# 自然語言處理與文字探勘

陳縕儂&教研處

# 課程內容

講師投影片
資料與投影片
影片播放列表
程式碼:~/courses-tpe/NLP

1. 文本特徵建模練習

2. NLP with RNN

3. BERT fine-tune

台灣人工智慧學校

# Code / Data 放在 hub 中的 courses 內

- 為維護課程資料 , courses 中的檔案皆為 read-only, 如需修改請 cp 至自身的環境中
- 打開 terminal, 輸入
  - [台北班]
    cp -r courses-tpe/NLP/part4 <存放至本機的名稱>
  - [新竹班]
    cp -r courses-hsi/NLP/part4 <存放至本機的名稱>
  - [台中班]
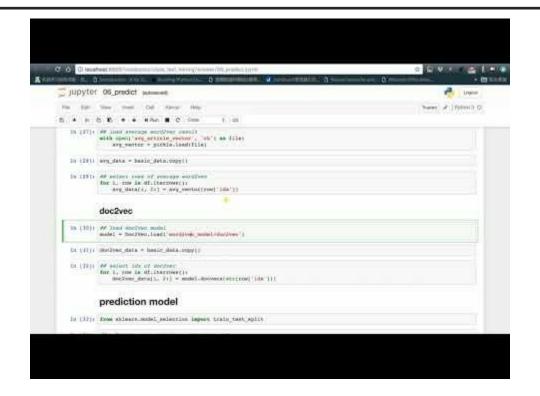    cp -r courses-txg/NLP/part4 <存放至本機的名稱>

台灣人工智慧學校

程式實作

# 文本特徵建模練習

# 程式練習時間

- 06_predict.ipynb
  - 定義 y: 發文的推文多 or 噓文多？
    - 絕對差異定義？> 20
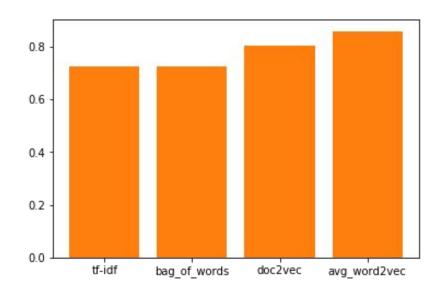  - 利用 bag of words, TF-IDF, average word2vec, doc2vec 各做一組 prediction model, 比較哪一組 features 最好
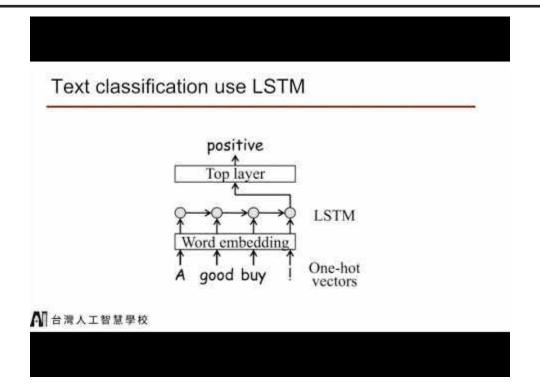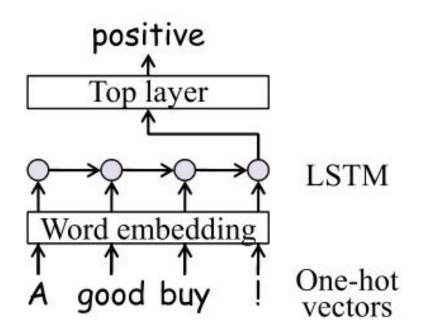
# 程式解說

# Each model's AUC

程式實作

# NLP with RNN

# 建模流程



Text classification use LSTM

positive
Top layer
LSTM
Word embedding
A good buy !
One-hot vectors

台灣人工智慧學校

# Text classification use LSTM

# how to feed text data into LSTM network ?

- 先把 text data 轉換成 id 格式
  - word to id
- 不在 word vector 的字，用別的 id 代表
  - e.g. len(word)+1

['韓瑜', '協志', '前妻', '正', '女演員', '周子', '瑜', 'TWICE', '團裡裡面', '台灣', '人', '正', '兩個', '要當', '鄉民', '老婆', '選', '五樓', '真', '勇氣']
[10461, 27588, 84244, 23278, 84491, 90934, 31569, 72550, 100035, 84284, 96798, 23278, 96689, 31004, 62798, 53752, 92708, 44764, 66642, 10179, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034, 100034]

# embedding lookup table

```python
self.inputs = tf.placeholder(tf.int32, [None, input_length], name='input_data')
self.targets = tf.placeholder(tf.float32, [None, 1], name='targets')
self.bz = tf.placeholder(tf.int32, [], name='batch_size')

## embedding lookup table
em_W = tf.Variable(wv.astype(np.float32), trainable=True)
x = tf.nn.embedding_lookup(em_W, self.inputs)
```
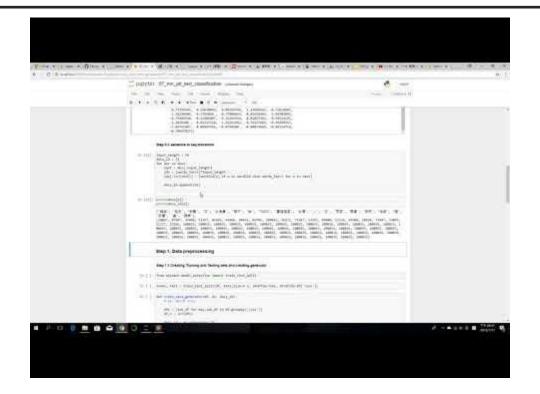
# 程式練習時間

- 07_rnn_ptt_text_classification.ipynb
  - 練習如何把 text data 轉換成 id 格式
  - 練習建構 embedding lookup table 接 LSTM network
  - 可嘗試
    - fine tune word vector (trainable = True)
    - 不 fine tune word vector (trainable = False)
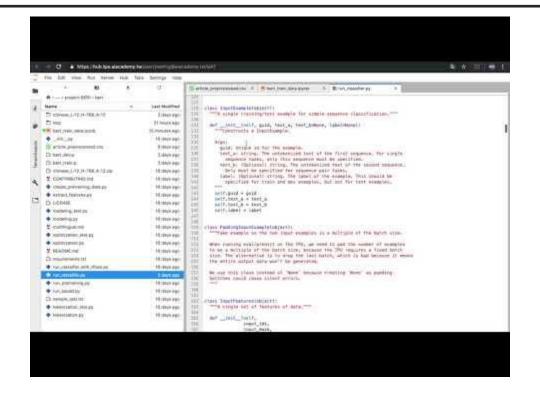    - randon initialize embedding lookup table

台灣人工智慧學校

# 程式解說

# BERT fine-tune

# BERT fine-tune

# 歡迎挑戰 kaggle text classification task

- [Toxic Comment Classification Challenge](#)
- [Mercari Price Suggestion Challenge](#)
- [Spooky Author Identification](#)
- [Personalized Medicine: Redefining Cancer Treatment](#)
- [Quora Question Pairs](#)
- ...

# 還不夠嗎？可以看看別的主題

- N-Gram
  - [Modelling Natural Language with N-Gram Models](#)
  - [自然語言處理中的 N-Gram 模型介紹](#)
- Topic model
  - [Begineer guide to topic modelling in python](#)
  - [Topic modelling with scikit-learn](#)
  - [Latent Dirichlet Allocation](#)
  - [LDA 數學八卦](#)
  - [手刻版 topic model](#)

台灣人工智慧學校