought hard about how bad actors can impact this technology, what improvements would you like to make? List them below.

Machine Learning Model for Classifying Tweet's Content

#### **Privacy**

Are you considering the privacy & personal data of the users of your technology?

This category is only partial filled.

## Does the technology register personal data? If yes, what personal data?

The hate speech detection technology inherently processes personal data as it analyzes content generated by individuals on social media platforms. This includes text from tweets, which may contain identifiable information about the user, potentially revealing their identity, location, or personal beliefs. Although the primary intent is to assess the content for hate speech or offensive language, the nature of social media data means that personal and sometimes sensitive information could be indirectly collected and analyzed.

To adhere to privacy laws such as the General Data Protection Regulation (GDPR), any handling of personal data by this technology must ensure privacy protections and data minimization principles. Special attention is needed to ensure that the data does not unintentionally capture or store sensitive categories of personal data, such as racial or ethnic origin, political opinions, or health information, without explicit consent and legal justification.

Do you think the technology invades the privacy of the stakeholders? If yes, in what way?

This question has not been answered yet.

Is the technology is compliant with prevailing privacy and data protection law? Can you indicate why? This question has not been answered yet.

Does the technology mitigate privacy and data protection risks/concerns (privacy by design)? Please indicate how. This question has not been answered yet.

In which way can you imagine a future impact of the collection of personal data?

This question has not been answered yet.

Now that you have thought hard about privacy and data protection,

Machine Learning Model for Classifying Tweet's Content

what improvements would you like to make? List them below. This question has not been answered yet.

Machine Learning Model for Classifying Tweet's Content

#### **Human values**

How does the technology affect your human values?

This category is only partial filled.

How is the identity of the (intended) users affected by the technology? The hate speech detection technology impacts user identity by potentially altering how individuals express themselves online. For friends using social media, it can enhance relationships by creating a safer communication space free from hate speech, which may encourage more open interaction. However, if the technology over-moderates, it could detract from authentic expression, leading to frustrations or misunderstanding about why certain content is censored.

The technology fills a role previously occupied by human moderators, automating content moderation which might lead to questions about the imposition of specific world views or biases. It aims to be neutral but could inadvertently enforce certain societal norms, which may be seen as imposing a belief system.

While the technology is designed to empower users by protecting them from harmful content, thereby maintaining their dignity, it could also change how they communicate by making them more mindful or cautious about the language they use online. This reflection on communication might lead to changes in social norms and personal behavior on digital platforms, aligning with a collective desire for more respectful and inclusive online communities.

How does the technology influence the users' autonomy? This question has not been answered yet.

What is the effect of the technology on the health and/or well-being of users?

This question has not been answered yet.

Now that you have thought hard about the impact of your technology on human values, what improvements would you like to make to the technology? List them below.

Machine Learning Model for Classifying Tweet's Content

#### **Stakeholders**

Have you considered all stakeholders?

This category is only partial filled.

Who are the main users/targetgroups/stakeholders for this technology? Think about the intended context by answering these questions.

Name of the stakeholder Social Media Platforms

How is this stakeholder affected?

-

Did you consult the stakeholder?

Are you going to take this stakeholder into account? No

Name of the stakeholder Blogs, Forums

How is this stakeholder affected?

**Did you consult the stakeholder?** No

**Are you going to take this stakeholder into account?** No

Did you consider all stakeholders, even the ones that might not be a user or target group, but still might be of interest?

Now that you have thought hard about all stakeholders, what improvements would you like to make? List them below. This question has not been answered yet.

https://www.tict.io

Machine Learning Model for Classifying Tweet's Content

#### Data

Is data in your technology properly used?

This category is only partial filled.

Are you familiar with the fundamental shortcomings and pitfalls of data and do you take this sufficiently into account in the technology? Yes, the technology is designed with an awareness of data limitations such as subjectivity, incompleteness, bias, and the complexity of language. We recognize that the dataset may not fully represent all nuances of language and cultural expressions, which can lead to inaccuracies in hate speech detection. To mitigate these issues, the model incorporates diverse data sources to balance representation, employs techniques to adjust for class imbalance, and includes regular updates and feedback mechanisms to refine its understanding and accuracy. Users are also informed of the model's potential limitations, ensuring they understand the contextual capabilities of the technology.

How does the technology organize continuous improvement when it comes to the use of data?

This question has not been answered yet.

How will the technology keep the insights that it identifies with data sustainable over time?

This question has not been answered yet.

In what way do you consider the fact that data is collected from the users?

This question has not been answered yet.

Now that you have thought hard about the impact of data on this technology, what improvements would you like to make? List them below.

Machine Learning Model for Classifying Tweet's Content

**Inclusivity** 

Is your technology fair for everyone?

This category is only partial filled.

Will everyone have access to the technology?

This question has not been answered yet.

Does this technology have a built-in bias?

Yes, the technology could potentially have built-in biases due to the nature of its data sources and the design of its algorithms. The dataset, primarily composed of Twitter data, may inherently reflect the biases of the users who generated the tweets, including geographical, linguistic, or cultural biases. Additionally, the annotators who labeled the dataset might have introduced subjective interpretations of what constitutes hate speech or offensive language. Recognizing these limitations, the design includes regular reviews and updates to the model to address and correct these biases. Moreover, efforts are made to include a diverse group of developers and feedback from a wide range of users to continuously evaluate and improve the algorithms fairness and inclusivity.

Does this technology make automatic decisions and how do you account for them?

This question has not been answered yet.

Is everyone benefitting from the technology or only a a small group? Do you see this as a problem? Why/why not?

This question has not been answered yet.

Does the team that creates the technology represent the diversity of our society?

This question has not been answered yet.

Now that you have thought hard about the inclusivity of the technology, what improvements would you like to make? List them below.

Machine Learning Model for Classifying Tweet's Content

#### **Transparency**

Are you transparent about how your technology works?

This category is only partial filled.

# Is it explained to the users/stakeholders how the technology works and how the business model works?

Yes, the operation and goals of the hate speech detection technology are broadly explained to users and stakeholders. We detail the types of data used, the purpose of the technology, and the methods employed in detecting hate speech and offensive language on our website and in user documentation. However, while we strive for transparency, the specific inner workings of the machine learning algorithmslike how they process data and arrive at specific determinations are not fully disclosed for two reasons: first, the complexity of the algorithms may not be easily understandable to all users; second, revealing detailed mechanics might compromise the effectiveness of the moderation tools. This balance aims to maintain user trust while protecting the systems integrity.

If the technology makes an (algorithmic) decision, is it explained to the users/stakeholders how the decision was reached? This question has not been answered yet.

Is it possible to file a complaint or ask questions/get answers about this technology?

This question has not been answered yet.

Is the technology (company) clear about possible negative consequences or shortcomings of the technology? This question has not been answered yet.

Now that you have thought hard about the transparency of this technology, what improvements would you like to make? List them below.

Machine Learning Model for Classifying Tweet's Content

#### Sustainability

Is your technology environmentally sustainable?

This category is only partial filled.

# In what way is the direct and indirect energy use of this technology taken into account?

The hate speech detection technology primarily operates in cloud environments, which do consume significant amounts of energy. However, we partner with cloud service providers committed to sustainability, utilizing energy-efficient data centers often powered by renewable energy sources. While the model is computationally intensive, efforts are made to optimize algorithms for better efficiency, reducing unnecessary computations and data transfers. We continuously explore advancements in technology that could further reduce the energy footprint, such as improved algorithm efficiency and deploying edge computing where feasible to lessen data center reliance.

# Do you think alternative materials could have been considered in the technology?

This question has not been answered yet.

Do you think the lifespan of the technology is realistic? This question has not been answered yet.

What is the hidden impact of the technology in the whole chain? This question has not been answered yet.

Now that you have thought hard about the sustainability of this technology, what improvements would you like to make? List them below.

Machine Learning Model for Classifying Tweet's Content

#### **Future**

Did you consider future impact?

This category is only partial filled.

What could possibly happen with this technology in the future? With widespread adoption, this hate speech detection technology could significantly alter online communication norms, encouraging more respectful interactions and possibly reducing the prevalence of harmful content on social media. As communities adapt to these moderated environments, there might be a shift towards more positive online behavior. However, there's also potential for overreliance on automated moderation, possibly leading to suppression of legitimate speech or failure to capture nuanced expressions. Additionally, if not managed correctly, the technology could be used to enforce biased norms or censorship, affecting freedom of expression. Hence, while it has the potential to make online spaces safer, careful management and continuous oversight are crucial to prevent misuse and unintended negative consequences.

Sketch a or some future scenario (s) (20-50 years up front) regarding the technology with the help of storytelling. Start with at least one utopian scenario.

This question has not been answered yet.

Sketch a or some future scenario (s) (20-50 years up front) regarding the technology with the help of storytelling. Start with at least one dystopian scenario.

This question has not been answered yet.

Would you like to live in one of this scenario's? Why? Why not? This question has not been answered yet.

What happens if the technology (which you have thought of as ethically well-considered) is bought or taken over by another party? This question has not been answered yet.

Impact Improvement: Now that you have thought hard about the future impact of the technology, what improvements would you like to make? List them below.