

Assignment_9(Decision Tree on Donors choose dataset)

1 message

Applied AI Course <team@appliedaicourse.com>

Sat, 17 Aug, 2019 at 00:47

To: SuGuru <sugurunaresh111@gmail.com>

Hi SuGuru,

Please go through this reference vectorization (https://colab.research.google.com/drive/1jBdNJdyO47mXt505PXM_RAHizGq8muLe)

```
1 # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
2 tfidf_model = TfidfVectorizer()
3 tfidf_model.fit(preprocessed_essays)
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
6 tfidf_words = set(tfidf_model.get_feature_names())

1 # average Word2Vec
2 # compute average word2vec for each review.
3 tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if (word in glove_words) and (word in tfidf_words):
9             vec = model[word] # getting the vector for each word
10            # here we are multiplying idf value(dictionary[word]) and the tf value
11            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
12            vector += (vec * tf_idf) # calculating tfidf weighted w2v
13            tf_idf_weight += tf_idf
14        if tf_idf_weight != 0:
15            vector /= tf_idf_weight
16        tfidf_w2v_vectors.append(vector)
17
18 print(len(tfidf_w2v_vectors))
19 print(len(tfidf_w2v_vectors[0]))
```

100% 5000

Please fit your tf-idf vectors only on train data.

Thank you