# 【资源】NLP多标签文本分类代码实现工具包

专知　2019-11-20

【导读】本文为大家推荐一份多标签文本分类代码实现工具包，希望对大家有所帮助。

**原文链接：**

https://github.com/RandolphVI/Multi-Label-Text-Classification

🖥 RandolphVI / **Multi-Label-Text-Classification**

About Muti-Label Text Classification Based on Neural Network.

\#text-classification  \#python3  \#tensorflow  \#sentence-classification  \#multi-label-classification

| 🕐 **2 commits** | 🔱 **1 branch** | 📦 **0 packages** | 🏷 **0 releases** | 👥 **1 contributor** | ⚖ **Apache-2.0** |
|---|---|---|---|---|---|

| Branch: master ▾ | New pull request | | | Create new file | Upload files | Find file | Clone or download ▾ |
|---|---|---|---|---|---|---|---|

| 🖼 **RandolphVI** Update README.md | | Latest commit eb9ff3a on 16 Apr |
|---|---|---|
| 📁 ANN | Initial commit | 8 months ago |
| 📁 CNN | Initial commit | 8 months ago |
| 📁 CRNN | Initial commit | 8 months ago |
| 📁 FastText | Initial commit | 8 months ago |
| 📁 HAN | Initial commit | 8 months ago |
| 📁 RCNN | Initial commit | 8 months ago |
| 📁 RNN | Initial commit | 8 months ago |
| 📁 SANN | Initial commit | 8 months ago |
| 📁 data | Initial commit | 8 months ago |
| 📁 utils | Initial commit | 8 months ago |
| 📄 .gitignore | Initial commit | 8 months ago |
| 📄 .travis.yml | Initial commit | 8 months ago |
| 📄 LICENSE | Initial commit | 8 months ago |
| 📄 README.md | Update README.md | 7 months ago |
| 📄 requirements.txt | Initial commit | 8 months ago |

📖 **README.md**

# Deep Learning for Multi-Label Text Classification

`language python3.6`  `build passing`  `code quality A`  `license Apache-2.0`  `issues 11 open`

This repository is my research project, and it is also a study of TensorFlow, Deep Learning (Fasttext, CNN, LSTM, etc.).

The main objective of the project is to solve the multi-label text classification problem based on Deep Neural Networks. Thus, the format of the data label is like [0, 1, 0, ..., 1, 1] according to the characteristics of such a problem.

## Requirements

- Python 3.6
- Tensorflow 1.1 +
- Numpy
- Gensim

## Innovation

### Data part

1. Make the data support **Chinese** and English (Which use `jieba` seems easy).
2. Can use **your own pre-trained word vectors** (Which use `gensim` seems easy).

3. Add embedding visualization based on the **tensorboard**.

## Model part

1. Add the correct **L2 loss** calculation operation.
2. Add **gradients clip** operation to prevent gradient explosion.
3. Add **learning rate decay** with exponential decay.
4. Add a new **Highway Layer** (Which is useful according to the model performance).
5. Add **Batch Normalization Layer**.

## Code part

1. Can choose to **train** the model directly or **restore** the model from the checkpoint in `train.py`.
2. Can predict the labels via **threshold** and **top-K** in `train.py` and `test.py`.
3. Can calculate the evaluation metrics --- **AUC** & **AUPRC**.
4. Add `test.py`, the **model test code**, it can show the predicted values and predicted labels of the data in Testset when creating the final prediction file.
5. Add other useful data preprocess functions in `data_helpers.py`.
6. Use `logging` for helping to record the whole info (including **parameters display**, **model training info**, etc.).
7. Provide the ability to save the best n checkpoints in `checkmate.py`, whereas the `tf.train.Saver` can only save the last n checkpoints.

# Data

See data format in `data` folder which including the data sample files.

## Text Segment

You can use `jieba` package if you are going to deal with the Chinese text data.

## Data Format

This repository can be used in other datasets (text classification) in two ways:

1. Modify your datasets into the same format of the sample.
2. Modify the data preprocess code in `data_helpers.py`.

Anyway, it should depend on what your data and task are.

😊Before you open the new issue about the data format, please check the `data_sample.json` and read the other open issues first, because someone maybe ask me the same question already. For example:

- 输入文件的格式是什么样子的?
- Where is the dataset for training?
- 在 data_helpers.py 中的 content.txt 与 metadata.tsv 是什么，具体格式是什么，能否提供一个样例?

## Pre-trained Word Vectors

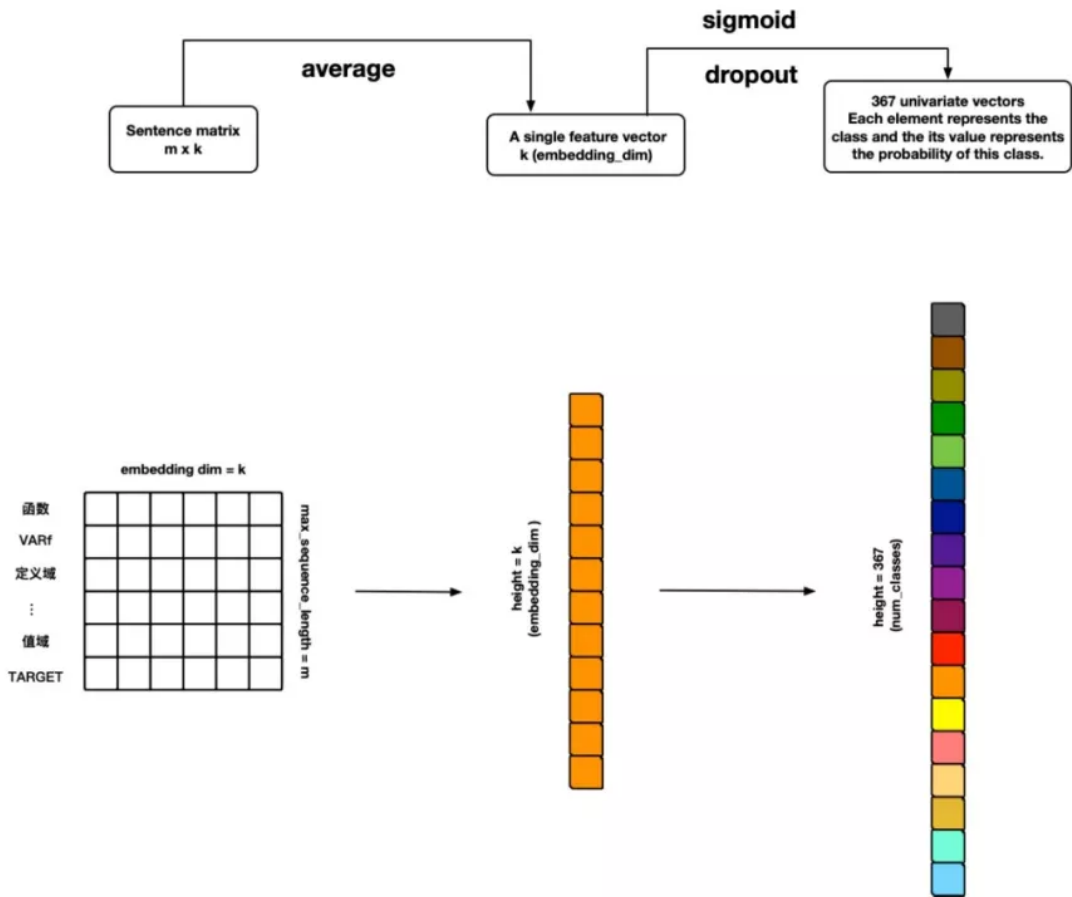You can pre-training your word vectors (based on your corpus) in many ways:

- Use `gensim` package to pre-train data.
- Use `glove` tools to pre-train data.
- Even can use a **fasttext** network to pre-train data.

# Network Structure

## FastText

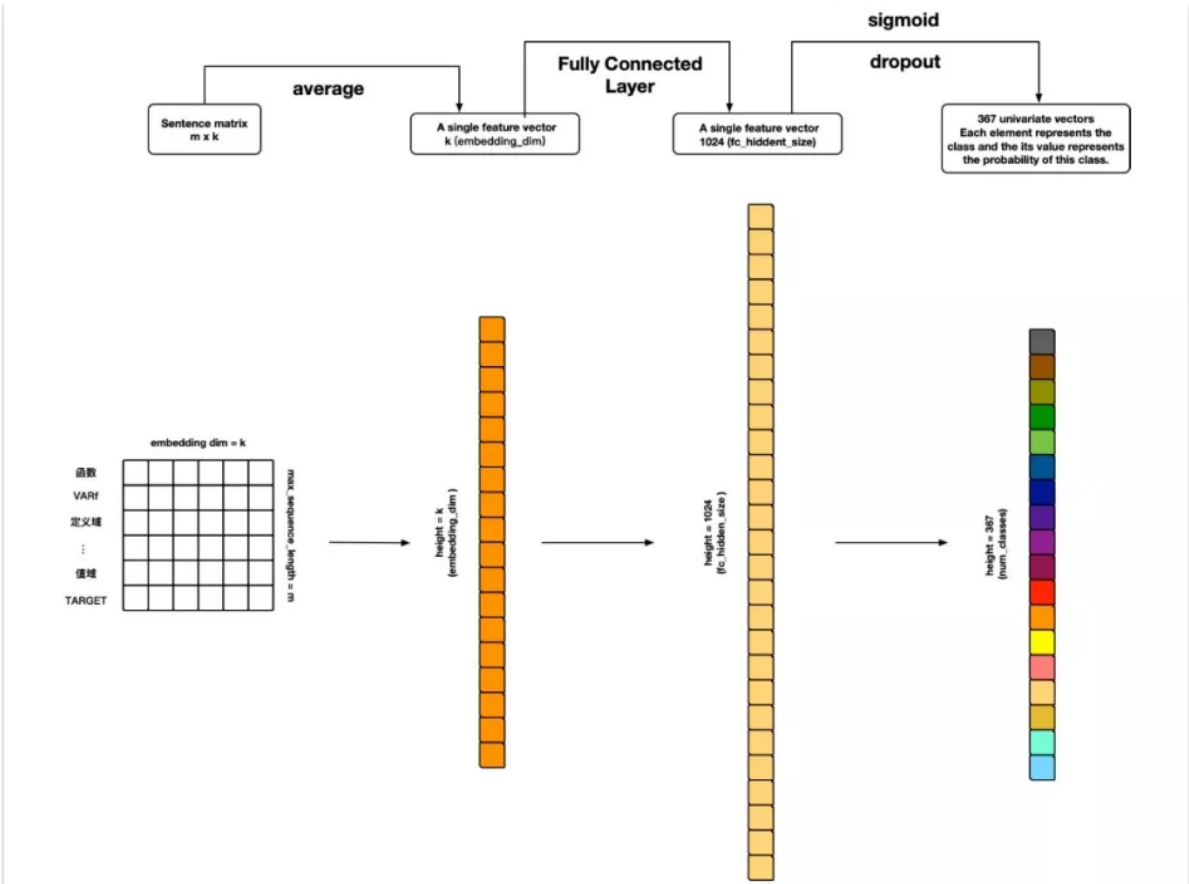2019/11/20    RandolphVl/Multi-Label-Text-Classification: About Muti-Label Text Classification Based on Neural Network.



References:
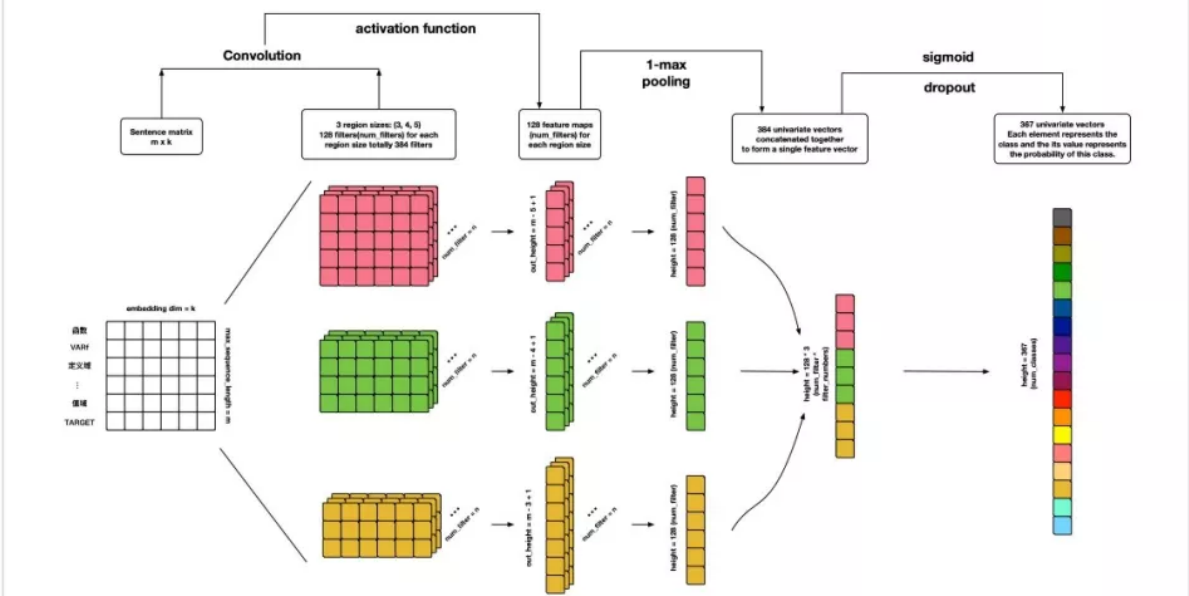
- Bag of Tricks for Efficient Text Classification

## TextANN

References:

- Personal ideas 😊

## TextCNN



References:

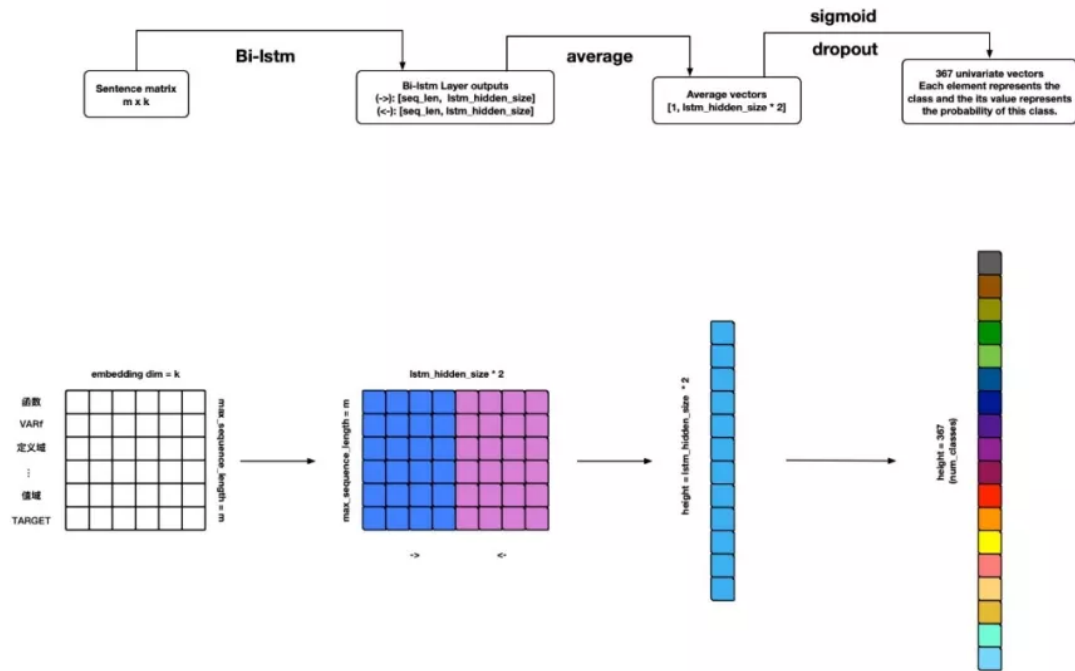- Convolutional Neural Networks for Sentence Classification

- A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification

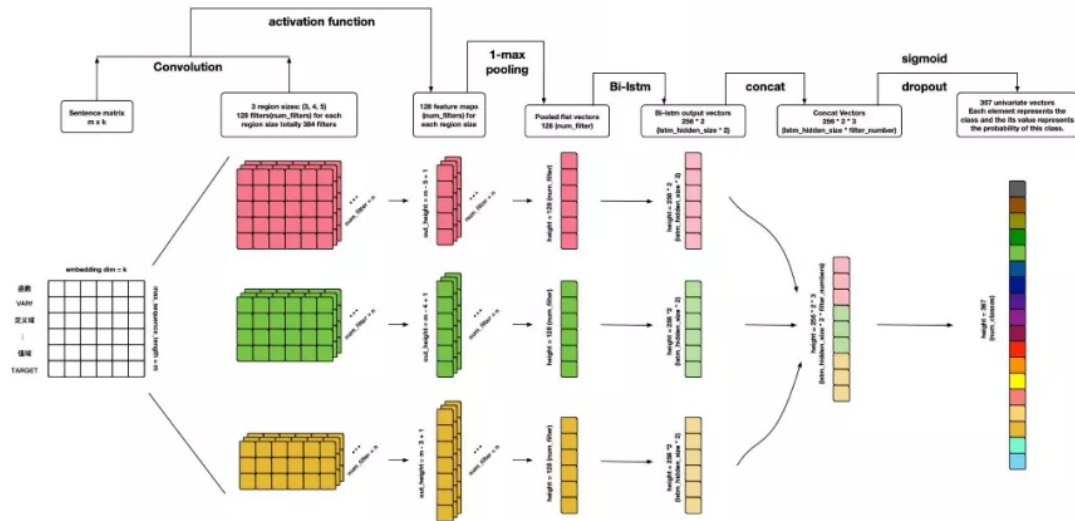## TextRNN

Warning: Model can use but not finished yet 😊!



## TODO

1. Add BN-LSTM cell unit.
2. Add attention.

References:

- Recurrent Neural Network for Text Classification with Multi-Task Learning
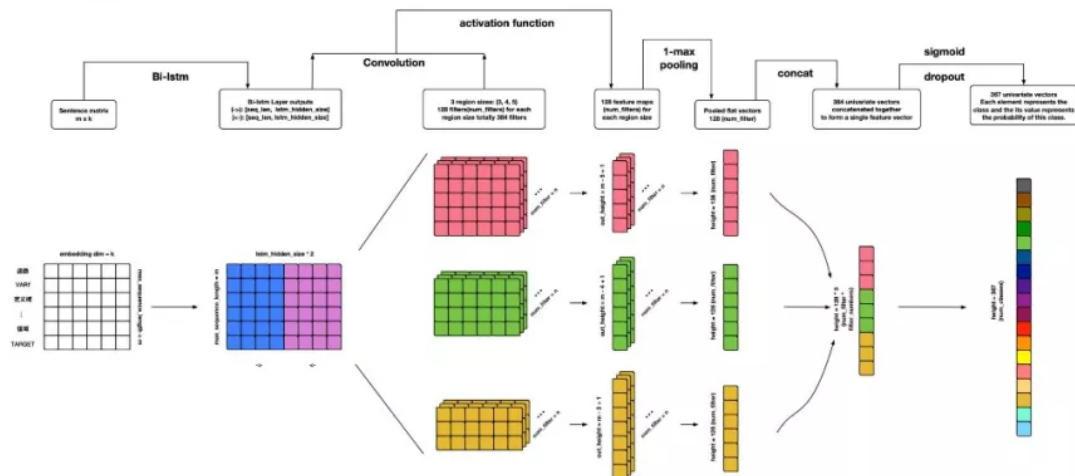
## TextCRNN



References:

- Personal ideas 🙂

## TextRCNN



References:

- Personal ideas 🙂

## TextHAN

References:

- Hierarchical Attention Networks for Document Classification

## TextSANN

Warning: Model can use but not finished yet 😅!

TODO

1. Add attention penalization loss.
2. Add visualization.

References:

- A STRUCTURED SELF-ATTENTIVE SENTENCE EMBEDDING

## About Me

黄威，Randolph

SCU SE Bachelor; USTC CS Master

Email: chinawolfman@hotmail.com

My Blog: randolph.pro

LinkedIn: randolph's linkedin

-END-

专·知