അവ AI共學社群 我的 AI共學社群 > Part1 - NLP經典機器學習馬拉松 > D17 計數方法詞向量實作介紹 D17 計數方法詞向量實作介紹 囯 簡報閱讀 範例與作業 問題討論 重要知識點 NLP自然語言學習實戰馬拉松 字詞前處理 ▶ Day17 - 計數方法詞向量實作介紹 共現矩陣 詞向量相似度 陪跑專家:Leo Liou 劉冠宏 重要知識點 重要知識點 • 透過實作更加了解計數方法詞向量 了解如何使用 SVD 對詞向量進行降維 (本日課程為前日課程的實作,同學可搭配前日課程的講義進行學習) 字詞前處理 在進行字詞向量化之前,我們需要針對文本資料進行前置處理,針對文本資料處理的方法很多,本日 課程會採用以下三種前處理程序: • 文本資料分割成字詞(斷詞) 字詞轉換成字詞 ID 清單 • 字詞移除標點符號 文本轉化為數字的 ID 形式表示 將所有字詞轉為小寫並移除標點符號: 1 #定義簡易文本資料(使用Ch17講義中的例子) corpus = ['You say goodbye and I say hello.'] word_dic = set() processed_sentence = [] for sentence in corpus: #將所有字詞轉為小寫 sentence = sentence.lower() 10 #移除標點符號(可以依據使用狀況決定是否要移除標點符號) 11 12 pattern = r'[^\W_]+' sentence = re.findall(pattern, sentence) 13 14 #添加字詞到字典中 15 word_dic |= set(sentence) 16 17 processed_sentence.append(sentence) 19 processed_sentence [['you', 'say', 'goodbye', 'and', 'i', 'say', 'hello']] 1 #定義簡易文本資料(使用Ch17講義中的例子) 複製程式碼 corpus = ['You say goodbye and I say hello.'] 4 word_dic = set() processed_sentence = [] for sentence in corpus: #將所有字詞轉為小寫 sentence = sentence.lower() 11 #移除標點符號(可以依據使用狀況決定是否要移除標點符號) 12 pattern = r'[^\W_]+' 建立 ID 字典清單並將文本轉換為 ID 形式: 1 #建立字詞ID清單 word2idx = dict() 3 idx2word = dict() 4 for word in word_dic: if word not in word2idx: idx = len(word2idx) word2idx[word] = idx idx2word[idx] = word 10 #將文本轉為ID型式 11 id_mapping = lambda x: word2idx[x] corpus = np.array([list(map(id_mapping, sentence)) for sentence in processed_sentence]) 14 print(corpus) 15 print(word2idx)
16 print(idx2word) 文本的 [[3 5 1 0 4 5 2]] {'and': 0, 'goodbye': 1, 'hello': 2, 'you': 3, 'i': 4, 'say': 5} {0: 'and', 1: 'goodbye', 2: 'hello', 3: 'you', 4: 'i', 5: 'say'} ——— ID 字典清單 1 #建立字詞ID清單 複製程式碼 word2idx = dict() 3 idx2word = dict() 4 for word in word_dic: if word not in word2idx: idx = len(word2idx)word2idx[word] = idx idx2word[idx] = word 10 #將文本轉為ID型式 id_mapping = lambda x: word2idx[x] corpus = np.array([list(map(id_mapping, sentence)) for sentence in processed_sentence]) 共現矩陣 將轉化處理過的文本資料轉化為共現矩陣 (字詞向量化) 以 window 為 1 建立共現矩陣: 1 #定義共現矩陣函式 vocab_size = len(word2idx) 3 window_size = 1 5 # initialize co-occurrence matrix co_matrix = np.zeros(shape=(vocab_size, vocab_size), dtype=np.int32) 8 for sentence in corpus: sentence_size = len(sentence) 10 for idx, word_id in enumerate(sentence): 11 12 for i in range(1, window_size+1): left_idx = idx - i 13 right_idx = idx + i 14 15 if left idx >= 0: 16 17 left_word_id = sentence[left_idx] 18 co_matrix[word_id, left_word_id] += 1 19 if right_idx < sentence_size:</pre> 20 21 right_word_id = sentence[right_idx] co_matrix[word_id, right_word_id] += 1 22 23 co_matrix array([[0, 1, 0, 0, 1, 0], [1, 0, 0, 0, 0, 1], [0, 0, 0, 0, 0, 1], ▶ 文本共現矩陣 [0, 0, 0, 0, 0, 1], [1, 0, 0, 0, 0, 1], [0, 1, 1, 1, 1, 0]], dtype=int32) 1 #定義共現矩陣函式 vocab_size = len(word2idx) 3 window_size = 1 # initialize co-occurrence matrix co_matrix = np.zeros(shape=(vocab_size, vocab_size), dtype=np.int32) for sentence in corpus: sentence_size = len(sentence) for idx, word_id in enumerate(sentence): for i in range(1, window_size+1): 詞向量相似度 當將字詞轉化為向量後,時常會需要比較字詞(如字詞為相似詞或反義詞),而比較轉換為向量的字詞的 方法有很多種,其中當要表示字詞的相似度時,最常使用的方法為餘弦相似度 (Cosine Similarity)。 計算 'i' 與 'you' 的字詞相似度: 1 # 定義餘弦相似度 def cos_similarity(x: np.ndarray, y: np.ndarray, eps: float=1e-8): nx = x / (np.sqrt(np.sum(x**2)) + eps)ny = y / (np.sqrt(np.sum(y**2)) + eps)return np.dot(nx,ny) 8 # calculate the similarity between I and you 9 cos_similarity(co_matrix[word2idx['i']], co_matrix[word2idx['you']]) 0.7071067691154799 字詞相似度 1 # 定義餘弦相似度 def cos_similarity(x: np.ndarray, y: np.ndarray, eps: float=1e-8): nx = x / (np.sqrt(np.sum(x**2)) + eps)ny = y / (np.sqrt(np.sum(y**2)) + eps)return np.dot(nx,ny) 8 # calculate the similarity between I and you cos_similarity(co_matrix[word2idx['i']], co_matrix[word2idx['you']]) 正向點間互資訊(PPMI) 由於共生矩陣在高頻字上的缺陷,而 PMI 中加入了兩字詞共同出現的機率與各自出現的機率的關係, 以此解決高頻詞在共生矩陣上的缺陷。而 PPMI 即將 PMI 內會產生負值的情況排除(若出現負值則賦予 0) ° 基於共現矩陣建立 PPMI 矩陣: 1 # 建立初始PPMI矩陣 2 M = np.zeros_like(co_matrix, dtype=np.float32) 3 # 取得所有字詞語個別字詞出現次數 4 N = np.sum(co_matrix) 5 S = np.sum(co_matrix, axis=0) 6 total = co_matrix.shape[0]*co_matrix.shape[1] 8 for i in range(co_matrix.shape[0]): 9 for j in range(co_matrix.shape[1]): pmi = np.log2(co_matrix[i, j]*N / (S[i]*S[j])) 11 M[i, j] = max(0, pmi)12 M /Users/admin/Documents/Coding/Python/Deep_Learning_Implement/Pytorch_Implement ckages/ipykernel_launcher.py:10: RuntimeWarning: divide by zero encountered i # Remove the CWD from sys.path while we load stuff. array([[0. , 1.5849625, 0. , 1.5849625, 0.], [1.5849625, 0. , 0. , 0. , 0. , 0.5849625], [0. , 0. , 0. , 0. , 0. , 1.5849625], PPMI (全為大 ← [0. , 0. , 0. , 0. , 0. , 1.5849625], [1.5849625, 0. , 0. , 0. , 0. , 0. , 0.5849625], 於等於0的值) [0. , 0.5849625, 1.5849625, 1.5849625, 0.5849625, 0.]], dtype=float32) 2 N = np.sum(co_matrix) 3 S = np.sum(co_matrix, axis=0) 4 total = co_matrix.shape[0]*co_matrix.shape[1] 6 for i in range(co_matrix.shape[0]): 7 for j in range(co_matrix.shape[1]): pmi = np.log2(co_matrix[i, j]*N / (S[i]*S[j])) M[i, j] = max(0, pmi)10 11 print(M) 奇異值分解(SVD) 對於稀疏矩陣因大部分的元素都為 0, 此矩陣包含了許多無法提供訊息的元素, 可利用奇異值分解, 一 方便將矩陣降維(減少計算複雜度),並保存重要的資訊。 使用 SVD 進行降維: 1 # 使用np的linalg.svd對PPMI矩陣進行奇異值分解 3 # SVD 4 U, S, V = np.linalg.svd(output_ppmi) 6 # 使用SVD將將原本的稀疏向量轉變為稠密向量 7 print(f'hello in co-occurrence matrix: {co_matrix[0]}') 8 print(f"hello in PPMI: {output_ppmi[0]}") 9 print(f"hello in SVD: {U[0]}") hello in co-occurrence matrix: [0 1 0 0 1 0 0] hello in PPMI: [0. 1.8073549 0. 0. 1.8073549 0. 0. hello in SVD: [-4.9782813e-01 3.3306691e-16 -3.3306691e-16 -6.8039632e-01 5.3779924e-01 0.0000000e+00 1.1102230e-16] 1 # 使用np的linalg.svd對PPMI矩陣進行奇異值分解 3 # SVD 4 U, S, V = np.linalg.svd(output_ppmi) 6 # 使用SVD將將原本的稀疏向量轉變為稠密向量 print(f'hello in co-occurrence matrix: {co_matrix[0]}') 8 print(f"hello in PPMI: {output_ppmi[0]}") print(f"hello in SVD: {U[0]}") 詳細使用操作 請參照使用技術方法詞向量實作.ipynb 檔進行更詳細的使用操作 知識點回顧

在這章節我們學習到 • 實作計數方法詞向量 • 使用 SVD 對詞向量進行降維

延伸閱讀

網站: 矩陣分解 針對大量資料的矩陣分解介紹

數學、人工智慧與蟒蛇 ML/DL 數學篇 ML/DL 應用篇



資料科學中的數據表達方式,大多是以矩陣的形式呈現的,舉個最簡單

 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_n x_n = X\beta$

 $[y_1]$ $[\beta_{01}]$ $[x_{11}$... $x_{1n}][\beta_1]$

的例子,我們考慮一個簡單線性迴歸方程式有如下表達方式:

下一步:閱讀範例與完成作業