AI共學社群 我的 AI共學社群 > Part1 - NLP經典機器學習馬拉松 > D20 KNN 實作 D20 KNN 實作 **=**~ 範例與作業 問題討論 學習心得(完成) 簡報閱讀 重要知識點 NLP自然語言學習實戰馬拉松 本日實作內容 ▶ Day20 - KNN 實作 資料讀取 資料清理 文字轉向量 陪跑專家:楊哲寧 重要知識點 重要知識點 • 了解如何實踐 KNN 算法,並使用 K-fold 選擇適合 K 值。 補充資料: K-Nearest Neighbors and Bias-Variance Tradeoff 本日實作內容 A. 資料讀取 B. 資料清理 C. 文字轉向量 D. KNN模型導入與參數設置 E. 模型驗證 F. 混淆矩陣 (Confusion matrix) G. K-fold尋找適合K值 資料讀取 這裡我們將 Spam(垃圾郵件標註為 1), 反之 (Ham) 為 0。 Importing the dataset • 從三個資料夾中讀取資料 In [2]: all\_data=[] paths =[r'spam\_data/spam', r'spam\_data/easy\_ham', r'spam\_data/hard\_ham'] for path in paths: for fn in glob.glob(path+"/\*"): if "ham" not in fn: is\_spam = 1 #codecs.open可以避開錯誤,用errors='ignore'
with codecs.open(fn,encoding='utf-8', errors='ignore') as file: for line in file: #這個line的開頭為Subject: if line.startswith("Subject:"):
 subject=re.sub(r"^Subject:","",line).strip() all\_data.append([subject,is\_spam]) all\_data = np.array(all\_data) 1 all\_data=[] paths =[r'spam\_data/spam', r'spam\_data/easy\_ham', r'spam\_data/hard\_ham'] 3 for path in paths: 4 for fn in glob.glob(path+"/\*"): 5 if "ham" not in fn: is\_spam = 1 else: is\_spam = 0 #codecs.open可以避開錯誤,用errors='ignore' with codecs.open(fn,encoding='utf-8', errors='ignore') as file: 11 for line in file: #這個line的開頭為Subject: 12 資料清理 細節可參考先前文字預處理章節,或是歸納出自己的清理方式也可以。 In [6]: from sklearn.metrics import confusion\_matrix from nltk.corpus import stopwords import nltk nltk.download('stopwords') # Lemmatize with POS Tag from nltk.corpus import wordnet from nltk.stem import WordNetLemmatizer ## ##Lemmatizer lemmatizer = WordNetLemmatizer() def get\_wordnet\_pos(word):
 """將pos\_tag結果mapping到lemmatizer中pos的格式""" tag = nltk.pos\_tag([word])[0][1][0].upper() tag dict = {"J": wordnet.ADJ, "N": wordnet.NOUN, "V": wordnet.VERB, "R": wordnet.ADV) return tag\_dict.get(tag, wordnet.NOUN) def clean\_content(X):
 # remove non-alphabet characters  $X_{clean} = [re.sub('[^a-zA-z]',' ', x).lower() for x in X]$ X\_word\_tokenize = [nltk.word\_tokenize(x) for x in X\_clean]
# stopwords\_lemmatizer X\_stopwords\_lemmatizer = []
stop\_words = set(stopwords.words('english')) for content in X\_word\_tokenize: content\_clean = [] for word in content: if word not in stop\_words: word = lemmatizer.lemmatize(word, get\_wordnet\_pos(word)) content\_clean.append(word) X\_stopwords\_lemmatizer.append(content\_clean) X\_output = [' '.join(x) for x in X\_stopwords\_lemmatizer] return X\_output [nltk\_data] Downloading package stopwords to
[nltk\_data] /Users/jeff.yang/nltk\_data...
[nltk\_data] Package stopwords is already up-to-date! In [7]: X = clean\_content(X) 1 from sklearn.metrics import confusion\_matrix 2 from nltk.corpus import stopwords 4 import nltk 6 nltk.download ('stopwords') 8 # Lemmatize with POS Tag 9 from nltk.corpus import wordnet 10 from nltk.stem import WordNet Lemmatizer 12 ## 創建Lemmatizer 1 def clean\_content(X): 2 # remove non-alphabet characters 3 X\_clean = [re.sub ('[a-zA-Z]',' ',x).lower() for x in X] 4 # tokenize 5 X\_word\_tokenize = [nltk.word\_tokenize(x) for x in X\_clean] 6 # stopwords\_lemmatizer 7 X\_stopwords\_lemmatizer = [] 8 stop\_words = set (stopwords.words ('english')) 9 for content in X\_word\_tokenize : 10 content\_clean = [] 11 for word in content: if word not in stop\_words: 1 X = clean\_content(X) 文字轉向量 這邊使用 BOW 方式,學員們也可以嘗試 TFIDF。 Bag of words cv=CountVectorizer(max\_features = 1500) X=cv.fit\_transform(X).toarray() In [9]: # 有 3423個樣本,每個樣本用1500維表示 Out[9]: (3423, 1500) 1 from sklearn.feature\_extraction.text import CountVectorizer 2 #max\_features是要建造幾個column,會按造字出現的頻率高低去篩選,1500並沒有特別含義,大家可以自己嘗試不同數值 3 cv=CountVectorizer(max\_features = 1500) 4 x=cv.fit\_transform(X).toarray() 1 # 有 3423個樣本,每個樣本用1500維表示 2 X.shape KNN 模型導入與參數設置 • 首先將資料拆解成 trainset/testset (詳細可見參考資料) • 從 sklearn 中呼叫 KNN classifier API 將資料拆成 train/test set In [10]: from sklearn.model\_selection import train\_test\_split # random\_state 是為了讓各為學員得到相同的結果,平時可以移除 X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, Y, test\_size = 0.2, random\_state = 0) Training the K-NN model on the Training set In [18]: from sklearn.neighbors import KNeighborsClassifier classifier = KNeighborsClassifier(n\_neighbors = 5, metric = 'minkowski', p = 2) classifier.fit(X\_train, y\_train) Out[18]: KNeighborsClassifier(algorithm='auto', leaf\_size=30, metric='minkowski', metric\_params=None, n\_jobs=None, n\_neighbors=5, p=2,
weights='uniform') 1 from sklearn.model\_selection import train\_test\_split 2 # random\_state 是為了讓各位學員得到相同的結果,平時可以移除 3 X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, Y, test\_size = 0.2, random\_state = 0) 1 from sklearn.neighbors import KNeighborsClassifier 2 classifier = KNeighborsClassifiler(n\_neighbors = 5, metric = 'minkowski', p = 2) 3 classifier.fit(X\_train, y\_train) • n\_neighbors : K 值,為 hyperparameters(超參數),可透過 K-fold 協助選擇最適合 K 值。 • metric:計算樣本距離的方式,默認設置為 Minkowski,其為歐式距離及曼哈頓距離兩種計算 距離的延伸。 • P:為 Minkowski metric 的超參數, p 為 2 時所使用歐式距離。 模型驗證 測試模型在訓練/測試集上的 Accuracy。 測試 train/testset的 Accuracy In [19]: print('Trainset Accuracy: {}'.format(classifier.score(X\_train, y\_train))) Trainset Accuracy: 0.91672753834916 In [20]: print('Testset Accuracy: {}'.format(classifier.score(X\_test, y\_test))) Testset Accuracy: 0.8875912408759125 print('Trainset Accuracy: {}" .format(classifier.score(X\_train, y\_train))) print('Testset Accuracy: {}" .format(classifier.score(X\_test, y\_test))) 混淆矩陣 (Confusion matrix) 矩陣對應 TP、FN、FP、TN Making the Confusion Matrix In [23]: from sklearn.metrics import confusion\_matrix, accuracy\_score cm = confusion\_matrix(y\_test, y\_pred) accuracy\_score(y\_test, y\_pred) [[581 6] [71 27]] Out[23]: 0.8875912408759125 1 from sklearn.metrics import confusion\_matrix, accuracy\_score 2 cm = confusion\_matrix(y\_test, y\_pred) 3 print(cm) 4 accuracy\_score(y\_test, y\_pred) predicted condition total population prediction positive prediction negative condition False Negative (FN) True Positive (TP) (type II error) positive true condition False Positive (FP) condition True Negative (TN) negative (Type I error) 協助我們更深入的了解正確與錯誤數據的分佈,主要可以分為以下四種 TP(True Positive): 正確預測成功的正樣本,例如在一個預測是不是貓的 圖像分類器,成功的把一張貓的照片標示成貓,即為TP TN(True Negative): 正確預測成功的負樣本,以上述例子,成功的把一張 狗的照片標示成不是貓,即為TN FP(False Positive): 錯誤預測成正樣本,實際上為負樣本,例如:錯誤的 把一張狗的照片標示成貓 FN(False Negative): 錯誤預測成負樣本(或者說沒能預測出來的正樣本), 例如:錯誤的把一張貓的照片標示成不是貓 資料來源:<u>混淆矩陣(confusion matrix)介紹</u> • 我們常用的 Accuracy(正確數量/樣品總數)能傳達的訊息並不夠豐富,常用的 metrics 還包含 Recall/Precision/F1,這些資訊可藉由混淆矩陣進一步推算獲得。 • 這裡我們不詳談其意義,有興趣的學員們可參考參考資料中所附上的連結。 實際 Positive 實際 Negative TP (True Positive) FP 預測 Positive (False Positive) 預測 Negative (True Negative) (False Negative) Precision =  $\frac{TP}{TP + FP}$ Recall =  $\frac{TP}{TP + FN}$ 資料來源: Precision、Recall、F1 三種評估模型的指標 K-fold 尋找適合 K 值 scikit-learn 提供相當方便的 API 直接合併 K-fold 與 Classifier。 運用K-fold尋找適合K值 In [13]: # Applying k-Fold Cross Validation #n-jobs=-1,是指cpu全開 from sklearn.model\_selection import cross\_val\_score from sklearn.neighbors import KNeighborsClassifier n\_neighbors = [5, 10, 25, 50, 100] ## 可自行嘗試不同K值 for k in n\_neighbors:
 classifier = KNeighborsClassifier(n\_neighbors = k, metric = 'minkowski', p = 2)
# cv = 10 代表切成10等分 accuracies = cross\_val\_score(estimator = classifier, X = X\_train, y = y\_train, cv = 10,n\_jobs=-1) print('設置K值:{}'.format(k))
print('Average Accuracy: {}'.format(accuracies.mean())) print('Accuracy STD: {}'.format(accuracies.std())) Average Accuracy: 0.880947378655408 Accuracy STD: 0.011561254412983857 Average Accuracy: 0.8623193886405567 Accuracy STD: 0.006510165116834898 設置K值:25 Average Accuracy: 0.8528182474459847 Accuracy STD: 0.0024837372232837055 Average Accuracy: 0.8520856467133839 Accuracy STD: 0.0014068907077920568 Average Accuracy: 0.8520856467133839 Accuracy STD: 0.0014068907077920568 1 # Applying k-Fold Cross Validation 2 #n-jobs=-1,是指cpu全開 3 from sklearn.model\_selection import cross\_val\_score 4 from sklearn.neighbors import KNeighborsClassifier 5 n\_neighbors = [5, 10, 25, 50, 100] ## 可自行嘗試不同K值 6 for k in n\_neighbors: 7 classifier = KNeighborsClassifier(n\_neighbors = k, metric = 'minkowski', p = 2) 8 # cv = 10 代表切成10等分 accuracies = cross\_val\_score(estimator = classifier, X\_train, y = y\_train, cv = 10,n\_jobs=-1) 11 pint('設置K值:{}'.format(k)) 12 print('Average Accuracy: {}'.format(accuracies.mean())) KNN 優缺點整理 優點 缺點 1.需要大量運算資源 1.原理簡單 2.需要較大的記憶體空間 2.在相對簡單的資料集上都能獲得不錯的精度 3.不適合用在大型資料集 3.資料並不需要特別的前處理(沒有統計假設) 4.不適合處理高維度資料(Curse of Dimensionality) 參考資料 trainset/testset/validation set 用意: 網站:訓練集、驗證集、測試集 一文看懂 AI 訓練集、驗證集、測試集(附:分割方法+交叉驗證) 區塊鏈遊戲研究院 · 2019年12月20日20:01 0 數據在人工智能技術裏是非常重要的!本篇文章將詳細給大家介紹 3 種數據集:訓練集 驗證集、測試集。 同時還會介紹如何更合理的講數據劃分爲 3 種數據集。最後給大家介紹一種充分利用有限 數據的方式:交叉驗證法。 先用一個不恰當的比喻來說明 3 種數據集之間的關係: • 訓練集相當於上課學知識 in • 驗證集相當於課後的的練習題,用來糾正和強化學到的知識 • 測試集相當於期末考試,用來最終評估學習效果 高維度資<mark>料對模型的影響:</mark> 網站: Curse of Dimensionality k-Nearest Neighbors and the Curse of Dimensionality Why the k-nearest neighbors algorithm is especially sensitive to additional dimensions Peter Grant Follow

emboth

有一種分類問題常用的指標稱之為Confusion Matrix,這個命名很有趣,這個表格的確是很讓人感到很

困惑啊!至少在看完這篇之前。Confusion Matrix是用於分類問題的一種常用的指標,它衍生很多不同

的指標,下面這張圖我將Confusion Matrix畫出來,並把一些比較重要的衍生指標給標出來。

Jul 23, 2019 · 6 min read ★

詞就攤在那邊,讓人無從下手。