

# Tweets Sentiment Analysis

Kuo-Wei Ho<sup>1</sup>, Hao-Chien Wang<sup>2</sup>

<sup>1</sup>NTUEE

<sup>2</sup>NTUPhys

Data Science Programing,  
July 2019

# Table of Contents

# Problem

Given a set of data containing 1,600,000 tweets and the sentiment of each tweets. Create a model that can analyze sentiment of new tweets.

Table: Data example

sentiment	Post ID	User ID	tweets
0	1467814192	Ljelli3166	blagh class at 8 tomorrow
0	1467821455	CiaraRenee	I need a hug
4	1677796507	FoodAllergyBuzz	@otibml Thx for the tweet!
4	1677796519	lakido	Sunshine.....I LOVE this weather!!!

0: Negative

4: Positive

Data: <https://www.kaggle.com/kazanova/sentiment140>

Github link: <https://github.com/b07901135/2019dsp-summer-project>

- Vectorizing text: *GloVe* (*Global Vectors for Word Representation* by Stanford University.)
- Neural network: *RNN* (*Recurrent Neural Network*)

# Steps: Overview

- ① Clean the data: remove non-UTF8 symbols, numbers and URLs.
- ② Combine all tweets into one string and tokenize.
- ③ Feed the tokens to GloVe to generate word vectors.
- ④ tokenize all tweets and search each words in the vectors to transform it into a list of matrices.
- ⑤ Train RNN with the list of word vectors.

# Steps: Data Cleaning and Vectorization

- 1 Replace URLs as “url”
- 2 Replace name tags ( e.g. @allen1234 ) as “names”
- 3 Remove other non-UTF8 characters (*stri\_enc\_toutf8()* doesn't help)
- 4 Combine tweets into a string, tokenize and remove stopwords.
- 5 Generate TCM, feed it to the neural network to fit the model.
- 6 Generate word vectors.

# Steps: Tweets Vectorization

- Discard data other than **sentiment** and **tweets text**
- Discard tweets containing more than **30 tokens** so that the matrices will not contain too much zeros.

Table: Data manipulation

sentiment	tweets
0	blagh class at 8 tomorrow
0	I need a hug
4	@otibml Thx for the tweet!
4	Sunshine.....I LOVE this weather!!!

# Dark Magic Functions

- `save()/load()`
- `pbapply`