游馬豐智在不易出的應用

姐員:謝奇容、唐祥勝

深度學習

深度學習與文字

自然語言處理

深度學習在文字上的應用

應用實例

實作

目錄

什麼是深意學習?

深度學習

- •機器學習的一種、現今人工智慧的主流
- •學習;經過大量的訓練過程,找出一個最佳函數,得出最佳解
- 針對特定需求設計
- •用"規則"及"大量學習資料",告訴機器什麼是對的
- 使用深度神经網路

深度學習

CNN

(Convolutional Neural Network)

RNN

(Recurrent Neural Network)

GAN

(Generative Adversarial Network)

處理空間上連續的資料 影像辨識

處理有時間序列、語意 结構的資料

讓電腦不會忘記自己做 過的事

深度神经網路

探度學習遇到文字的困難

深度學習 VS 文字

- •電腦無法讀取文字,僅能讀取數字
- 語文不像數學,需要根據前後文才能理解意義
- •根據地域性會有當地的方言、用法

自然語言處理NLP

自然語言處理

- •程式語言:人類為了與電腦溝通而設計的語言
- •自然語言:人類與人類溝通所發展、創造出的語言



自然語言處理

自然語言處理

- 文本分詞
- 將文字轉成數字序列
- 建立模型
- 定義衡量指標
- •訓練模型
- 用已訓練模型對新數據做預測

文本分詞

- 將非結構化的文本數據,轉成結構化數據
- 詞是表達含意的最小單位
- 中英文分詞區別:
 - ■分詞方式不同
 - ■詞的型態分別不同
- 中文分詞的難點:
 - ■分詞沒有統一標準
 - ■歧義詞切分
 - ■新詞識別

文本分詞

- 中文分詞工具:
 - Hanlp
 - Stanford 分詞
 - ansj 分詞器
 - ■哈工大 LTP
 - KCWS分詞器
 - jieba
 - IK
 - ■清華大學THULAC
 - ICTCLAS

Jieba



精確模式

分詞模式: 全模式

搜尋引擎模式



自定義字典



探發學

- TensorFlow
- Keras
- PyTorch
- CRF
- Kashgari

Kashgari

- 集中多種模型
- 擁有多種預訓練模型
- 方便快速訓練一個序列標註或 文字分類模型
- 擁有統一的 和説明文件

Kashgari

- 無監督式訓練
- 開發的語言代表模型
- · 先以 的方式預先訓 鍊出一個對自然語言有一定「理解」 的通用模型
- 再將該模型拿來做特徵擷取或是下游的(監督式)任務



應用實例



玉山人工智慧公開挑戰賽比賽目的:透過,抓出 提升銀行效率 夏季賽 焦點人物,

```
bert_embed = BERTEmbedding('code/chinese_L-12_H-768_A-12',
                                task=kashgari.LABELING,
                                sequence length=100)
model = LSTM CNN Model(bert embed)
# This step will build token dict, label dict and model structure
model.build model(train x, train y, valid x, valid y)
# Compile model with custom optimizer, you can also customize loss and metrics.
# optimizer = RAdam()
# model.compile model(optimizer=optimizer, loss=categorical focal loss(gamma=2.0, alpha=0.25))
eval callback = EvalCallBack(kash model=model,
                             valid_x=valid_x,
                             valid y=valid y,
                             step=1)
model.compile model()
# Train model
model.fit(train_x, train_y, valid_x, valid_y, batch_size=128, epochs=2, callbacks=[eval_callback])
```

Kashgari實作

資料來源

- https://kknews.cc/zh-tw/code/rlyeezn.html
- https://blog.kennycoder.io/2020/02/12/Python-%E7%9F%A5%E5%90%8DJieba%E4%B8%AD%E6%96%87 %E6%96%B7%E8%A9%9E%E5%B7%A5%E5%85%B7%E6% 95%99%E5%AD%B8/
- https://research.sinica.edu.tw/nlp-natural-language-processing-chinese-knowledge-information/
- https://medium.com/@leemeng/%E9%80%B2%E5%85% A5-nlp-%E4%B8%96%E7%95%8C%E7%9A%84%E6%9C%80%E4% BD%B3%E6%A9%8B%E6%A8%91-1e90e21f3838
- https://medium.com/%E5%AD%B8%E4%BB%A5%E5%BB %A3%E6%89%8D/nlp%E7%99%BC%E5%B1%95%E5%8F% B2%E6%91%98%E8%A6%81-1-a30af62cbcec
- https://www.itread01.com/content/1565680922.html
- https://panx.asia/archives/53209
- https://research.sinica.edu.tw/deep-learning-2017-aimonth/

分工

姓名 謝奇容 唐祥勝 工作 資料整理 文本分詞 製作

7 HANK