

# Text Normalization Challenge - English Language

Presented by

Dahou Abdelhalim Hafedh

Pavithra POORNACHANDRAN

Tahani FENNIR



**UNIVERSITÉ  
DE LORRAINE**

November 24th, 2020

**idmC** Institut des  
sciences du Digital  
Management & Cognition

# Table of Content

## Introduction

- Definition of text normalization
- Examples

## Data Analysis

## Previous Work

## Our state

# Introduction:

- Text normalization is the process of transforming a text into a standard form.

## Examples:

- Transform word numerals into numbers (eg.: 'twenty three' → '23')
- Acronym normalization (e.g.: 'US' → 'United States'/'U.S.A')
- Removal or substitution of special characters (e.g.: remove hashtags).
- Removal of duplicate whitespaces and punctuation.

# Data Analysis:

We have an overview of training data set and testing data set.

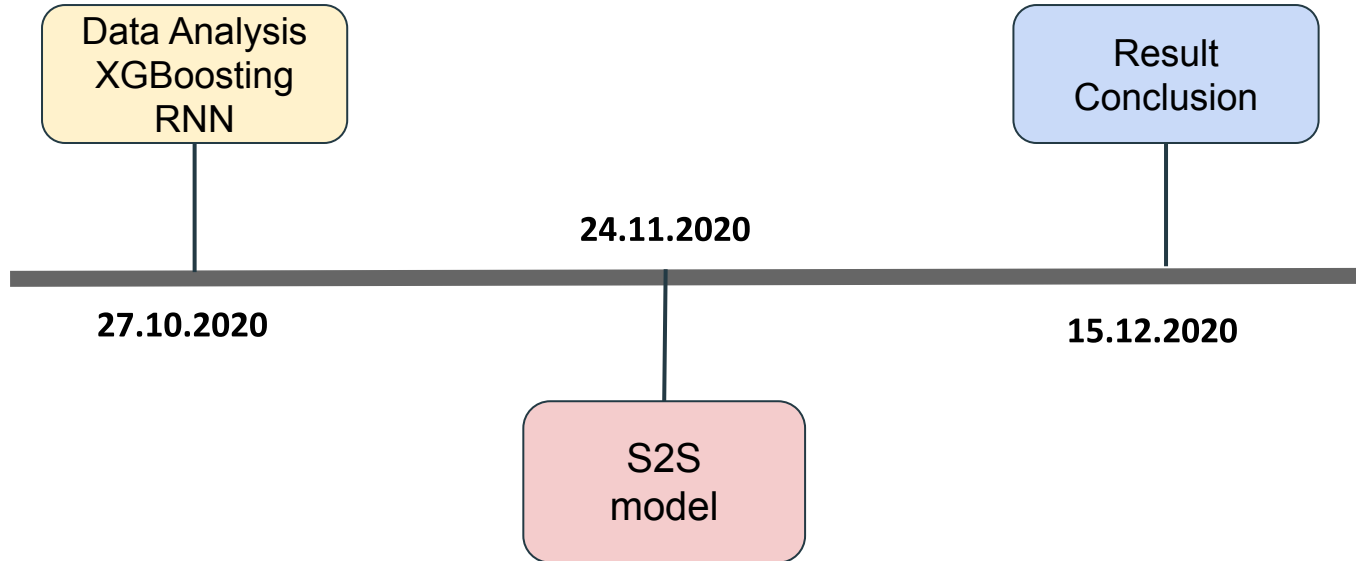
	sentence_id	token_id	class	before	after
0	0	0	PLAIN	Brillantaisia	Brillantaisia
1	0	1	PLAIN	is	is
2	0	2	PLAIN	a	a
3	0	3	PLAIN	genus	genus
4	0	4	PLAIN	of	of

( Training dataset )

	sentence_id	token_id	before
0	0	0	Another
1	0	1	religious
2	0	2	family
3	0	3	is
4	0	4	of

( Testing dataset  
)

## Previous work:



# Previous Work

## XGBoosting

- XGBoost stands for eXtreme Gradient Boosting.
- XGBoost is an algorithm that has recently been dominating applied machine learning and Kaggle competitions for structured or tabular data.
- We built XGboost model to use the context to label data.

```
[0]    valid-merror:0.00759    train-merror:0.00778
Multiple eval metrics have been passed: 'train-merror' will be used for early stopping.
```

```
Will train until train-merror hasn't improved in 20 rounds.
```

```
[10]    valid-merror:0.00547    train-merror:0.00354
[20]    valid-merror:0.00487    train-merror:0.00207
[30]    valid-merror:0.00478    train-merror:0.00099
[40]    valid-merror:0.00475    train-merror:0.00050
[49]    valid-merror:0.00478    train-merror:0.00026
```

- We choose 320,000 samples instead of 960,000 to save time and lack powerful computers.
- The accuracy of our model is 99.52%

# Neural Network

- We used recurrent neural networks to predict the outputs of the data.
- We used Keras to build our model. we used Sequential API.
- Sequential API: It is used to build models as a simple stack of layers. First, we instantiate our Sequential model object and then, you add layers to it one by one using the `add()` method.



# Neural Network

Testing :

Number of instances : 20000

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

Accuracy : 79.205

precision\_score : 0.7158730832445231

recall\_score : 0.7027046408699041

f1\_score : 0.706442272615788

Training :

Number of instances : 80000

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

Accuracy : 79.16

precision\_score : 0.70427659684227

recall\_score : 0.6851092957685128

f1\_score : 0.6903384701603337

# Current State

## Sequence-to-Sequence Model

- It can encode the entire sequence data with hidden neurons that would naturally capture any useful context information in a sentence to improve text normalization performance.

### How does it work?

- It reads the informal text sequences in encoder and transforms them in a continuous-space representation that is passed on the decoder to generate a target normalized sequence.

## Challenges :

Such an approach faces a major challenge of high percentage of OOV tokens , the model need :

- Enough training data to handle complex normalization
- Hybrid

- After exploring and analyzing different existing works in that task , we chose a good architecture to work on it and try to adapted on our data.
- Architecture is built based on **the encoder-decoder framework** both of which are parameterized by **attention-based recurrent neural networks (RNN)**.
- Proposed in : “Adapting Sequence to Sequence models for Text Normalization in Social Media” paper, 2019 by AAAI.

**Tested data :** LexNorm dataset from the 2015 ACL-IJCNLP Workshop on Noisy User-generated Text (W-NUT).

Datase t	Tweets	Tokens	Noisy	1:1	1:N	N:1	Our Vocab
Train	2950	44385	3942	2875	1043	10	10,084
Test	1967	29421	2776	2024	704	10	7389

**Result:**

Model name	Precision	Recall	F1	Method highlights
HS2S	90.66	78.14	83.94	Hybrid word-char Seq2Seq
S2S	93.39	75.75	83.65	Word-level Seq2Seq

## Train and test our dataset :

We tried to train the same model on our data , but we faced a challenge in the format of input so we developed a script that takes as input the .csv data and produce like output .json data with another organization of independent and dependent features.

We used this model to predict the outputs of the our data , but the results are very low compared with the expected values.

## Results :

Number of instance for train : 90000

Precision	Recall	F1 score
0.9598	0.8542	0.9039

Number of instance for test : 90000

Precision	Recall	F1 score
0.4054	0.070	0.1208

## Further work :

After those results , we decide to :

- Re-modified the model to deal correctly with this task.
- Add a rule-based system or a dictionary-based
- Find the good hyper-parameters to achieve high accuracy.
- Train it and test it for the whole data.



Thank you!