# BiasShield - Toxicity & Bias Detection Report
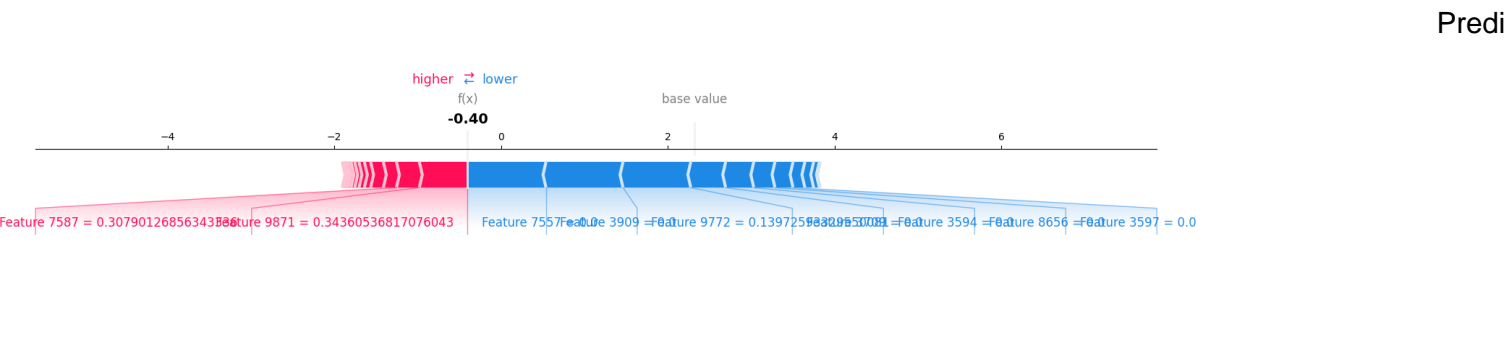
## Prompt 1:

I hate you, go to hell.
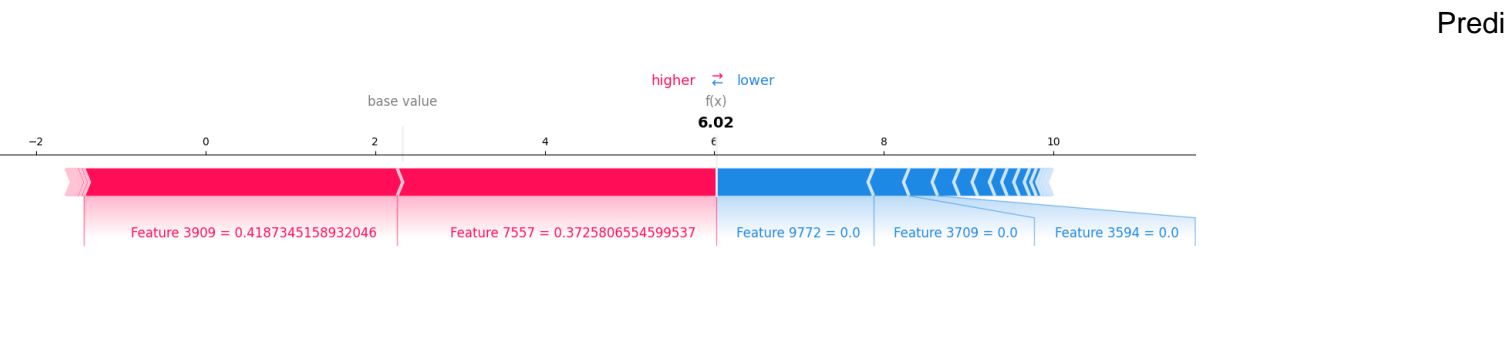
Predi

higher ⇄ lower
base value            f(x)
                     4.48

−10    −8    −6    −4    −2    0    2    4    6    8    10

Feature 3438 = 0.233625...  = 0.462143...  = 0.460368...  = 0.349263...  = 0.34520929018063556    Feature 9772 = 0.1022...    Feature 3909 = 0.0

## Prompt 2:

You're such a kind and wonderful person.

Predi

higher ⇄ lower
f(x)              base value
-0.40

−4              −2              0              2              4              6

Feature 7587 = 0.30790126856343...  Feature 9871 = 0.34360536817076043    Feature 7557 = 0.0  Feature 3909 = 0.0  Feature 9772 = 0.13972...    Feature 3594 = 0.0  Feature 8656 = 0.0  Feature 3597 = 0.0

## Prompt 3:

Stupid idiot, learn to speak properly!

Predi

higher ⇄ lower
base value            f(x)
                     6.02

−2         0         2         4         6         8         10

Feature 3909 = 0.4187345158932046    Feature 7557 = 0.3725806554599537    Feature 9772 = 0.0    Feature 3709 = 0.0    Feature 3594 = 0.0

## Prompt 4:

Let's have a respectful conversation.

higher ⇄ lower
f(x)
**-3.19**

−4    −3    −2    −1    0    1    2    3    4    5
                                        base value

Feature 9772 = 0.0    Feature 7557 = 0.0    Feature 3909 = 0.0    Feature 3709 = 0.0    Feature 3594 = 0.0    Feature 8656 = 0.0    Feature 3597 = 0.0

## Prompt 5:

You filthy scumbag!

base value

higher ⇄ lower
f(x)
**4.72**

0    2    4    6    8    10    12

Feature 9772 = 1.0    Feature 7557 = 0    Feature 3909 = 0    Feature 3709 = 0    Feature 3594 = 0    Feature 8656 = 0    Feature 3597 = 0.0

# Exploratory Visualizations

higher ⇄ lower

f(x)      base value
0.26

−3      −2      −1      0      1      2      3      4      5      6

wanna = 0.6402836240270403    rape = 0.6471084313405574    fuck = 0.0    hi = 0.3848309609547878    anal = 0.0    crap = 0.0

## Comment Length by Toxicity



Comment Length by Toxicity
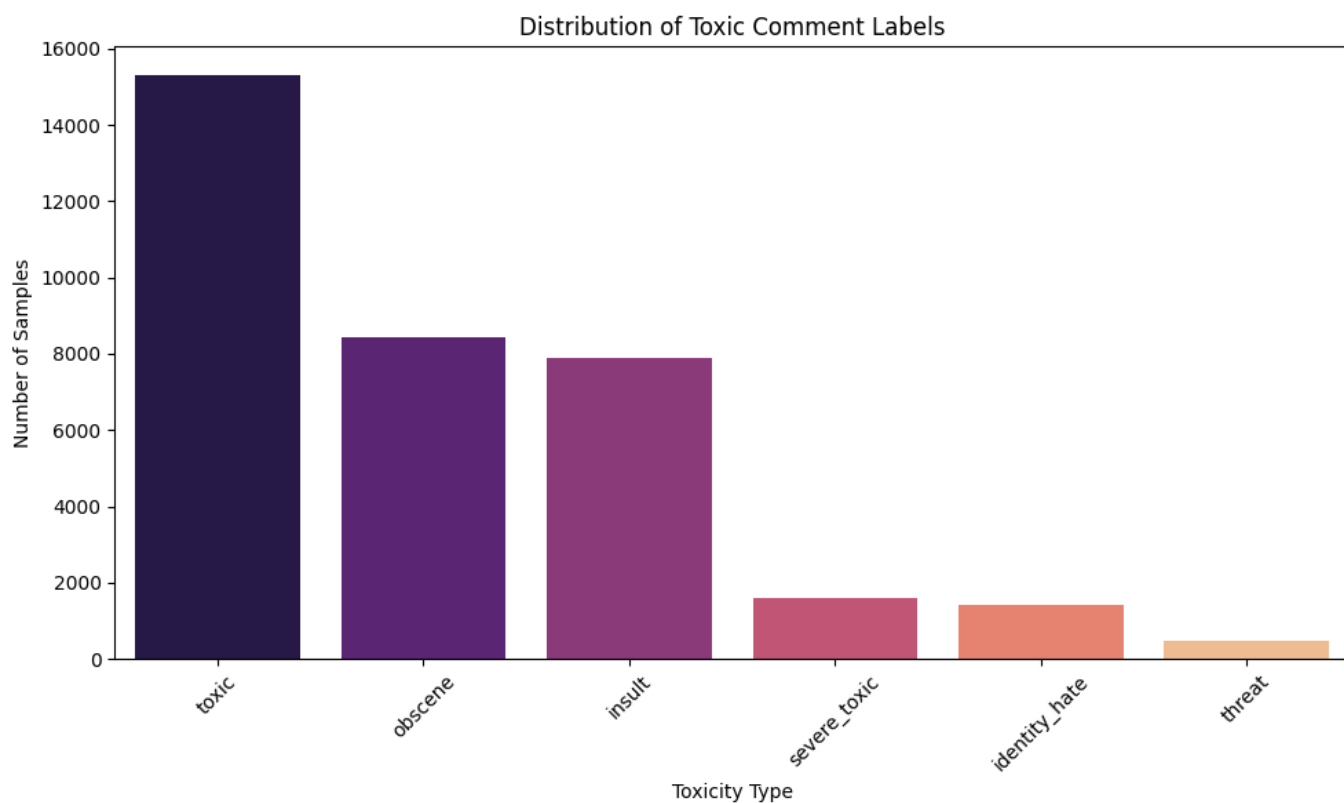
## Distribution of Toxicity Labels



## Top Identity Terms in Identity Hate Comments

Confusion Matrix - Logistic Regression (Toxicity)



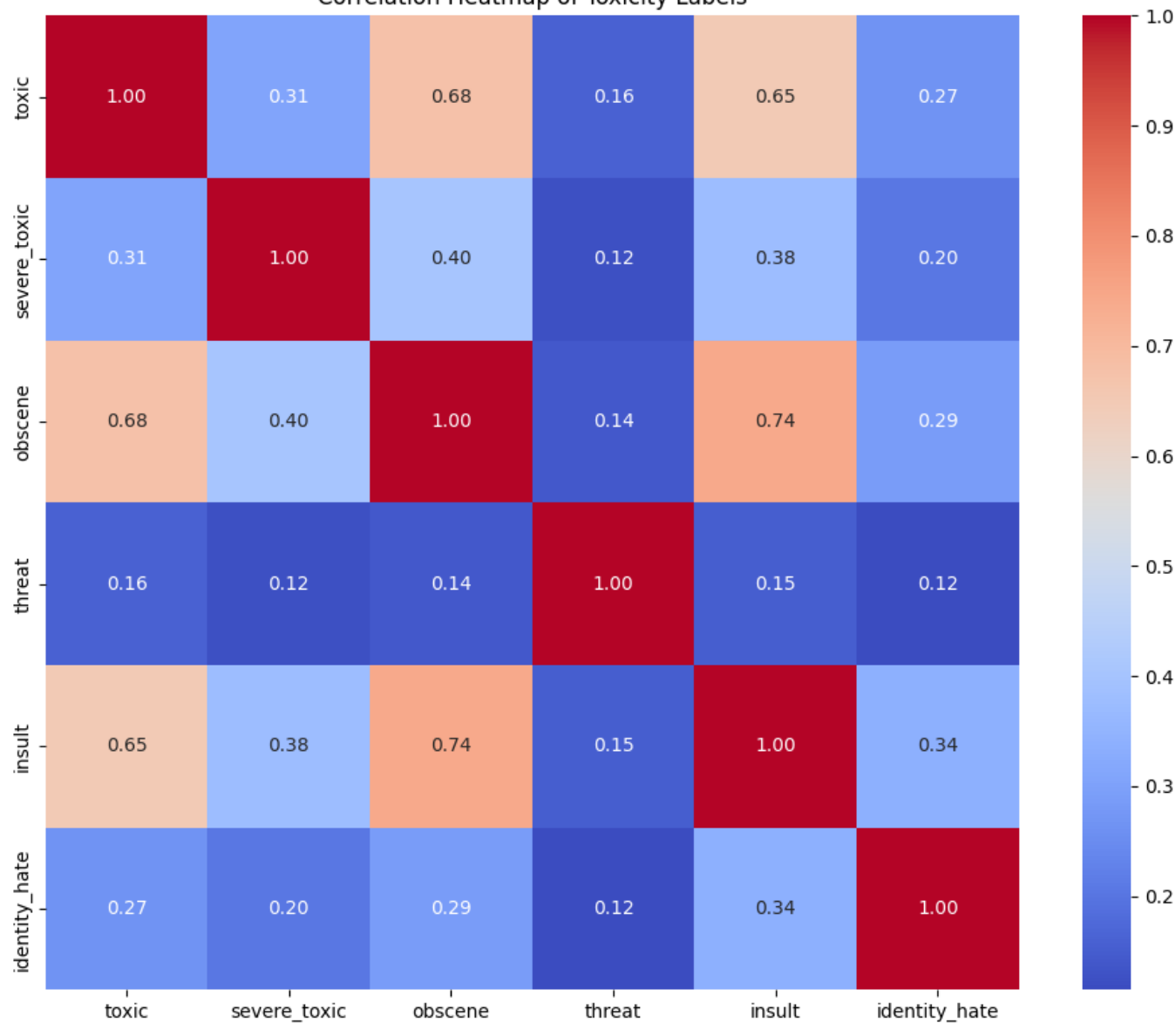Distribution of Toxic Comment Labels

## WordCloud - Toxic Comments



## WordCloud - Non-Toxic Comments



## Pairplot of Toxicity Labels

Correlation Heatmap of Toxicity Labels

|              | toxic | severe_toxic | obscene | threat | insult | identity_hate |
|--------------|-------|--------------|---------|--------|--------|---------------|
| toxic        | 1.00  | 0.31         | 0.68    | 0.16   | 0.65   | 0.27          |
| severe_toxic | 0.31  | 1.00         | 0.40    | 0.12   | 0.38   | 0.20          |
| obscene      | 0.68  | 0.40         | 1.00    | 0.14   | 0.74   | 0.29          |
| threat       | 0.16  | 0.12         | 0.14    | 1.00   | 0.15   | 0.12          |
| insult       | 0.65  | 0.38         | 0.74    | 0.15   | 1.00   | 0.34          |
| identity_hate| 0.27  | 0.20         | 0.29    | 0.12   | 0.34   | 1.00          |

Word Count Distribution of Comments