

Toxic Comment Classification

****Disclaimer****

The Presentation Slides contain text
that may be considered profane, vulgar, or offensive

서대원, 김동욱, 김호현, 이민규, 임희진



CONTENTS

1 Data Description

- Overview
- Example
- Y_Visualization
- X_Visualization

2 Preprocessing

- Tokenization
- TF-IDF

3 Model

- Model1
- Model2
- Model3
- Model4

4 Result

- Lime
Package

CONTENTS



1 Data Description

- Overview
- Example
- Y_Visualization
- X_Visualization

2 Preprocessing

- Tokenization
- TF-IDF

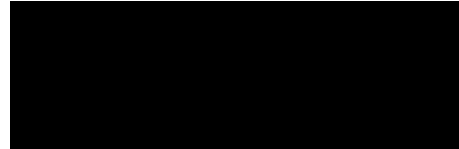
3 Model

- Model1
- Model2
- Model3
- Model4

4 Result

- Lime
Package

OVERVIEW



powered by Google

Toxic Comment Classification Challenge

Purpose

To make health environment for the Internet Users

Data Source

Wikipedia corpus dataset
collected through Google Human-labeling

Approximately 63M comments from web pages/
article from 2004 to 2015

OVERVIEW : Multiclass-Multilabel Classification

X_label

```
x.head(5)
```

	comment_text
0	ExplanationWhy the edits made under my usern...
1	D'aww! He matches this background colour I'm s...
2	Hey man, I'm really not trying to edit war. It...
3	"MoreI can't make any real suggestions on ...
4	You, sir, are my hero. Any chance you remember...

Y_label

```
y.head(5)
```

	toxic	severe_toxic	obscene	threat	insult	identity_hate
0	0	0	0	0	0	0
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0

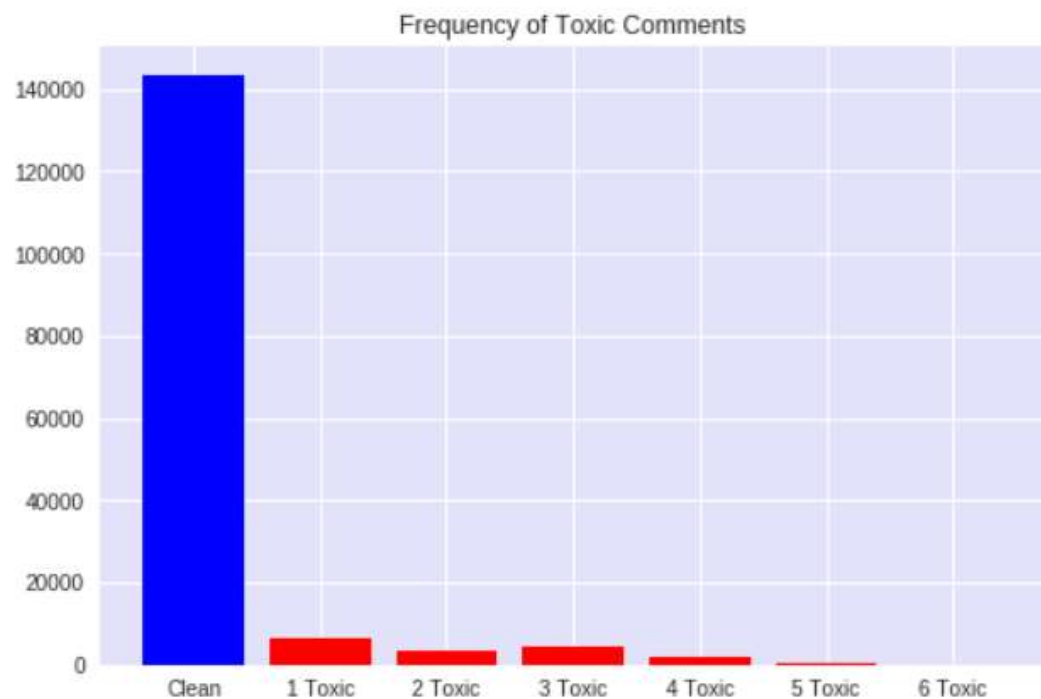
SAMPLES : Stereotypical Sample

comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
████████ JCKER BEFORE YOU ████████ AROUND ON MY WORK	1	1	1	0	1	0

id	8b20912530eebd56
comment_text	"WnSURPRISE!Wn / " "WWn ...
toxic	0
severe_toxic	0
obscene	0
threat	0
insult	0
identity_hate	0
Name: 152298, dtype: object	

Visualization : The vast majority of comments are Non-toxic

The Proportion of Clean comments > 95%



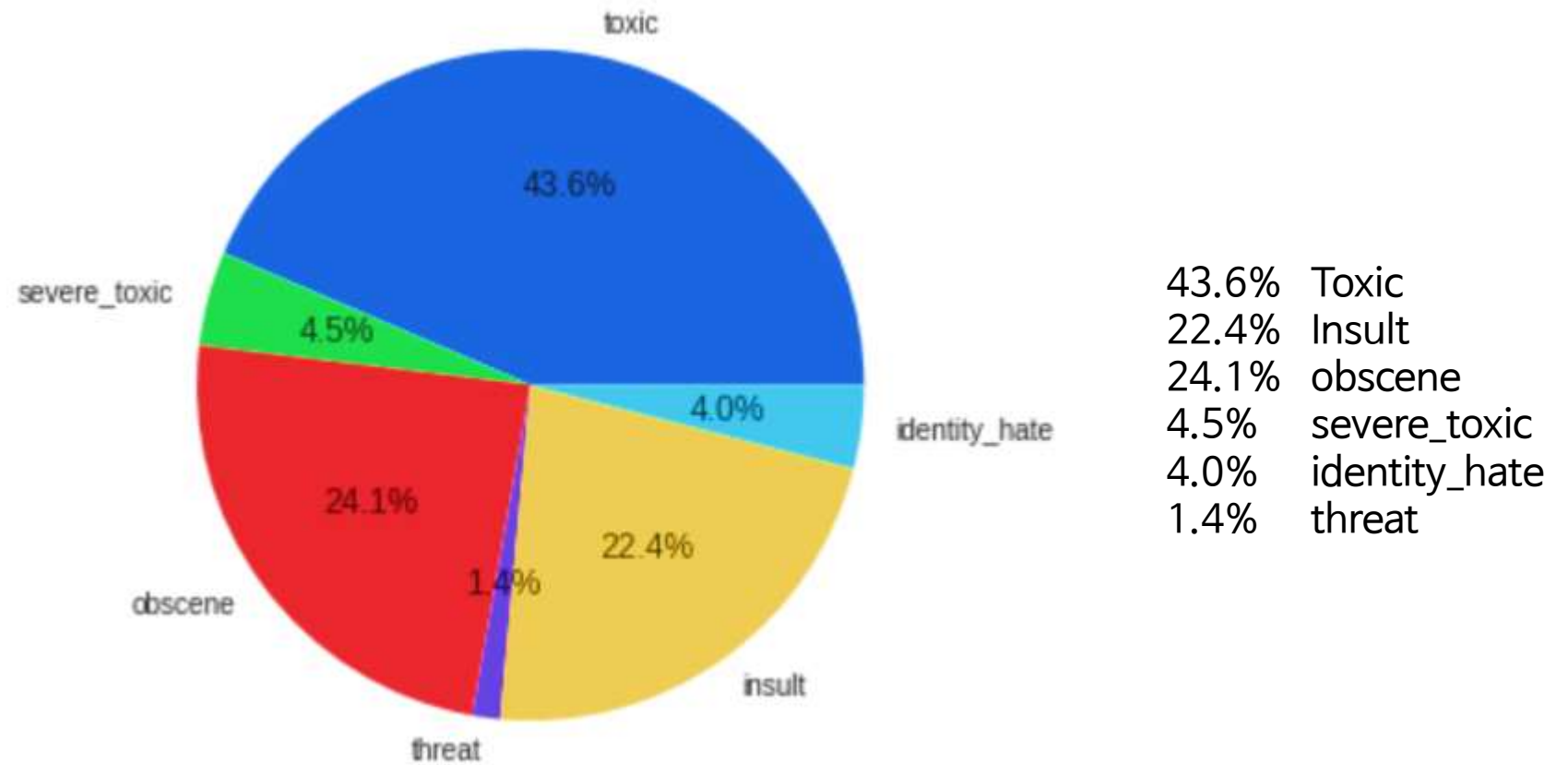
```
import numpy as np
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score

print("bare minimum acc_score:%0.3f"
      %(accuracy_score(y, np.zeros(y.shape))))
print('bare minimum roc_auc_score:%0.3f'
      %(roc_auc_score(y, np.random.rand(159571,6))))
```

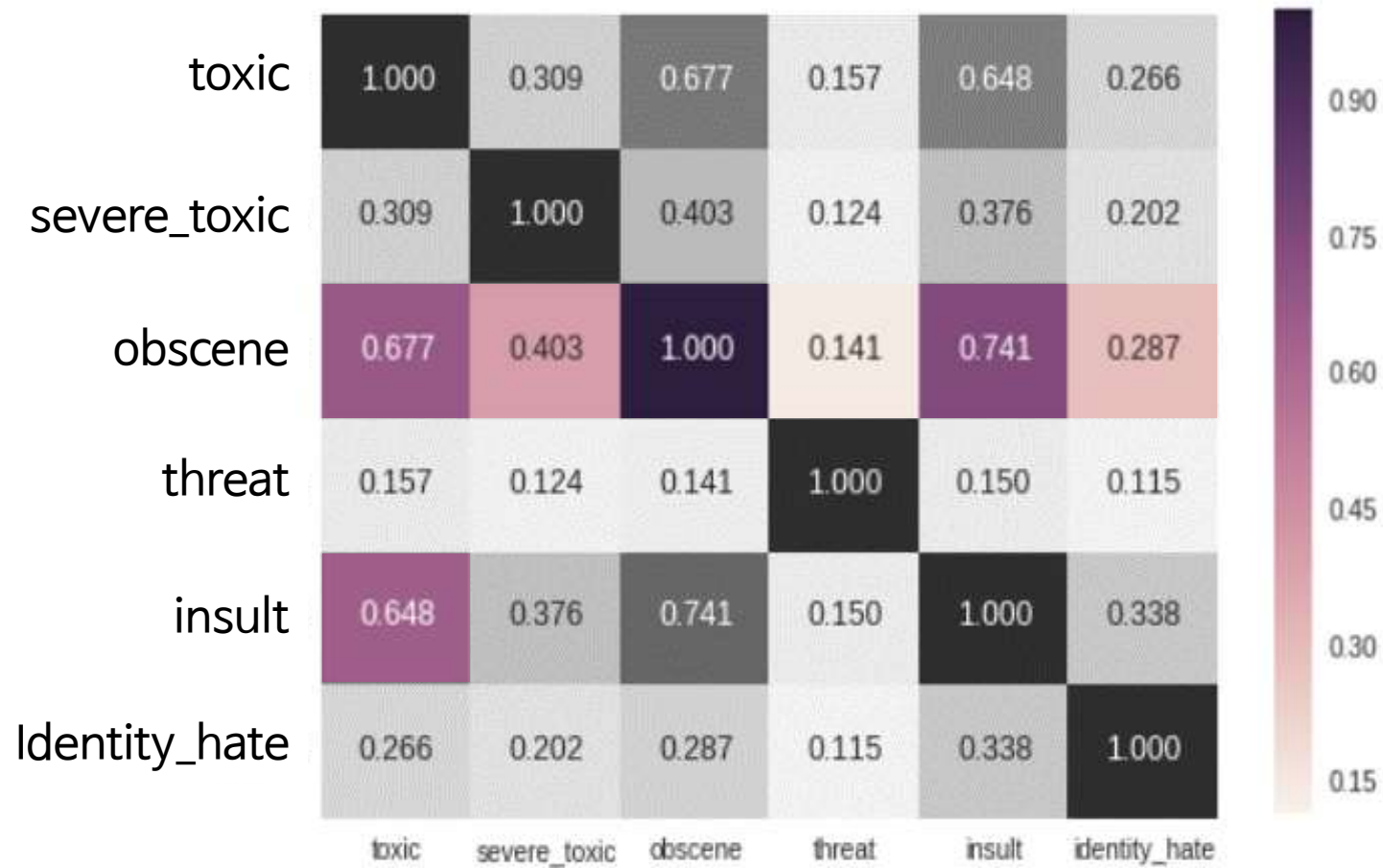
bare minimum acc_score:0.898

bare minimum roc_auc_score:0.503

Visualization : the Proportion of Each Label

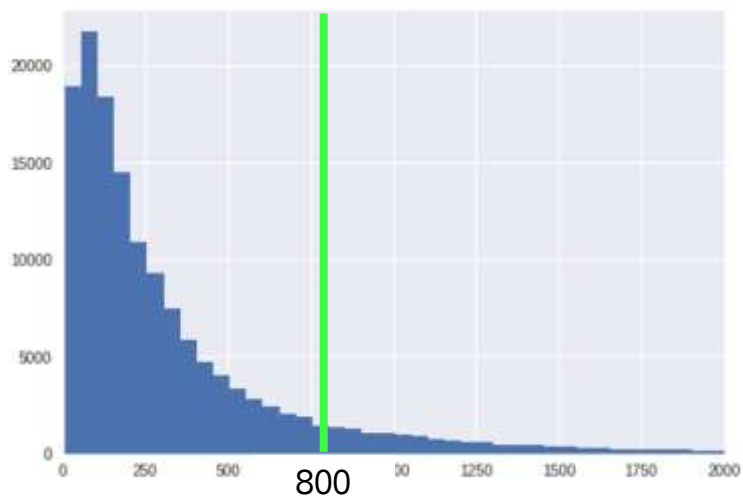


Visualization : the Correlation between Y_label



Visualization : the Length of Text

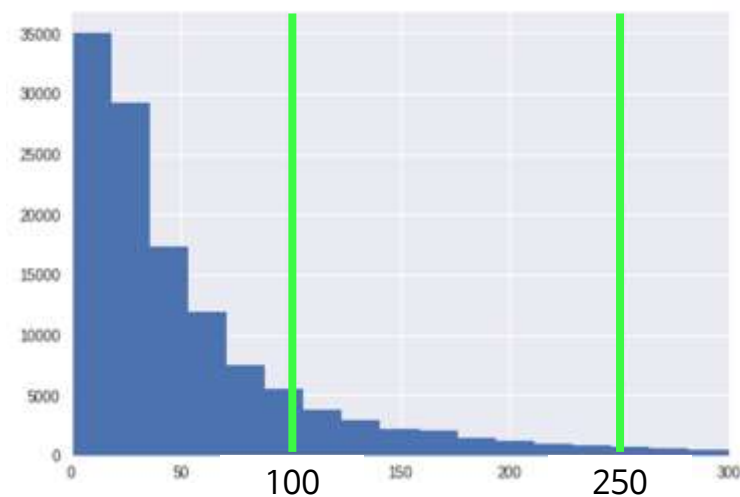
WORD_wise



Tokenization max length = 800

김호현

CHARCATER_wise



Tokenization max length = 250

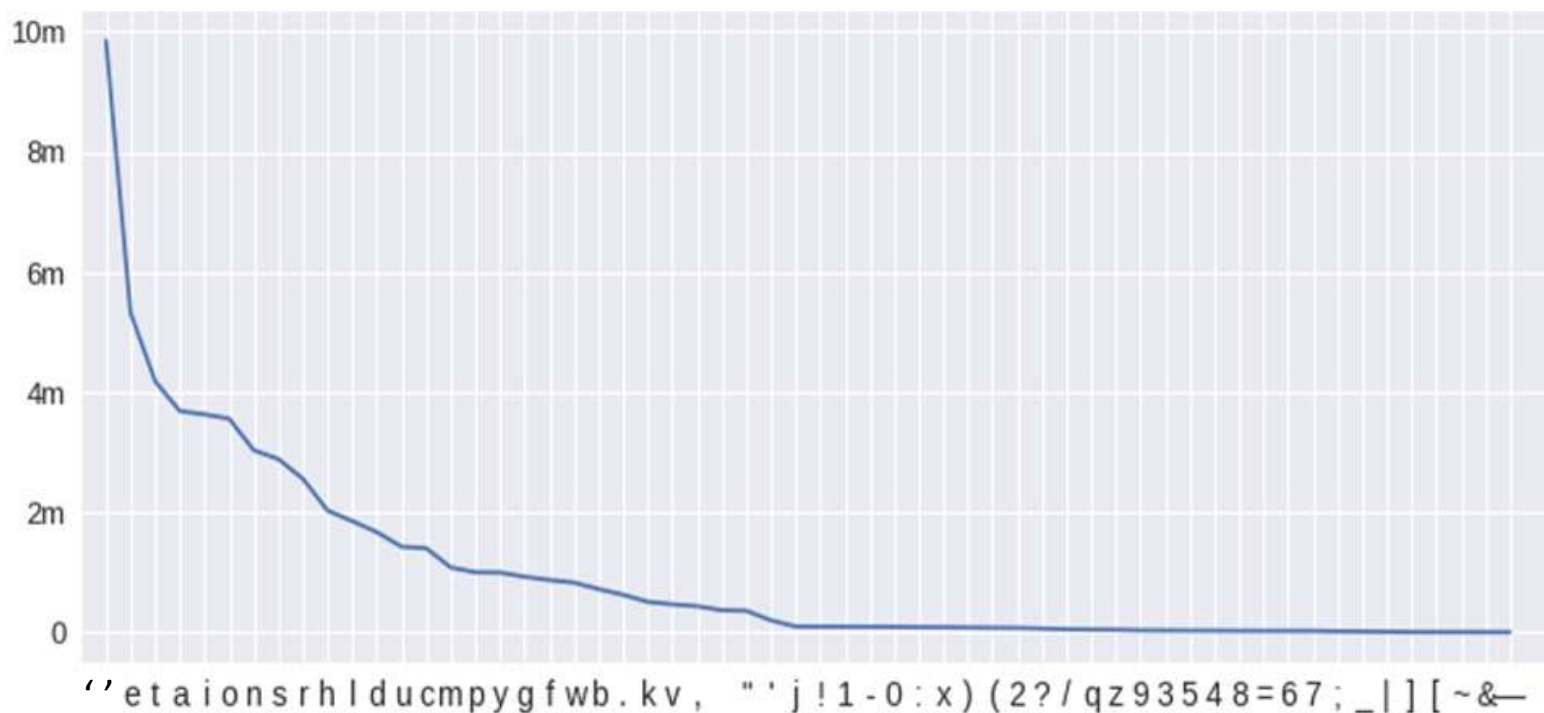
김호현

Tokenization max length = 100

서대원

Visualization : the Frequency of Character

Many Special Characters are Placed in High Rank



총 문자 수 : 1934

문자 출현 빈도 순위

‘‘ : 9850190

‘e’ : 5320996

‘t’ : 4179970

‘a’ : 3685571

Visualization : Special Characters

[illegible]

Visualization : Word Cloud

Clean



Toxic



Insult



Threat



Obscene



Severe Toxic



Identity_hate





1 Data Description

- Overview
- Example
- Y_Visualization
- X_Visualization

2 Preprocessing

- Tokenization
- TF-IDF

3 Model

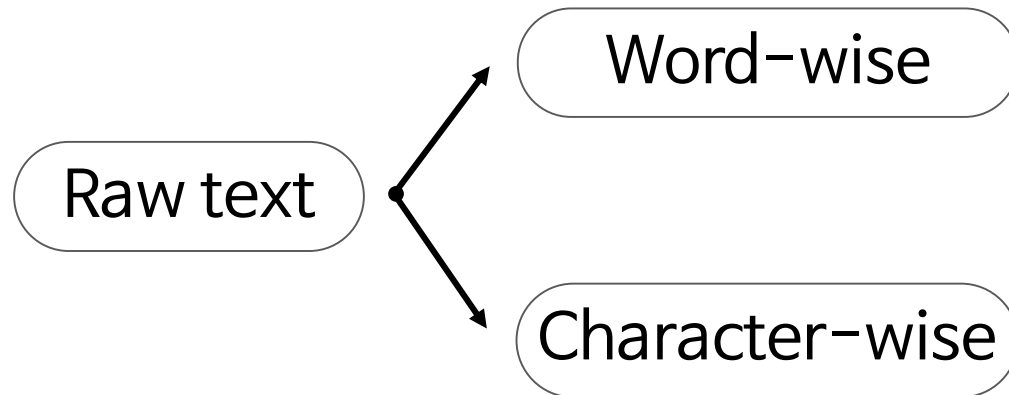
- Model1
- Model2
- Model3
- Model4

4 Result

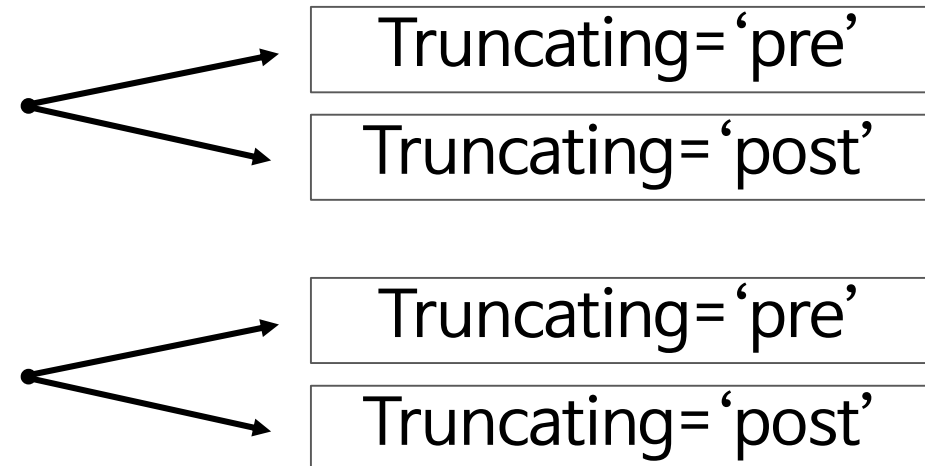
- Lime
Package

Tokenization

Step1. Tokenization



Step2. Padding



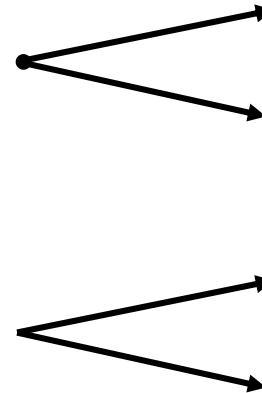
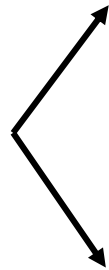
Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text



Tokenization

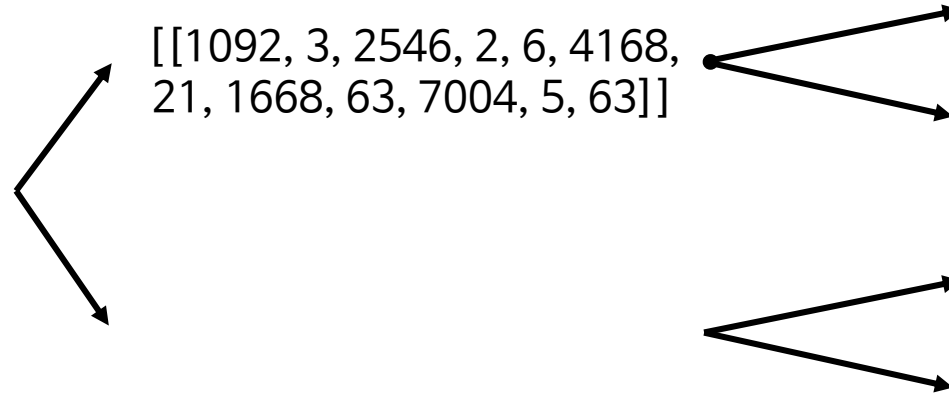
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]



Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14,
6, 15, 16, 4, 9, 2, 1, 3, 10, 2, 2,
1, 3, 6, 1, 4, 1, 8, 13, 15, 15, 2,
9, 28, 8, 1, 12, 4, 17, 39, 1, 3,
10, 6, 13, 1, 4, 9, 3, 1, 15, 6, 9,
2, 1, 11, 6, 24, 2, 11, 17, 1, 4,
7, 12, 1, 15, 6, 9, 2, 1, 3, 2, 15,
16, 2, 9, 4, 3, 2]]

Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]

[[2546 2 6 4168 21
1668 63 7004 5 63]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14,
6, 15, 16, 4, 9, 2, 1, 3, 10, 2, 2,
1, 3, 6, 1, 4, 1, 8, 13, 15, 15, 2,
9, 28, 8, 1, 12, 4, 17, 39, 1, 3,
10, 6, 13, 1, 4, 9, 3, 1, 15, 6, 9,
2, 1, 11, 6, 24, 2, 11, 17, 1, 4,
7, 12, 1, 15, 6, 9, 2, 1, 3, 2, 15,
16, 2, 9, 4, 3, 2]]

Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14,
6, 15, 16, 4, 9, 2, 1, 3, 10, 2, 2,
1, 3, 6, 1, 4, 1, 8, 13, 15, 15, 2,
9, 28, 8, 1, 12, 4, 17, 39, 1, 3,
10, 6, 13, 1, 4, 9, 3, 1, 15, 6, 9,
2, 1, 11, 6, 24, 2, 11, 17, 1, 4,
7, 12, 1, 15, 6, 9, 2, 1, 3, 2, 15,
16, 2, 9, 4, 3, 2]]

[[2546 2 6 4168 21
1668 63 7004 5 63]]

[[1092 3 2546 2 6 4168
21 1668 63 7004]]

Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14,
6, 15, 16, 4, 9, 2, 1, 3, 10, 2, 2,
1, 3, 6, 1, 4, 1, 8, 13, 15, 15, 2,
9, 28, 8, 1, 12, 4, 17, 39, 1, 3,
10, 6, 13, 1, 4, 9, 3, 1, 15, 6, 9,
2, 1, 11, 6, 24, 2, 11, 17, 1, 4,
7, 12, 1, 15, 6, 9, 2, 1, 3, 2, 15,
16, 2, 9, 4, 3, 2]]

[[2546 2 6 4168 21
1668 63 7004 5 63]]

[[1092 3 2546 2 6 4168
21 1668 63 7004]]

[[1 3 2 15 16 2 9 4 3 2]]

Tokenization

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate.

Step1. Tokenization

Step2. Padding

Raw text

[[1092, 3, 2546, 2, 6, 4168,
21, 1668, 63, 7004, 5, 63]]

[[2546 2 6 4168 21
1668 63 7004 5 63]]

[[1092 3 2546 2 6 4168
21 1668 63 7004]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14,
6, 15, 16, 4, 9, 2, 1, 3, 10, 2, 2,
1, 3, 6, 1, 4, 1, 8, 13, 15, 15, 2,
9, 28, 8, 1, 12, 4, 17, 39, 1, 3,
10, 6, 13, 1, 4, 9, 3, 1, 15, 6, 9,
2, 1, 11, 6, 24, 2, 11, 17, 1, 4,
7, 12, 1, 15, 6, 9, 2, 1, 3, 2, 15,
16, 2, 9, 4, 3, 2]]

[[1 3 2 15 16 2 9 4 3 2]]

[[8, 10, 4, 11, 11, 1, 5, 1, 14, 6]]

TF-IDF

Term Frequency

Number of times term t
appears in a document

*

1/Total number of terms
in the document

X

Inverse Document Frequency

log(

Total number of documents

*

1/Number of documents
with term t in it

)



1 Data Description

- Overview
- Example
- Y_Visualization
- X_Visualization

2 Preprocessing

- Tokenization
- TF-IDF

3 Model

- Model1
- Model2
- Model3
- Model4

4 Result

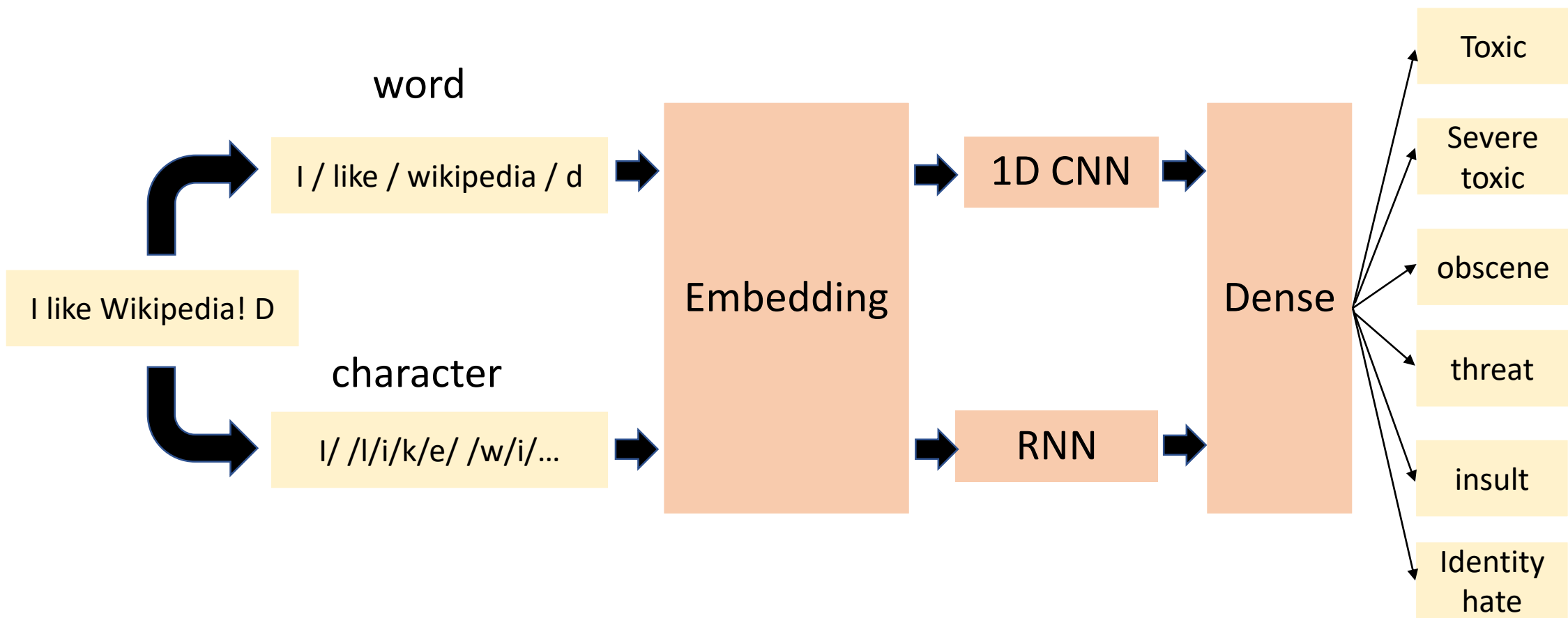
- Lime
Package

Data Description

Preprocessing

Model

Result



Data Description

Preprocessing

Model

Result

Model 1

InputLayer

Embedding

Bidirectional(LSTM)

Bidirectional(LSTM)

Dense

'i'

8

 $\begin{pmatrix} -0.08 \\ 0.17 \\ 0.04 \\ 0.06 \\ \dots \end{pmatrix}$

'like'

50

 $\begin{pmatrix} -0.14 \\ -0.08 \\ 0.14 \\ 0.04 \\ \dots \end{pmatrix}$

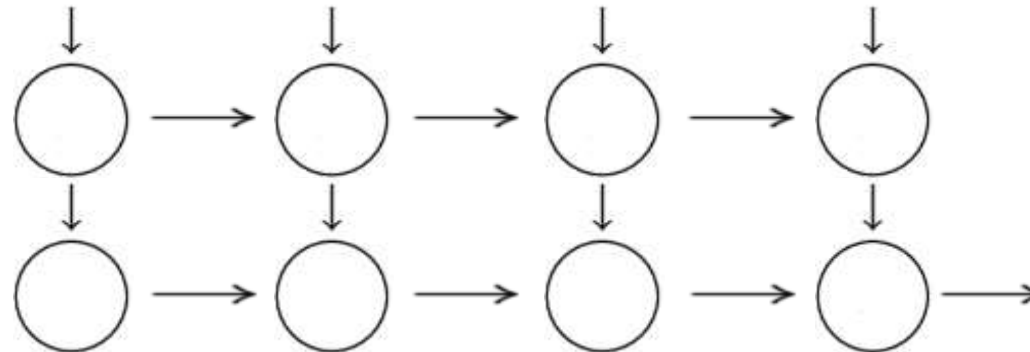
'wikipedia'

29

 $\begin{pmatrix} -0.1 \\ -0.09 \\ 0.05 \\ 0.08 \\ \dots \end{pmatrix}$

'd'

578

 $\begin{pmatrix} -0.14 \\ -0.25 \\ 0.11 \\ 0.06 \\ \dots \end{pmatrix}$ 

1.3

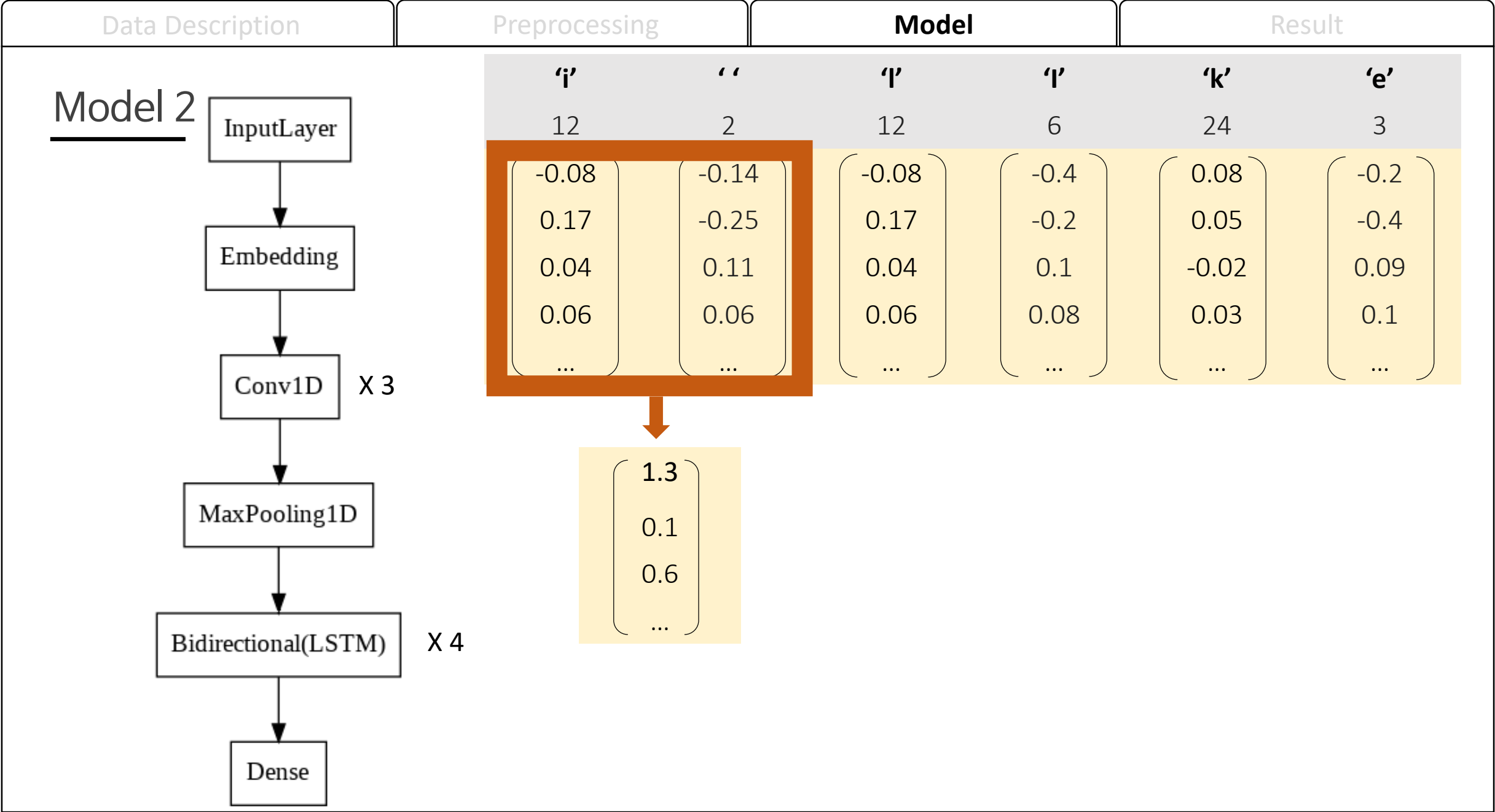
0.1

-2.4

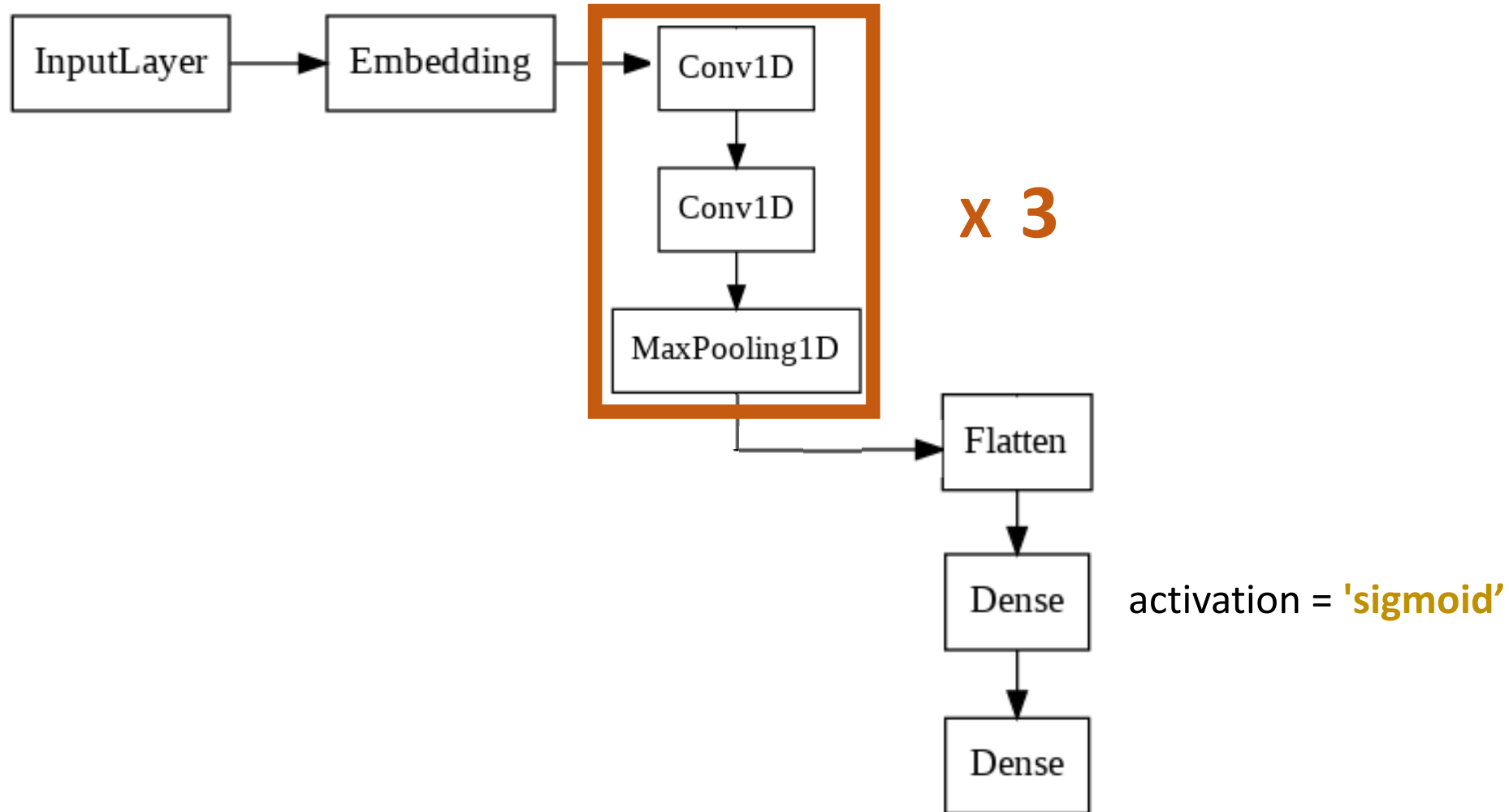
0.09

...

Roc_auc_score: 0.974

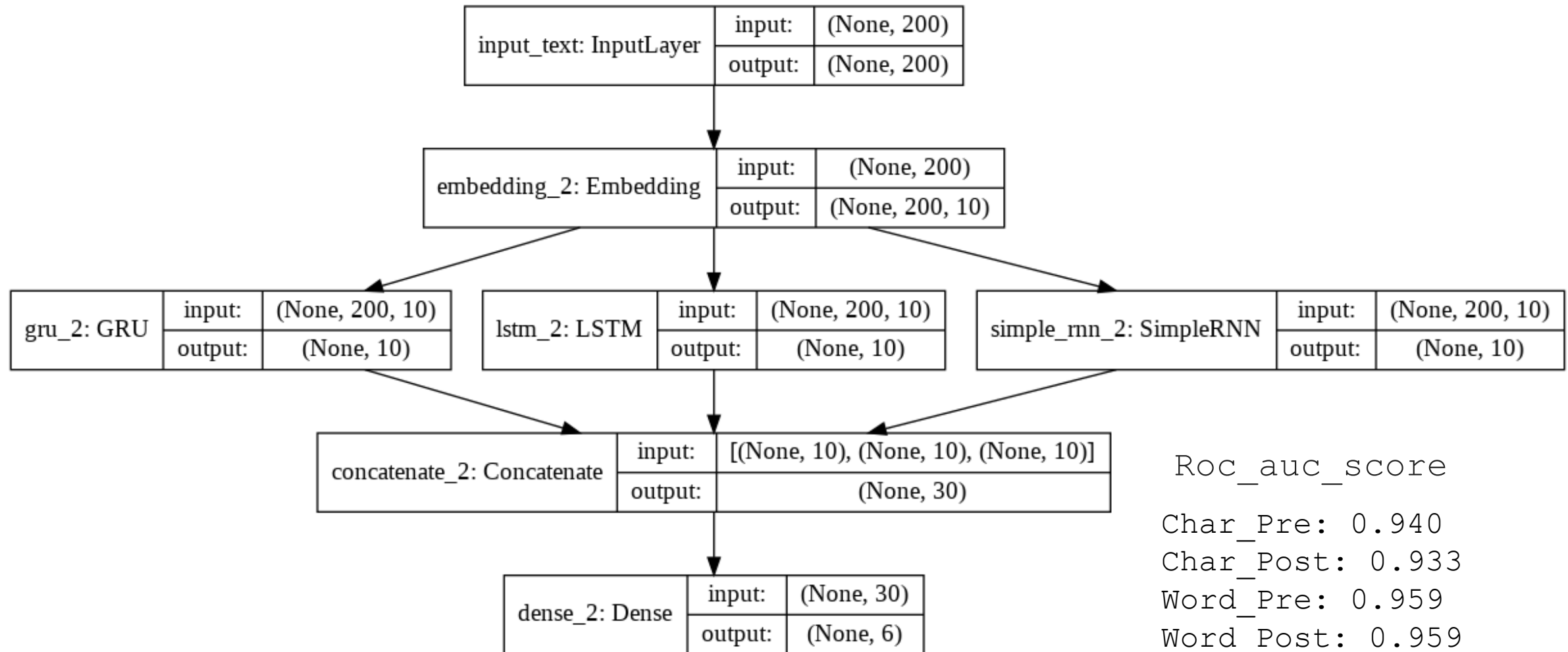


Model 3



Roc_auc_score: 0.977

Model 4



Roc_auc_score

Char_Pre: 0.940

Char_Post: 0.933

Word_Pre: 0.959

Word_Post: 0.959

Ensemble: 0.971

1 Data Description

- Overview
- Example
- Y_Visualization
- X_Visualization

2 Preprocessing

- Tokenization
- TF-IDF

3 Model

- Model1
- Model2
- Model3
- Model4



4 Result

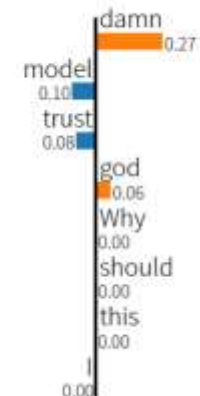
- Lime
Package

Why should I trust this god damn model?

Prediction probabilities



acceptable not_acetable



Model Interpretation : Unboxing the Black box

Which Part of the Text did make it Toxic?

Credibility

Can we rely on the model?

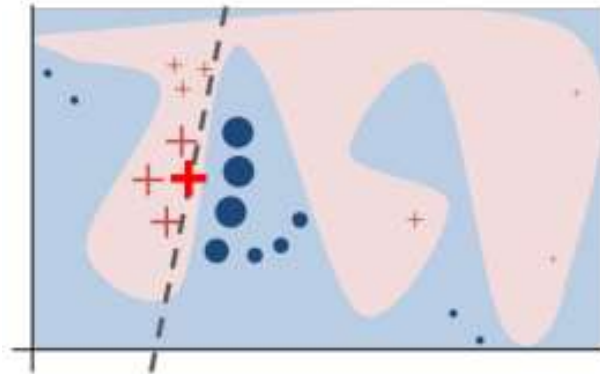
Masking

Where should we mask if necessary?

```
from LIME.lime_text import LimeTextExplainer
```

Local Interpretable

Instead of analyzing the whole model,
Let's figure out how the model work
in the neighbor of boundaries



Model-agnostic Explanation

Both Machine-Learning and Deep-Learning models
are explainable

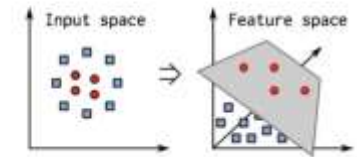
Neural Networks



Random Forests



SVM with any kernel



Data Description

Preprocessing

Model

Result

Please educate me as to how this article is a copyright infringement and, if so, from which source. Once this has been established could you explain how best to rectify the situation. Thank you

Prediction probabilities



Stupid peace of shit stop deleting my stuff asshole go die and fall in a hole go to hell!

Prediction probabilities



Data Description

Preprocessing

Model

Result

!!!roflmfao! OMFG!!! WTF??? THIS SITE IS ****ing Ga*!!! OMFG!!! roflmfao! WTF??? !!!**

acceptable

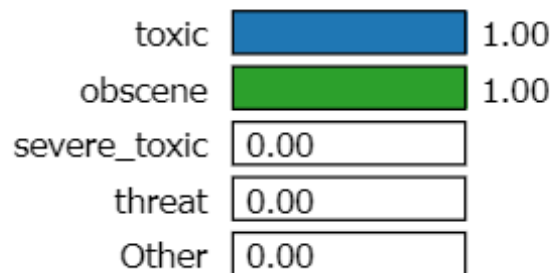
not_accetable

Prediction probabilities



Fuck you, Smith. Please have me notified when you die. I want to dance on your grave.

Prediction probabilities



NOT obscene

obscene



NOT toxic

toxic



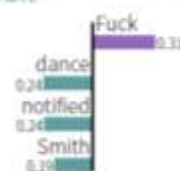
NOT identity_hate

identity_hate



NOT insult

insult



NOT threat

threat



NOT severe_toxic

severe_toxic



Q & A