

# Introduction

Natural Language Inferencing (NLI) task is one of the most important subsets of Natural Language Processing (NLP). NLI is also known as Recognizing Textual Entailment (RTE). NLI is a task of determining whether the given "hypothesis" and "premise" logically follow (entailment) or unfollow (contradiction) or are undetermined (neutral) to each other. In other words, the task of NLI is also defined as determining whether a "hypothesis" is true (entailment), false (contradiction), or undetermined (neutral) given a "premise". Table 1 gives some examples of these three cases.

Text (Premise)	Hypothesis	Judgement (Label)
A soccer game with multiple males playing. ---	Some men are playing a sport. ---	entailment ---
A black race car starts up in front of a crowd of people. ---	A man is driving down a lonely road. ---	contradiction ---
An older and a younger man smiling. ---	Two men are smiling and laughing at the cats playing on the floor. ---	neutral ---

## Reading and preparing data

In [328]:

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import os
6 import re
7 from nltk.corpus import stopwords
8 import pickle
9 from tqdm import tqdm
10 from wordcloud import WordCloud
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 import tensorflow as tf
13 import random
14 import math
15 from textblob import TextBlob
16 from collections import Counter
17 from sklearn.metrics.pairwise import cosine_similarity
18 from sklearn.feature_extraction.text import TfidfVectorizer
19 import nltk, string
20 from collections import Counter
21 from nltk.stem import WordNetLemmatizer
22 import Levenshtein
23 %matplotlib inline

```

In [2]: 1 nltk.download('punkt')

```
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\anike\AppData\Roaming\nltk_data...
[nltk_data]     Package punkt is already up-to-date!
```

Out[2]: True

In [3]: 1 # Setting seed  
2 global\_seed = 42  
3 np.random.seed(42)  
4 random.seed(42)  
5 tf.random.set\_seed(42)

In [3]: 1 train\_df = pd.read\_csv('data/snli\_1.0/snli\_1.0\_train.txt', sep='\t')  
2 val\_df = pd.read\_csv('data/snli\_1.0/snli\_1.0\_dev.txt', sep='\t')  
3 test\_df = pd.read\_csv('data/snli\_1.0/snli\_1.0\_test.txt', sep='\t')

In [4]: 1 print(f"train data shape: {train\_df.shape}\nvalidation data shape:{val\_df.sh  
train data shape: (550152, 14)  
validation data shape:(10000, 14)  
test data shape:(10000, 14)

In [5]: 1 print(f"all the columns in the dataset:{list(train\_df.columns)}")

```
all the columns in the dataset:['gold_label', 'sentence1_binary_parse', 'sentence2_binary_parse', 'sentence1_parse', 'sentence2_parse', 'sentence1', 'sentence2', 'captionID', 'pairID', 'label1', 'label2', 'label3', 'label4', 'label5']
```

**Here sentence1 represents premise, sentence2 represents hypothesis and gold\_label represents our target variable i.e. whether hypothesis entails or contradicts the premise. neutral gold\_label suggests no relationship between premise and hypothesis**

In [6]: 1 train\_df.head(5)

	gold_label	sentence1_binary_parse	sentence2_binary_parse	sentence1_parse	sentence2_parse
0	neutral	(( A person ) ( on ( a horse ) ) ( jump...))	(( A person ) ( ( is ( training ( his horse...)))	(ROOT (S (NP (NP (DT A) (NN person)) (PP (IN o...)))	(ROOT (S (NP (DT A) (NN person)) (VP (VBZ is) ...
1	contradiction	(( A person ) ( on ( a horse ) ) ( jump...))	(( A person ) ( ( ( is ( at ( a diner ) ) ...))	(ROOT (S (NP (NP (DT A) (NN person)) (PP (IN o...)))	(ROOT (S (NP (DT A) (NN person)) (VP (VBZ is) ...
2	entailment	(( A person ) ( on ( a horse ) ) ( jump...))	(( A person ) ( ( ( is outdoors ), ) ( on ...))	(ROOT (S (NP (NP (DT A) (NN person)) (PP (IN o...)))	(ROOT (S (NP (DT A) (NN person)) (VP (VBZ is) ...
3	neutral	( Children ( ( ( smiling and ) waving ) ( at c...))	( They ( are ( smiling ( at ( their parents ) ...))	(ROOT (NP (S (NP (NNP Children)) (VP (VBP smil...))	(ROOT (S (NP (PRP They)) (VP (VBP are) (VP (VB...))
4	entailment	( Children ( ( ( smiling and ) waving ) ( at c...))	( There ( ( are children ) present ))	(ROOT (NP (S (NP (NNP Children)) (VP (VBP smil...))	(ROOT (S (NP (EX There)) (VP (VBP are) (NP (NN...))

◀ ▶

**we only need sentence1, sentence2 and gold\_label columns. Therefore let us drop rest of the columns**

In [7]:

```
1 to_keep = ['sentence1', 'sentence2', 'gold_label']
2 train_df = train_df[to_keep]
3 val_df = val_df[to_keep]
4 test_df = test_df[to_keep]
```

In [8]: 1 train\_df.head(5)

	sentence1	sentence2	gold_label
0	A person on a horse jumps over a broken down a...	A person is training his horse for a competition.	neutral
1	A person on a horse jumps over a broken down a...	A person is at a diner, ordering an omelette.	contradiction
2	A person on a horse jumps over a broken down a...	A person is outdoors, on a horse.	entailment
3	Children smiling and waving at camera	They are smiling at their parents	neutral
4	Children smiling and waving at camera	There are children present	entailment

## Checking and dealing with null values

In [9]: 1 train\_df[train\_df.isna().any(axis=1)]

Out[9]:

		sentence1	sentence2	gold_label
91479	Cannot see picture to describe.	NaN	neutral	
91480	Cannot see picture to describe.	NaN	entailment	
91481	Cannot see picture to describe.	NaN	contradiction	
311124	Jumping with purple balls is so much fun!	NaN	contradiction	
311125	Jumping with purple balls is so much fun!	NaN	neutral	
311126	Jumping with purple balls is so much fun!	NaN	entailment	

In [10]: 1 val\_df[val\_df.isna().any(axis=1)]

Out[10]: sentence1 sentence2 gold\_label

In [11]: 1 test\_df[test\_df.isna().any(axis=1)]

Out[11]: sentence1 sentence2 gold\_label

**6 rows in train data have NaN values. It is better to drop these rows because we can not find appropriate sentences here and inserting some sentences on our own based on target would be inappropriate**

In [12]: 1 train\_df = train\_df.drop(train\_df.index[train\_df.isna().any(axis=1)]

In [13]: 1 train\_df['gold\_label'].unique()

Out[13]: array(['neutral', 'contradiction', 'entailment', '-'], dtype=object)

In [14]: 1 train\_df.groupby(by = 'gold\_label').count()

Out[14]: sentence1 sentence2

	gold_label	
-	785	785
contradiction	183185	183185
entailment	183414	183414
neutral	182762	182762

In [15]: 1 train\_df.shape

Out[15]: (550146, 3)

```
In [16]: 1 train_df[train_df['gold_label'].values != '-'].shape
```

```
Out[16]: (549361, 3)
```

we can see that there are some sentence pairs for which the label is '-' i.e. it is undetermined. Out of 550152 such pairs 785 pairs have undetermined label. Since target variable is unknown for these pairs and number of such pairs is too low compared to the number of pairs for which target label is present, we can drop the pairs where label is undetermined

```
In [17]: 1 train_df = train_df[train_df['gold_label'].values != '-']
2 val_df = val_df[val_df['gold_label'].values != '-']
3 test_df = test_df[test_df['gold_label'].values != '-']
```

```
In [18]: 1 train_df['gold_label'].unique()
```

```
Out[18]: array(['neutral', 'contradiction', 'entailment'], dtype=object)
```

```
In [19]: 1 train_df.groupby(by = 'gold_label').count()
```

```
Out[19]:
```

	sentence1	sentence2
<b>gold_label</b>		
<b>contradiction</b>	183185	183185
<b>entailment</b>	183414	183414
<b>neutral</b>	182762	182762

gold_label	sentence1	sentence2
<b>contradiction</b>	183185	183185
<b>entailment</b>	183414	183414
<b>neutral</b>	182762	182762

```
In [ ]:
```

Since this is a standard benchmark dataset, there is almost same number of training examples for all target classes. So there is no imbalance in the dataset and hence accuracy can be taken as key performance index

Let us look at some random pair of sentences

In [20]:

```
1 for i,row in train_df.sample(15).iterrows():
2     print(f"index: {i}\nsentence1: {row['sentence1']}\nsentence2: {row['sent
3     print('*'*100)

index: 71202
sentence1: A woman plays a violin outdoors.
sentence2: The ball room dancer slipped on a banana peel.
label: contradiction

*****
index: 178279
sentence1: The red panted cyclist is amongst nature.
sentence2: The cyclist is outdoors.
label: entailment

*****
index: 478971
sentence1: A bicyclist is doing a trick in midair.
sentence2: The bicycle is slowly rolling down the straight.
label: contradiction

*****
index: 7164
sentence1: Two motorcyclists racing neck and neck around a corner.
sentence2: The two motorcyclists are racing each other.
label: entailment

*****
index: 237603
sentence1: A middle-age man in black suit speaking into the mic, behind brown podium.
sentence2: a guy is dancing on a table
label: contradiction

*****
index: 351768
sentence1: Young women playing a intense game of lacrosse.
sentence2: The girls are playing hard at lacrosse.
label: entailment

*****
index: 329609
sentence1: Young boy in a brown shirt doing a back flip
sentence2: The boy does a front flip then a back flip.
label: neutral

*****
index: 117774
sentence1: A smoking man in a green shirt throws a ball as others watch.
sentence2: The mans yellow shirt reflected the sun and made him miss the catch.
```

label: contradiction

\*\*\*\*\*  
\*\*\*\*\*

index: 172126

sentence1: A little boy wipes his nose with his blue shirt.

sentence2: The kid licks the snot from his nose.

label: contradiction

\*\*\*\*\*  
\*\*\*\*\*

index: 185666

sentence1: Picture is of the inside of an abandoned building and has name graffiti on it.

sentence2: The photo is of a cat in the grass.

label: contradiction

\*\*\*\*\*  
\*\*\*\*\*

index: 320636

sentence1: a person in pink pants playing a guitar

sentence2: A guy is walking his dog.

label: contradiction

\*\*\*\*\*  
\*\*\*\*\*

index: 531519

sentence1: A man wearing a bright neon orange vest sitting by a parked truck with his laptop.

sentence2: A man wearing a bright neon orange vest sitting by a parked truck with his laptop is waiting for his friend.

label: neutral

\*\*\*\*\*  
\*\*\*\*\*

index: 107301

sentence1: A fisherman drags a red net through the water close to the beach.

sentence2: A fisherman's net is nearly full as he drags it through the water.

label: neutral

\*\*\*\*\*  
\*\*\*\*\*

index: 116716

sentence1: Townspeople in a downtown market selling goods.

sentence2: the people are at a market

label: entailment

\*\*\*\*\*  
\*\*\*\*\*

index: 277120

sentence1: A woman with glasses is sitting in front of her laptop in a library.

sentence2: A woman sits at a computer looking at porn.

label: neutral

\*\*\*\*\*  
\*\*\*\*\*

Since it is a human labelled data prepared by researchers at stanform using amazon mechanical turk crowdsourcing service, the data is pretty clean. We only need to remove punctuation and numerical values if present in the data, make all character lower case.

## Cleaning and preprocessing text

Along with classification task, once we get good performance on classification task, we will also try to use the trained model to generate universal sentence embeddings. Therefore keeping all the words in your corpora in necessary. For this reason we will not remove stopwords

```
In [21]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"\n\t", " not", phrase)
11    phrase = re.sub(r"\re", " are", phrase)
12    phrase = re.sub(r"\s", " is", phrase)
13    phrase = re.sub(r"\d", " would", phrase)
14    phrase = re.sub(r"\ll", " will", phrase)
15    phrase = re.sub(r"\t", " not", phrase)
16    phrase = re.sub(r"\ve", " have", phrase)
17    phrase = re.sub(r"\m", " am", phrase)
18    return phrase
```

```
In [22]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentence in tqdm(text_data):
7         if type(sentence) != str:
8             sent = ''
9         else:
10             sent = decontracted(sentence)
11             sent = sent.replace('\r', ' ')
12             sent = sent.replace('\n', ' ')
13             sent = sent.replace('"', ' ')
14             sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
15             # https://gist.github.com/sebleier/554280
16             sent = ' '.join(e for e in sent.split())
17             preprocessed_text.append(sent.lower().strip())
18     return preprocessed_text
```

Type *Markdown* and *LaTeX*:  $\alpha^2$

In [23]:

```
1 train_df['sentence1_preprocessed'] = preprocess_text(train_df['sentence1'].v
2 val_df['sentence1_preprocessed'] = preprocess_text(val_df['sentence1'].value
3 test_df['sentence1_preprocessed'] = preprocess_text(test_df['sentence1'].val
4 train_df['sentence2_preprocessed'] = preprocess_text(train_df['sentence2'].v
5 val_df['sentence2_preprocessed'] = preprocess_text(val_df['sentence2'].value
6 test_df['sentence2_preprocessed'] = preprocess_text(test_df['sentence2'].val
```

```
100%|██████████| 5
49361/549361 [00:08<00:00, 64610.90it/s]
100%|██████████| 5
█| 9842/9842 [00:00<00:00, 69286.03it/s]
100%|██████████| 5
█| 9824/9824 [00:00<00:00, 66364.80it/s]
100%|██████████| 5
49361/549361 [00:06<00:00, 88990.73it/s]
100%|██████████| 5
█| 9842/9842 [00:00<00:00, 84092.00it/s]
100%|██████████| 5
█| 9824/9824 [00:00<00:00, 84667.25it/s]
```

In [24]:

```
1 for i,row in train_df.sample(5).iterrows():
2     print(f"index: {i}\nsentence1: {row['sentence1']}\npreprocessed sentence
3     {row['sentence1_preprocessed']}\\nsentence2: {row['sentence2']}\npreproce
4     {row['sentence2_preprocessed']}\\nlabel: {row['gold_label']}\\n")
5     print('*'*100)
```

```
index: 172202
sentence1: Dancers in white are performing on a blacked out stage with a pink
lighted hue illuminating them.
preprocessed sentence1:      dancers in white are performing on a blacked out
stage with a pink lighted hue illuminating them
sentence2: Dancers are dancing.
preprocessed sentence2:      dancers are dancing
label: entailment

*****
*****
index: 222740
sentence1: A man in a black jacket is riding a horse on a public sidewalk.
preprocessed sentence1:      a man in a black jacket is riding a horse on a pu
blic sidewalk
sentence2: A man is in a parade.
preprocessed sentence2:      a man is in a parade
label: neutral

*****
*****
index: 75709
sentence1: A baseball pitcher wearing a blue, red and white uniform is about
to throw a baseball.
preprocessed sentence1:      a baseball pitcher wearing a blue red and white u
niform is about to throw a baseball
sentence2: A baseball player prepares to throw a ball.
preprocessed sentence2:      a baseball player prepares to throw a ball
label: entailment

*****
*****
index: 226118
sentence1: Two females walking down the street with shopping bags.
preprocessed sentence1:      two females walking down the street with shopping
bags
sentence2: Two women are holding bags.
preprocessed sentence2:      two women are holding bags
label: entailment

*****
*****
index: 526529
sentence1: The snowboarder is in the middle of a very tall jump
preprocessed sentence1:      the snowboarder is in the middle of a very tall j
ump
sentence2: Snowboarder competing in jump competition.
preprocessed sentence2:      snowboarder competing in jump competition
label: neutral
```

```
*****  
*****
```



In [25]:

```

1  for i,row in val_df.sample(5).iterrows():
2      print(f"index: {i}\nsentence1: {row['sentence1']}\npreprocessed sentence
3      {row['sentence1_preprocessed']} }\nsentence2: {row['sentence2']}\npreproce
4      {row['sentence2_preprocessed']} }\nlabel: {row['gold_label']}\n")
5      print('*'*100)

```

```

index: 231
sentence1: A middle-aged man is sitting indian style outside holding a folded p
aper in his hands.
preprocessed sentence1: a middle aged man is sitting indian style outside h
olding a folded paper in his hands
sentence2: A middle aged man is outdoors.
preprocessed sentence2: a middle aged man is outdoors
label: entailment

*****
*****index: 9542
sentence1: little boy running in the street.
preprocessed sentence1: little boy running in the street
sentence2: A boy is running to school in the street.
preprocessed sentence2: a boy is running to school in the street
label: neutral

*****
*****index: 4964
sentence1: A person is walking down a cobblestone street.
preprocessed sentence1: a person is walking down a cobblestone street
sentence2: The person walking has shoes on.
preprocessed sentence2: the person walking has shoes on
label: neutral

*****
*****index: 9604
sentence1: A man swims in the pool with a child in a life jacket.
preprocessed sentence1: a man swims in the pool with a child in a life jack
et
sentence2: The man is in the water.
preprocessed sentence2: the man is in the water
label: entailment

*****
*****index: 3241
sentence1: a man with a tattoo on his arm cooking something in a frying pan.
preprocessed sentence1: a man with a tattoo on his arm cooking something in
a frying pan
sentence2: A man cooks food.
preprocessed sentence2: a man cooks food
label: entailment

*****
*****
```

In [26]:

```

1  for i,row in test_df.sample(5).iterrows():
2      print(f"index: {i}\nsentence1: {row['sentence1']}\npreprocessed sentence
3      {row['sentence1_preprocessed']} }\nsentence2: {row['sentence2']}\npreproce
4      {row['sentence2_preprocessed']} }\nlabel: {row['gold_label']}\n")
5      print('*'*100)

```

```

index: 3109
sentence1: A man sits at a table in a room.
preprocessed sentence1: a man sits at a table in a room
sentence2: A woman sits.
preprocessed sentence2: a woman sits
label: contradiction

*****
index: 82
sentence1: Three people sit on a bench at a station, the man looks oddly at the
two women, the redhead women looks up and forward in an awkward position, and
the yellow blond girl twiddles with her hair.
preprocessed sentence1: three people sit on a bench at a station the man lo
oks oddly at the two women the redhead women looks up and forward in an awkwa
rd position and the yellow blond girl twiddles with her hair
sentence2: People run together.
preprocessed sentence2: people run together
label: contradiction

*****
index: 9089
sentence1: Two boys in a canoe in front of a larger boat with a storm coming i
n.
preprocessed sentence1: two boys in a canoe in front of a larger boat with
a storm coming in
sentence2: Two boys in a sailboat
preprocessed sentence2: two boys in a sailboat
label: contradiction

*****
index: 9378
sentence1: A city street corner with two passing bicyclists and several pedestri
ans.
preprocessed sentence1: a city street corner with two passing bicyclists an
d several pedestrians
sentence2: The traffic is busy.
preprocessed sentence2: the traffic is busy
label: neutral

*****
index: 36
sentence1: Three firefighter come out of subway station.
preprocessed sentence1: three firefighter come out of subway station
sentence2: Three firefighters putting out a fire inside of a subway station.
preprocessed sentence2: three firefighters putting out a fire inside of a s
ubway station

```

```
label: neutral
```

```
*****  
*****
```

```
In [27]: 1 train_df['gold_label'].value_counts()[0]
```

```
Out[27]: 183414
```

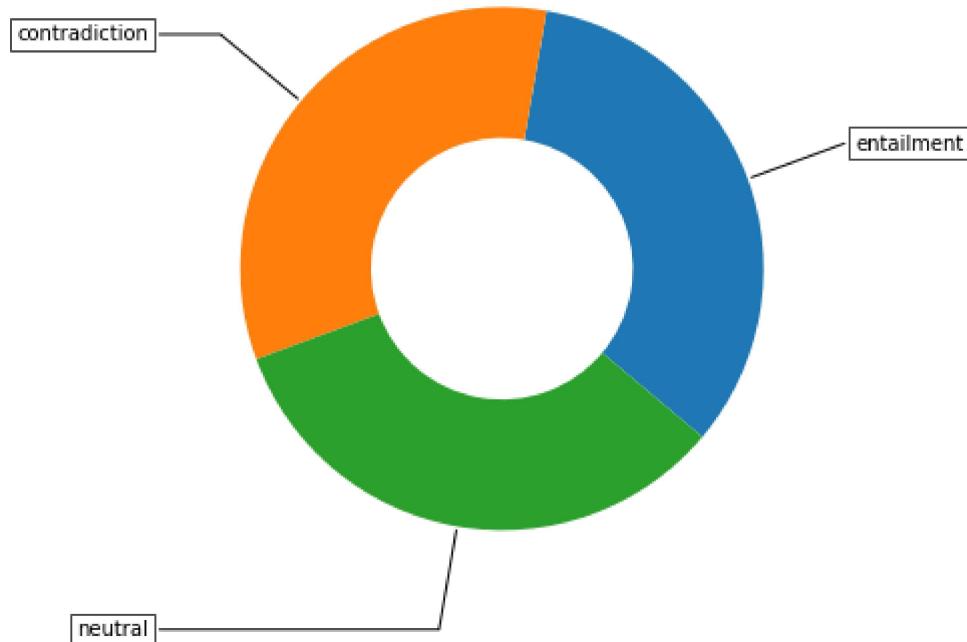
In [28]:

```

1 # this code is taken from
2 # https://matplotlib.org/gallery/pie_and_polar_charts/pie_and_donut_labels.h
3
4
5 y_value_counts = train_df['gold_label'].value_counts()
6 # print("Number of entailments ", y_value_counts[1], ", (", (y_value_counts[
7 # print("Number of projects that are not approved for funding ", y_value_cou
8
9 fig, ax = plt.subplots(figsize=(6, 6), subplot_kw=dict(aspect="equal"))
10 recipe = ["entailment", "contradiction", "neutral"]
11
12 data = [y_value_counts[0], y_value_counts[1], y_value_counts[2]]
13
14 wedges, texts = ax.pie(data, wedgeprops=dict(width=0.5), startangle=-40)
15
16 bbox_props = dict(boxstyle="square,pad=0.3", fc="w", ec="k", lw=0.72)
17 kw = dict(xycoords='data', textcoords='data', arrowprops=dict(arrowsstyle="-"
18             bbox=bbox_props, zorder=0, va="center"))
19
20 for i, p in enumerate(wedges):
21     ang = (p.theta2 - p.theta1)/2. + p.theta1
22     y = np.sin(np.deg2rad(ang))
23     x = np.cos(np.deg2rad(ang))
24     horizontalalignment = {-1: "right", 1: "left"}[int(np.sign(x))]
25     connectionstyle = "angle,angleA=0,angleB={}".format(ang)
26     kw["arrowprops"].update({"connectionstyle": connectionstyle})
27     ax.annotate(recipe[i], xy=(x, y), xytext=(1.35*np.sign(x), 1.4*y),
28                 horizontalalignment=horizontalalignment, **kw)
29
30 ax.set_title("distribution of label")
31
32 plt.show()

```

distribution of label



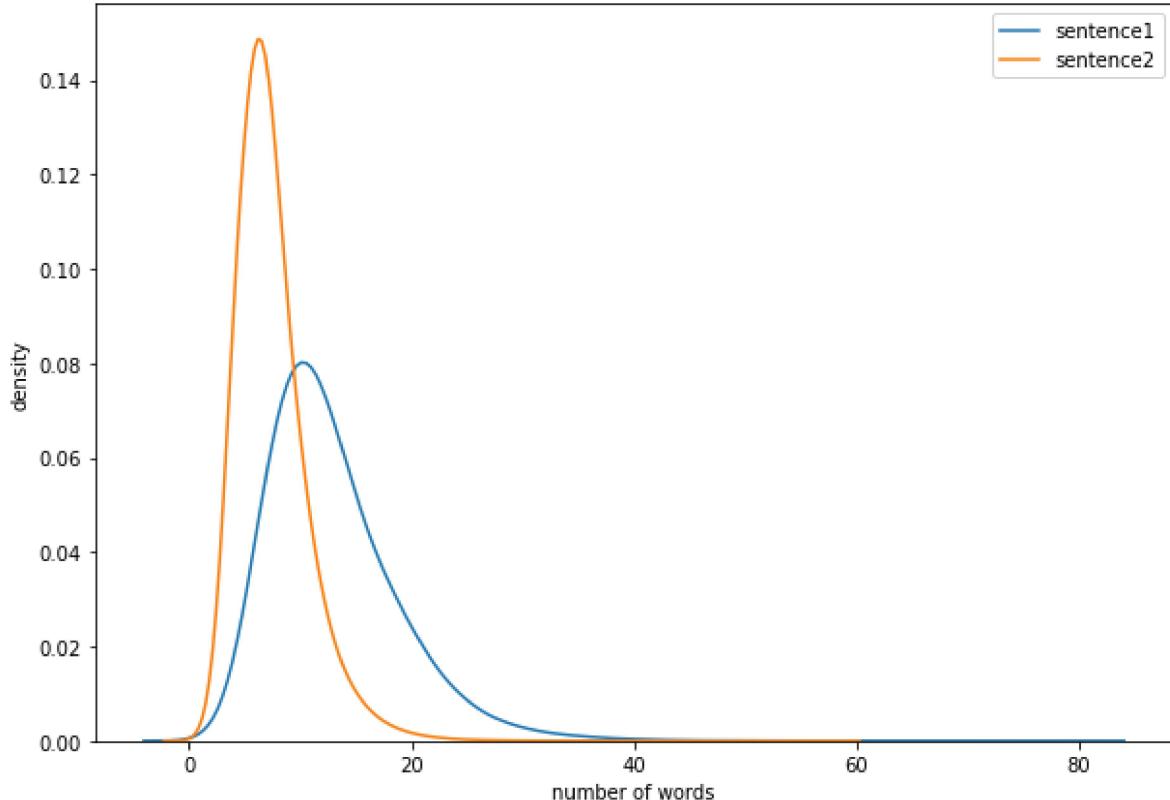
**As we saw earlier, the dataset is balanced from the point of view of target variable**

```
In [45]: 1 def get_num_words(sentence):
2     return(len(sentence.split(' ')))
```

```
In [46]: 1 train_df['sentence1_num_words'] = train_df['sentence1_preprocessed'].apply(get_n
2 val_df['sentence1_num_words'] = val_df['sentence1_preprocessed'].apply(get_n
3 test_df['sentence1_num_words'] = test_df['sentence1_preprocessed'].apply(get_n
4 train_df['sentence2_num_words'] = train_df['sentence2_preprocessed'].apply(get_n
5 val_df['sentence2_num_words'] = val_df['sentence2_preprocessed'].apply(get_n
6 test_df['sentence2_num_words'] = test_df['sentence2_preprocessed'].apply(get_n
```

```
In [47]: 1 fig,ax = plt.subplots(1,1, figsize=(10,7))
2 sns.kdeplot(x = train_df['sentence1_num_words'], ax = ax, label = "sentence1")
3 sns.kdeplot(x = train_df['sentence2_num_words'], ax = ax, label = "sentence2")
4 ax.set_xlabel('number of words')
5 ax.set_ylabel('density')
6 ax.set_title('distribution of count of words in sentence 1 and sentence 2',
7 plt.legend()
8 plt.show()
```

**distribution of count of words in sentence 1 and sentence 2**



**In general the number of words in premise(sentence1) is slightly greater than number of words in hypothesis(sentence2)**

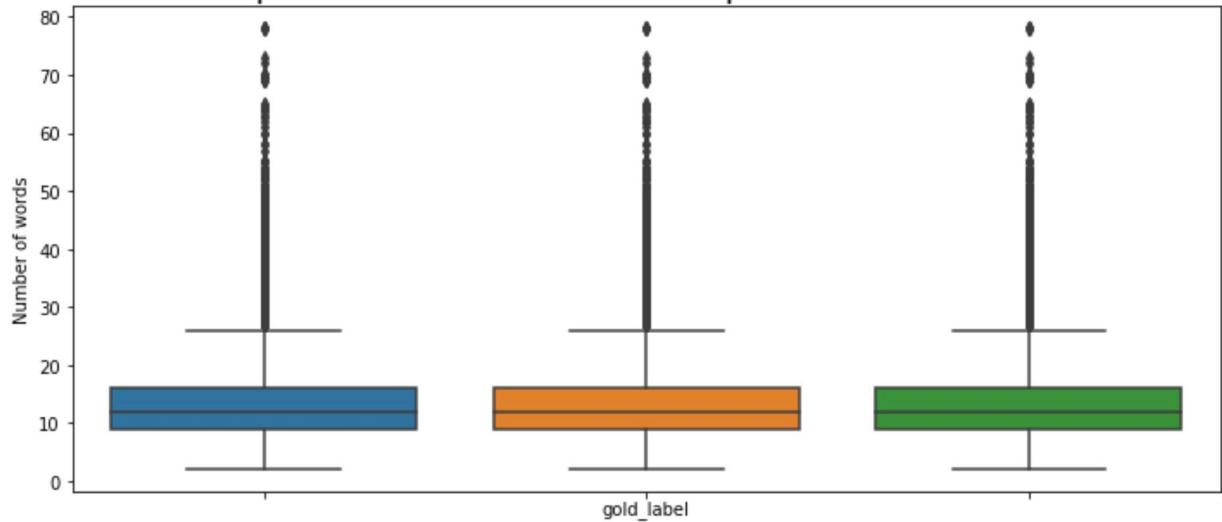
In [32]:

```

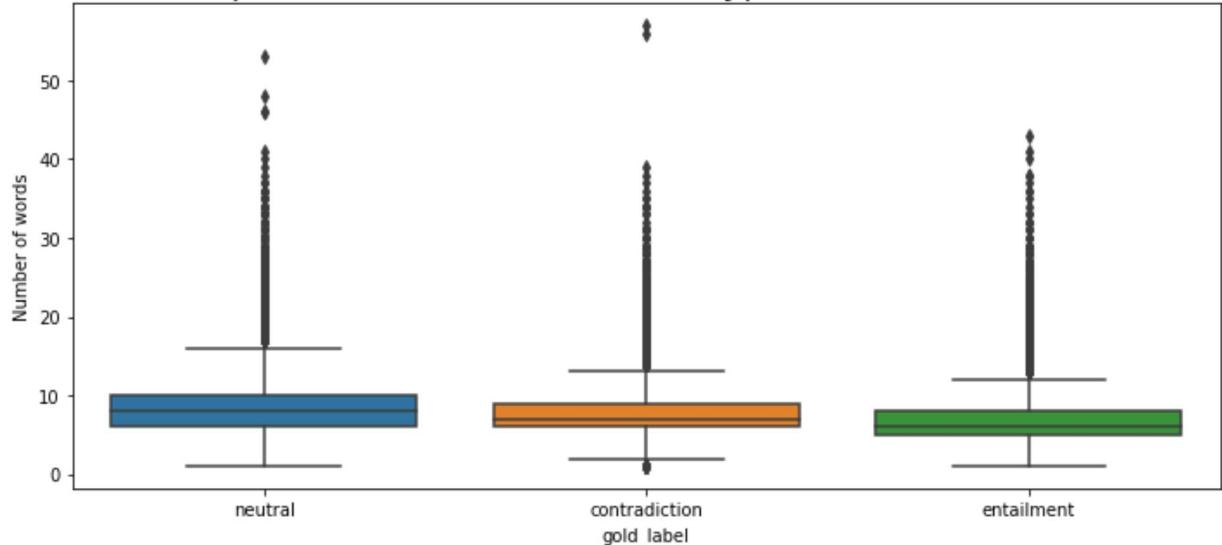
1 fig,ax = plt.subplots(2,1, figsize=(10,10), sharex=True)
2 sns.boxplot(x = train_df['gold_label'], y = train_df['sentence1_num_words'],
3 sns.boxplot(x = train_df['gold_label'], y = train_df['sentence2_num_words'],
4 ax[0].set_title('box plot of number of words in premise for each label', fontweight='bold')
5 ax[0].set_ylabel('Number of words')
6 ax[1].set_title('box plot of number of words in hypothesis for each label', fontweight='bold')
7 ax[1].set_ylabel('Number of words')
8 plt.tight_layout()
9 plt.subplots_adjust(hspace=0.3)
10 plt.show()

```

box plot of number of words in premise for each label



box plot of number of words in hypothesis for each label



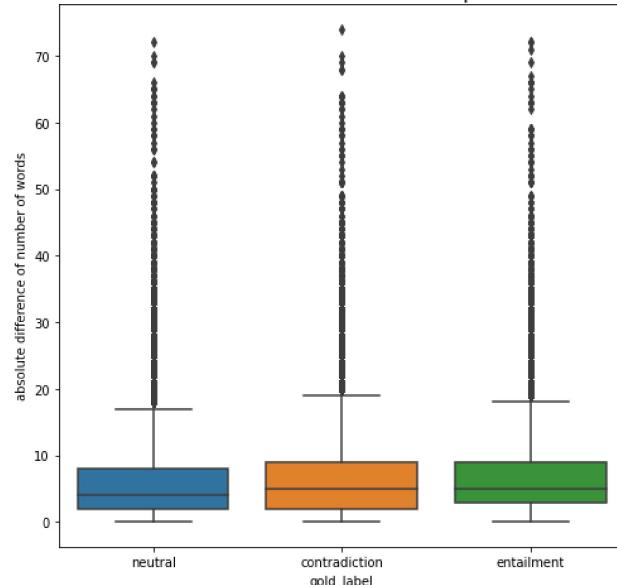
In general there is no pattern between number of words and label in case of premise. But there is a visible pattern in case of hypothesis. The number of words in hypothesis are highest for neutral, it decreases for contradiction and it is lowest for entailment

## Difference in num of words in sentence1 and sentence2

```
In [48]: 1 train_df['diff_num_words'] = (train_df['sentence1_num_words'] - train_df['se
```

```
In [49]: 1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['diff_num_words'], ax =
3 ax.set_title('box plot of absolute difference of num of words between premises')
4 ax.set_ylabel('absolute difference of number of words')
5 plt.tight_layout()
6 plt.show()
```

box plot of absolute difference of num of words between premise and hypothesis for each label



Although there is not much difference in general difference in number of words between pairs for neutral, contradiction and entailment sentences, it seems that the ascending order of difference in number of words between two sentences is neutral-->contradiction-->entailment

```
In [50]: 1 val_df['diff_num_words'] = (val_df['sentence1_num_words'] - val_df['sentence
2 test_df['diff_num_words'] = (test_df['sentence1_num_words'] - test_df['sente
```

```
In [ ]:
```

## Plotting Wordcloud for premise and hypothesis to see dominant words in these sentences

```
In [33]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df['sentence1_preprocessed']
3 X = vectorizer.fit_transform(data)
```

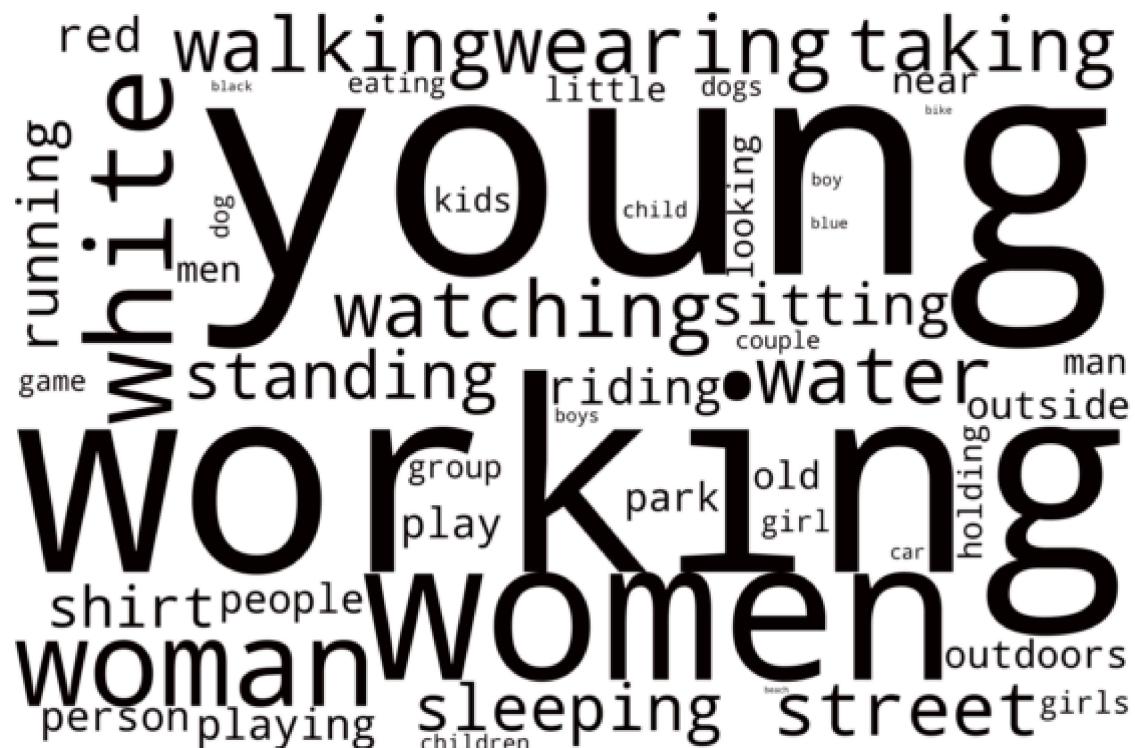
```
In [34]: 1 def black_color_func(word, font_size, position,orientation,random_state=None  
2     return("hsl(0,100%, 1%)")  
3 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_  
4 wordcloud.recolor(color_func = black_color_func)  
5 plt.figure(figsize=[10,10])  
6 plt.imshow(wordcloud, interpolation="bilinear")  
7 plt.axis("off")  
8 plt.show()
```



In premise words like yellow, young, women, woman, walking, talking, street, working, wearing, smiling etc. are dominant

```
In [35]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df['sentence2_preprocessed']
3 X = vectorizer.fit_transform(data)
```

```
In [36]: 1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
 2 wordcloud.recolor(color_func = black_color_func)
 3 plt.figure(figsize=[10,10])
 4 plt.imshow(wordcloud, interpolation="bilinear")
 5 plt.axis("off")
 6 plt.show()
```



**Hypothesis also contains words similar to premise**

## Wordcloud for premise for each label

- ## 1. premise--->entailment

```
In [37]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'entailment']['sentence1_preproces_
3 X = vectorizer.fit_transform(data)
```

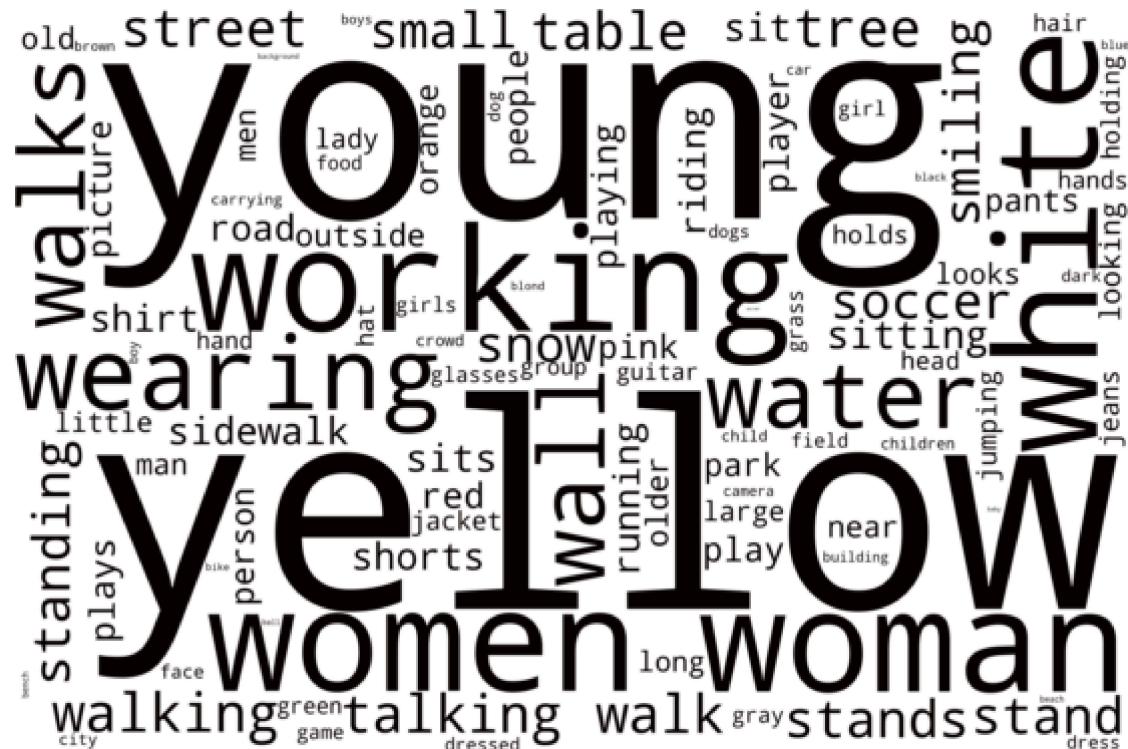
```
In [38]: 1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
2 wordcloud.recolor(color_func = black_color_func)
3 plt.figure(figsize=[10,10])
4 plt.imshow(wordcloud, interpolation="bilinear")
5 plt.axis("off")
6 plt.show()
```



## 2. premise--->contradiction

```
In [39]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'contradiction']['sentence1_prepro_
3 X = vectorizer.fit_transform(data)
```

```
In [40]: 1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
 2 wordcloud.recolor(color_func = black_color_func)
 3 plt.figure(figsize=[10,10])
 4 plt.imshow(wordcloud, interpolation="bilinear")
 5 plt.axis("off")
 6 plt.show()
```



### 3. premise--->neutral

```
In [41]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'neutral']['sentence1_preprocessed']
3 X = vectorizer.fit_transform(data)
```

```
In [42]: 1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
 2 wordcloud.recolor(color_func = black_color_func)
 3 plt.figure(figsize=[10,10])
 4 plt.imshow(wordcloud, interpolation="bilinear")
 5 plt.axis("off")
 6 plt.show()
```



## Wordcloud for hypothesis for each label

## 1. hypothesis--->entailment

```
In [43]: 1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'entailment']['sentence2_preproces
3 X = vectorizer.fit_transform(data)
```

In [44]:

```

1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
2 wordcloud.recolor(color_func = black_color_func)
3 plt.figure(figsize=[10,10])
4 plt.imshow(wordcloud, interpolation="bilinear")
5 plt.axis("off")
6 plt.show()

```



## 2. hypothesis--->contradiction

In [45]:

```

1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'contradiction']['sentence2_prepro
3 X = vectorizer.fit_transform(data)

```

In [46]:

```

1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
2 wordcloud.recolor(color_func = black_color_func)
3 plt.figure(figsize=[10,10])
4 plt.imshow(wordcloud, interpolation="bilinear")
5 plt.axis("off")
6 plt.show()

```



### 3. hypothesis--->neutral

In [47]:

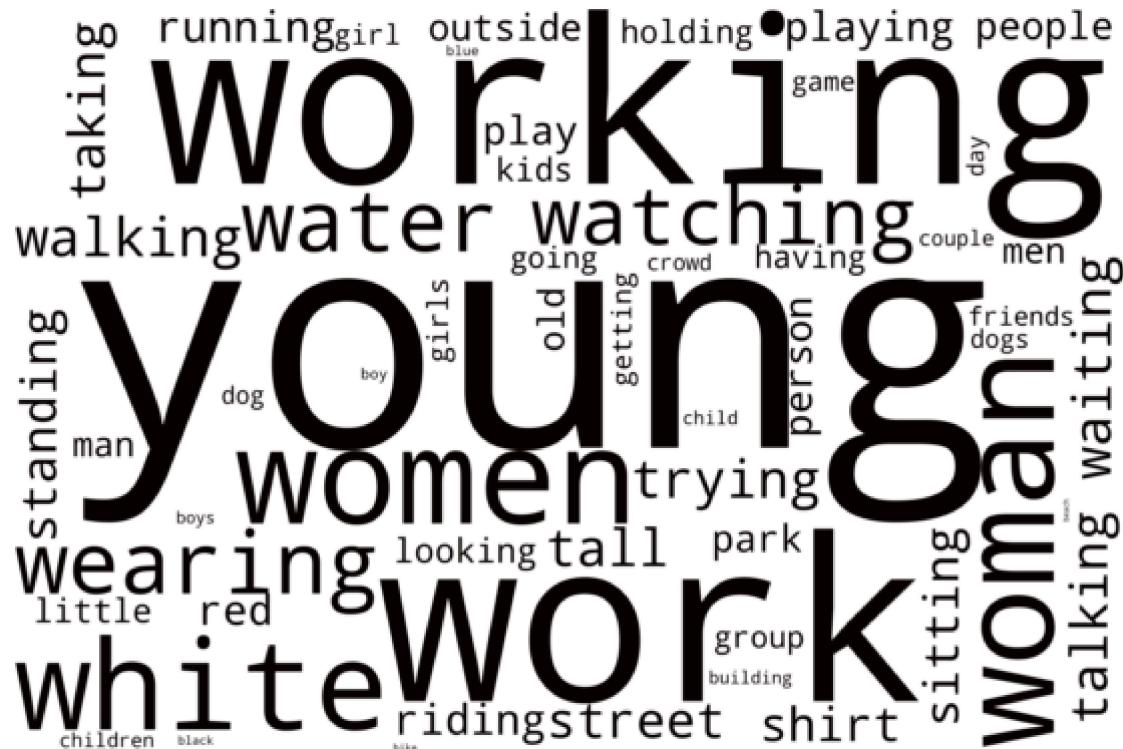
```

1 vectorizer = TfidfVectorizer(stop_words='english', ngram_range = (1,1), max_
2 data = train_df[train_df['gold_label'] == 'neutral']['sentence2_preprocessed']
3 X = vectorizer.fit_transform(data)

```

In [48]:

```
1 wordcloud = WordCloud(background_color="white", width=3000, height=2000, max_
2 wordcloud.recolor(color_func = black_color_func)
3 plt.figure(figsize=[10,10])
4 plt.imshow(wordcloud, interpolation="bilinear")
5 plt.axis("off")
6 plt.show()
```



It seems that dominant words in each label for each sentence do not give us much predictive power. We need to engineer more features

## Bleu score between two sentences

In [49]:

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 import nltk.translate.bleu_score as bleu
4 def get_bleu_score(sentence1, sentence2):
5     import nltk.translate.bleu_score as bleu
6     sentence1 = sentence1.split(' ')
7     sentence2 = sentence2.split(' ')
8     bleu_score = bleu.sentence_bleu([sentence1], sentence2)
9     return bleu_score

```

In [50]:

```

1 train_bleu_lst = list()
2 for i, row in train_df.iterrows():
3     sentence1 = row['sentence1_preprocessed']
4     sentence2 = row['sentence2_preprocessed']
5     bleu = get_bleu_score(sentence1, sentence2)
6     train_bleu_lst.append(bleu)

```

In [51]:

```

1 val_bleu_lst = list()
2 for i, row in val_df.iterrows():
3     sentence1 = row['sentence1_preprocessed']
4     sentence2 = row['sentence2_preprocessed']
5     bleu = get_bleu_score(sentence1, sentence2)
6     val_bleu_lst.append(bleu)

```

In [52]:

```

1 test_bleu_lst = list()
2 for i, row in test_df.iterrows():
3     sentence1 = row['sentence1_preprocessed']
4     sentence2 = row['sentence2_preprocessed']
5     bleu = get_bleu_score(sentence1, sentence2)
6     test_bleu_lst.append(bleu)

```

In [53]:

```
1 len(val_bleu_lst)
```

Out[53]: 9842

In [54]:

```
1 len(test_bleu_lst)
```

Out[54]: 9824

In [55]:

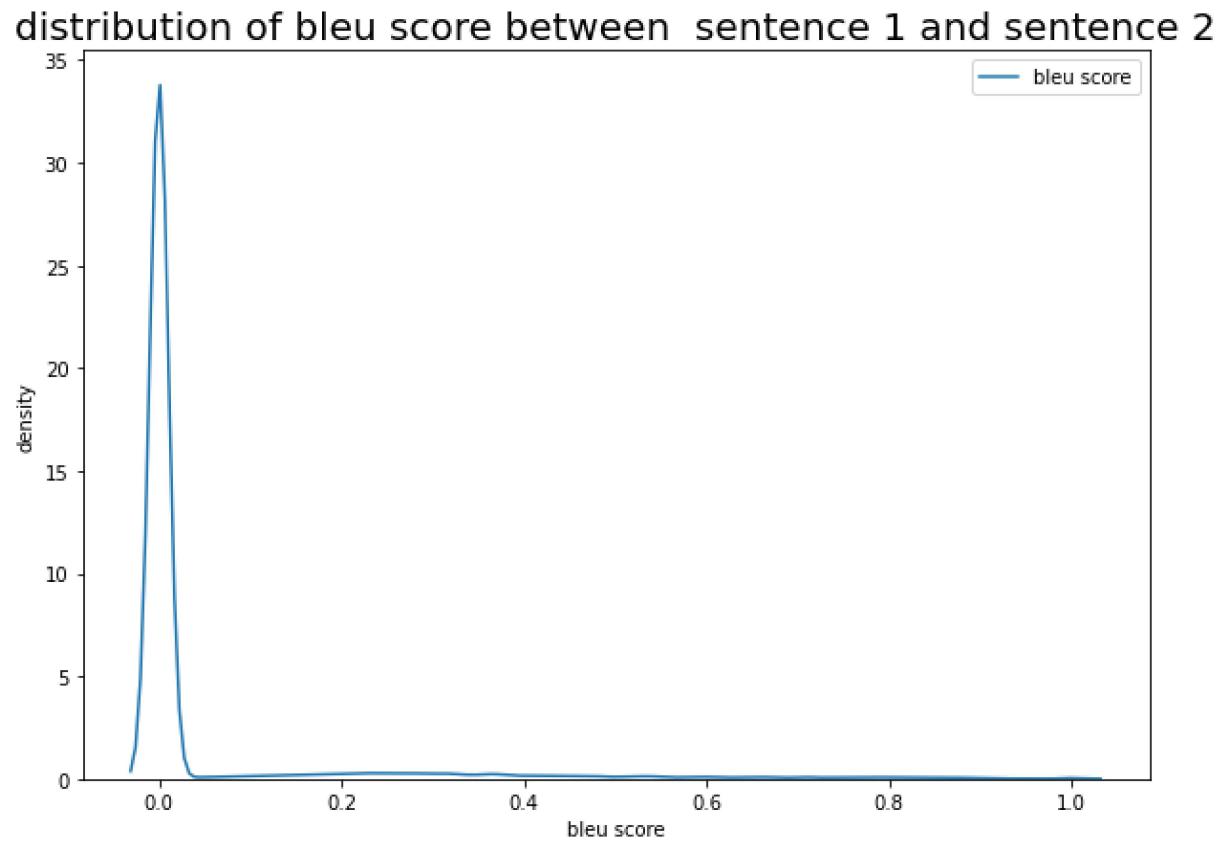
```

1 train_df['bleu_score'] = train_bleu_lst
2 val_df['bleu_score'] = val_bleu_lst
3 test_df['bleu_score'] = test_bleu_lst

```

In [56]:

```
1 fig,ax = plt.subplots(1,1, figsize=(10,7))
2 sns.kdeplot(x = train_df['bleu_score'], ax = ax, label = "bleu score")
3 ax.set_xlabel('bleu score')
4 ax.set_ylabel('density')
5 ax.set_title('distribution of bleu score between sentence 1 and sentence 2')
6 plt.legend()
7 plt.show()
```



**Most of the bleu scores are close to 0 few are greater than 0 and even fewer reach 1**

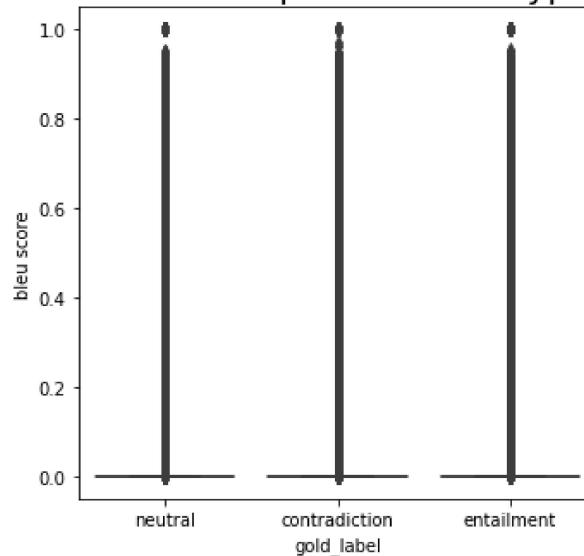
In [57]:

```

1 fig,ax = plt.subplots(1,1, figsize=(5,5), sharex=True)
2 sns.boxplot(x = train_df['gold_label'], y = train_df['bleu_score'], ax = ax)
3 ax.set_title('box plot of bleu score between premise and hypothesis for each')
4 ax.set_ylabel('bleu score')
5 plt.tight_layout()
6 plt.subplots_adjust(hspace=0.3)
7 plt.show()

```

box plot of bleu score between premise and hypothesis for each label



## Let us plot on log scale

In [52]:

```

1 def get_log(x):
2     if x != 0:
3         return math.log(x)
4     else:
5         return 0

```

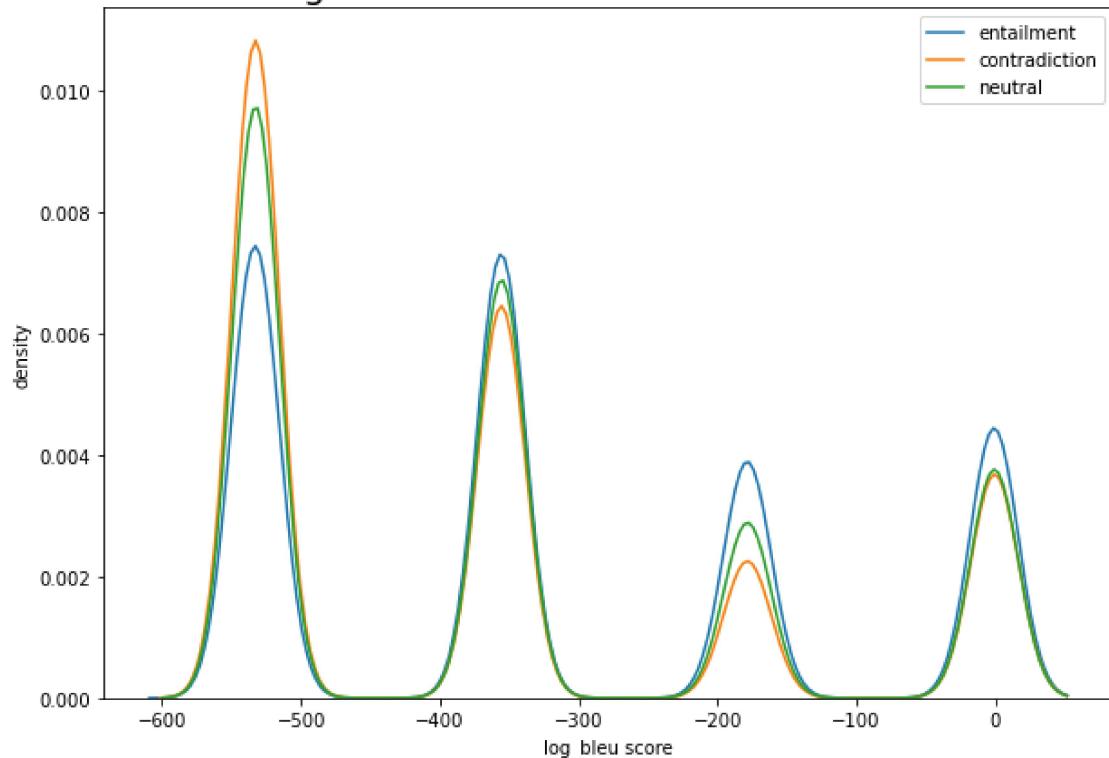
In [59]:

```
1 train_df['log_bleu_score'] = train_df['bleu_score'].apply(get_log)
```

In [60]:

```
1 fig,ax = plt.subplots(1,1, figsize=(10,7))
2 sns.kdeplot(x = train_df[train_df['gold_label'] == 'entailment']['log_bleu_s
3 sns.kdeplot(x = train_df[train_df['gold_label'] == 'contradiction']['log_ble
4 sns.kdeplot(x = train_df[train_df['gold_label'] == 'neutral']['log_bleu_scor
5 ax.set_xlabel('log_bleu score')
6 ax.set_ylabel('density')
7 ax.set_title('distribution of log bleu score between sentence 1 and sentenc
8 plt.legend()
9 plt.show()
```

distribution of log bleu score between sentence 1 and sentence 2



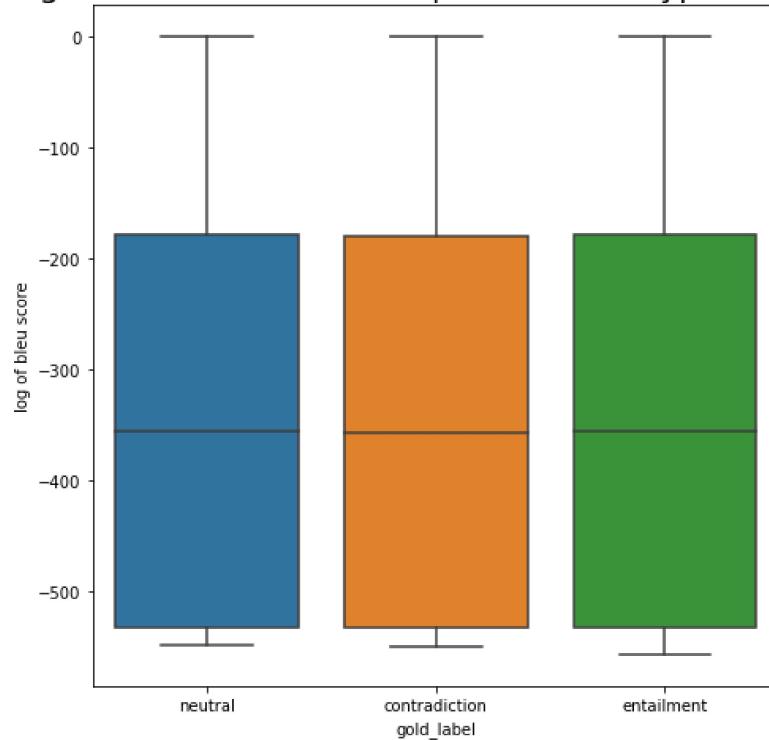
In [61]:

```

1 fig,ax = plt.subplots(1,1, figsize=(7,7), sharex=True)
2 sns.boxplot(x = train_df['gold_label'], y = train_df['log_bleu_score'], ax =
3 ax.set_title('box plot of log of bleu score between premise and hypothesis f
4 ax.set_ylabel('log of bleu score')
5 plt.tight_layout()
6 plt.show()

```

box plot of log of bleu score between premise and hypothesis for each label



it seems that log of bleu score has some discriminating power in some ranges

In [53]:

```
1 val_df['log_bleu_score'] = val_df['bleu_score'].apply(get_log)
```

In [54]:

```
1 test_df['log_bleu_score'] = test_df['bleu_score'].apply(get_log)
```

## Sentiment score

```
In [62]: 1 s = 'this food is too bad'
2 s = TextBlob(s)
3 print(s.sentiment.polarity)
4 print(s.sentiment.subjectivity)
```

-0.6999999999999998  
0.6666666666666666

```
In [63]: 1 def get_sentiment(x, option):
2
3     """
4         returns polarity score if option = 'polarity'
5         returns subjectivity score if option = 'subjectivity'
6         x = input sentence
7     """
8
9     x = TextBlob(x)
10    if option == 'polarity':
11        return x.sentiment.polarity
12    if option == 'subjectivity':
13        return x.sentiment.subjectivity
14
```

## calculating polarity score sentence1

```
In [64]: 1 train_df['sentence1_polarity'] = train_df['sentence1'].apply(lambda x: get_s
2 val_df['sentence1_polarity'] = val_df['sentence1'].apply(lambda x: get_senti
3 test_df['sentence1_polarity'] = test_df['sentence1'].apply(lambda x: get_sen
```

## calculating polarity score sentence2

```
In [65]: 1 train_df['sentence2_polarity'] = train_df['sentence2'].apply(lambda x: get_s
2 val_df['sentence2_polarity'] = val_df['sentence2'].apply(lambda x: get_senti
3 test_df['sentence2_polarity'] = test_df['sentence2'].apply(lambda x: get_sen
```

## calculating subjectivity score sentence1

```
In [66]: 1 train_df['sentence1_subjectivity'] = train_df['sentence1'].apply(lambda x: g
2 val_df['sentence1_subjectivity'] = val_df['sentence1'].apply(lambda x: get_s
3 test_df['sentence1_subjectivity'] = test_df['sentence1'].apply(lambda x: get
```

## calculating subjectivity score sentence2

```
In [67]: 1 train_df['sentence2_subjectivity'] = train_df['sentence2'].apply(lambda x: g
2 val_df['sentence2_subjectivity'] = val_df['sentence2'].apply(lambda x: get_s
3 test_df['sentence2_subjectivity'] = test_df['sentence2'].apply(lambda x: get
```

## Let us conduct some analysis in cases where polarity scores for two sentences are equal and not equal

In [68]: 1 train\_df.shape

Out[68]: (549361, 13)

In [69]: 1 train\_df[train\_df['sentence1\_polarity'] == train\_df['sentence2\_polarity']].s

Out[69]: (250594, 13)

We can see that out of 549361 training examples, 250594 examples have same exactly polarity score for sentence1 and sentence2. Let us see how that relates to our target variable

In [70]: 1 train\_df[train\_df['sentence1\_polarity'] == train\_df['sentence2\_polarity']].g

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_r
gold_label					
contradiction	82970	82970	82970	82970	82970
entailment	91746	91746	91746	91746	91746
neutral	75878	75878	75878	75878	75878

In [71]: 1 82970+91746+75878

Out[71]: 250594

In [72]: 1 train\_df[train\_df['sentence1\_polarity'] != train\_df['sentence2\_polarity']].g

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_r
gold_label					
contradiction	100215	100215	100215	100215	100215
entailment	91668	91668	91668	91668	91668
neutral	106884	106884	106884	106884	106884

there are almost equal number of training examples for entailment when polarity score is equal and not equal for two sentences. For contradiction number of training examples are more when polarity score of two sentences are not equal. Similar case can be seen for neutral. So equality of polarity score may have some predictive power.

## Let us conduct some analysis in cases where subjectivity scores for two sentences are equal and not equal

In [73]: 1 train\_df.shape

Out[73]: (549361, 13)

In [74]: 1 train\_df[train\_df['sentence1\_subjectivity'] == train\_df['sentence2\_subjectivity']]

Out[74]: (208181, 13)

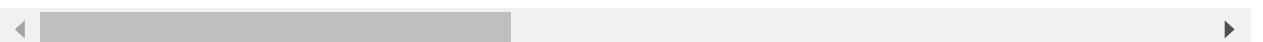
We can see that out of 549361 training examples, 208181 examples have same exactly subjectivity score for sentence1 and sentence2. Let us see how that relates to our target variable

In [75]: 1 train\_df[train\_df['sentence1\_subjectivity'] == train\_df['sentence2\_subjectivity']]

Out[75]:

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_label
<b>gold_label</b>					
<b>contradiction</b>	69715	69715	69715	69715	69715
<b>entailment</b>	75200	75200	75200	75200	75200
<b>neutral</b>	63266	63266	63266	63266	63266

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_label
<b>gold_label</b>					
<b>contradiction</b>	69715	69715	69715	69715	69715
<b>entailment</b>	75200	75200	75200	75200	75200
<b>neutral</b>	63266	63266	63266	63266	63266

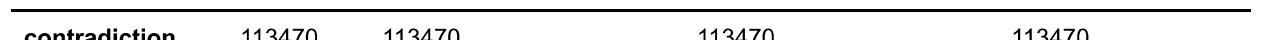


In [76]: 1 train\_df[train\_df['sentence1\_subjectivity'] != train\_df['sentence2\_subjectivity']]

Out[76]:

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_label
<b>gold_label</b>					
<b>contradiction</b>	113470	113470	113470	113470	113470
<b>entailment</b>	108214	108214	108214	108214	108214
<b>neutral</b>	119496	119496	119496	119496	119496

	sentence1	sentence2	sentence1_preprocessed	sentence2_preprocessed	sentence1_label
<b>gold_label</b>					
<b>contradiction</b>	113470	113470	113470	113470	113470
<b>entailment</b>	108214	108214	108214	108214	108214
<b>neutral</b>	119496	119496	119496	119496	119496



There is a lot of difference between number of training examples when subjectivity is equal and when it is not equal. Even though this does not contribute to discriminative power between entailment, contradiction and neutral labels it can still be of some help in predictions. We will have to see this when proceed for actual modelling

In [77]: 1 train\_df.to\_csv('train\_df.csv')  
2 val\_df.to\_csv('val\_df.csv')  
3 test\_df.to\_csv('test\_df.csv')

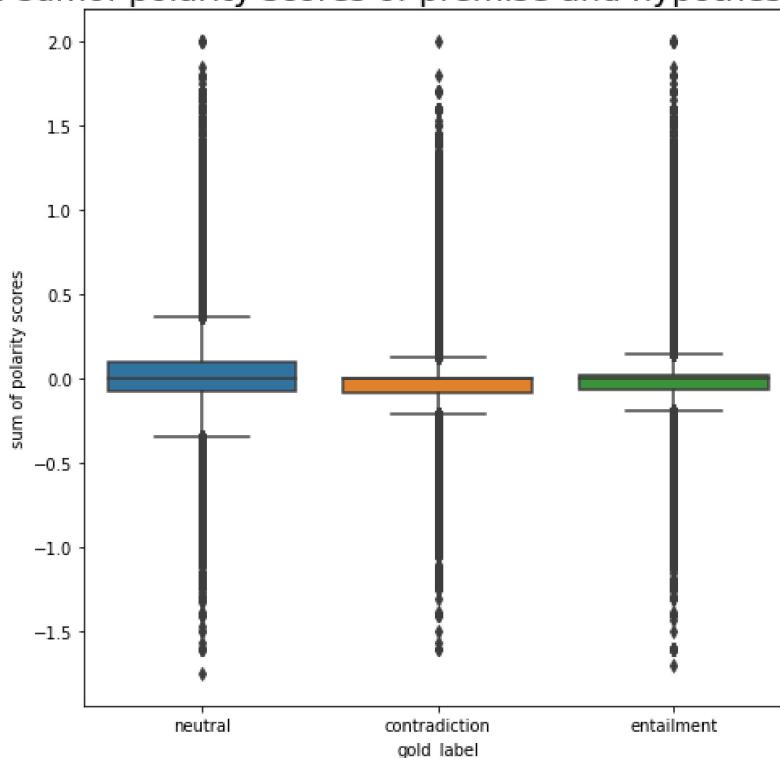
In [4]: 1 train\_df = pd.read\_csv('train\_df.csv', index\_col = 0)  
2 val\_df = pd.read\_csv('val\_df.csv', index\_col = 0)  
3 test\_df = pd.read\_csv('test\_df.csv', index\_col = 0)

## Sum of polarity scores of sentence1 and sentence2

```
In [5]: 1 train_df['sum_polarity'] = train_df['sentence1_polarity'] + train_df['senten
```

```
In [6]: 1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['sum_polarity'], ax = a
3 ax.set_title('box plot of sumof polarity scores of premise and hypothesis fo
4 ax.set_ylabel('sum of polarity scores')
5 plt.tight_layout()
6 plt.show()
```

box plot of sumof polarity scores of premise and hypothesis for each label



sum of polarity scores seems to be somewhat helpful in classification. polarity scores lie between -1(negative) and 1(positive). So if both sentences have similar polarity then adding those will further make these polarity scores extreme. For example if both sentences have polarity score of -1 then adding both will give -2. But if the two sentences have different polarity then adding these will give(most probably) a score near to 0. This behaviour is what is visible in the box plots. Most of the added polarity scores lie near to 0. Then there are no added scores around the boxes and added polarity scores can be seen further away from the boxes(extreme values). Also the boxes are not perfectly aligned for all the three labels. So this feature might be helpful in classification

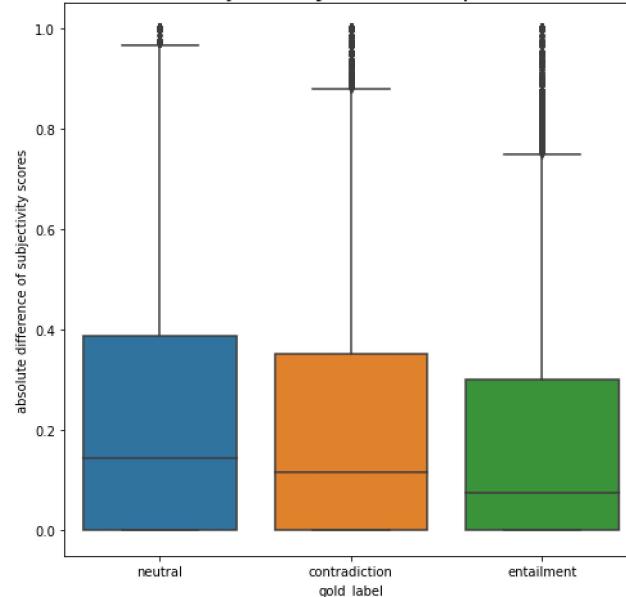
```
In [7]: 1 val_df['sum_polarity'] = val_df['sentence1_polarity'] + val_df['sentence2_po
2 test_df['sum_polarity'] = test_df['sentence1_polarity'] + test_df['sentence2_
```

## absolute difference of subjectivity scores of two sentences

```
In [8]: 1 train_df['diff_subjectivity'] = (train_df['sentence1_subjectivity'] - train_
```

```
In [9]: 1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['diff_subjectivity'], a
3 ax.set_title('box plot of absolute difference of subjectivity scores of prem
4 ax.set_ylabel('absolute difference of subjectivity scores')
5 plt.tight_layout()
6 plt.show()
```

box plot of absolute difference of subjectivity scores of premise and hypothesis for each label



**absolute difference between subjectivity scores for neutral seem to be higher than contradiction and entailment. Further the same score for contradiction is higher than entailment.**

```
In [10]: 1 val_df['diff_subjectivity'] = (val_df['sentence1_subjectivity'] - val_df['se
2 test_df['diff_subjectivity'] = (test_df['sentence1_subjectivity'] - test_df[
```

## word count based cosine similarity

```
In [11]: 1 #function to calculate cosine similarity between two sentences
2 WORD = re.compile(r"\w+")
3 def get_cosine(vec1, vec2):
4     intersection = set(vec1.keys()) & set(vec2.keys())
5     numerator = sum([vec1[x] * vec2[x] for x in intersection])
6
7     sum1 = sum([vec1[x]**2 for x in list(vec1.keys())])
8     sum2 = sum([vec2[x]**2 for x in list(vec2.keys())])
9     denominator = math.sqrt(sum1) * math.sqrt(sum2)
10
11    if not denominator:
12        return 0.0
13    else:
14        return float(numerator) / denominator
15
16 def text_to_vector(text):
17     words = WORD.findall(text)
18     return Counter(words)
19
```

```
In [12]: 1 text1 = train_df['sentence1'][0]
2 text2 = train_df['sentence2'][0]
3
4 vector1 = text_to_vector(text1)
5 vector2 = text_to_vector(text2)
6
7 cosine = get_cosine(vector1, vector2)
```

## Train cosine similarity

```
In [13]: 1 train_cosine_sim = list()
2 for i, row in train_df.iterrows():
3     text1 = row['sentence1']
4     text2 = row['sentence2']
5     vector1 = text_to_vector(text1)
6     vector2 = text_to_vector(text2)
7     cosine = get_cosine(vector1, vector2)
8     train_cosine_sim.append(cosine)
9
10
```

```
In [14]: 1 train_df['cosine_similarity'] = train_cosine_sim
```

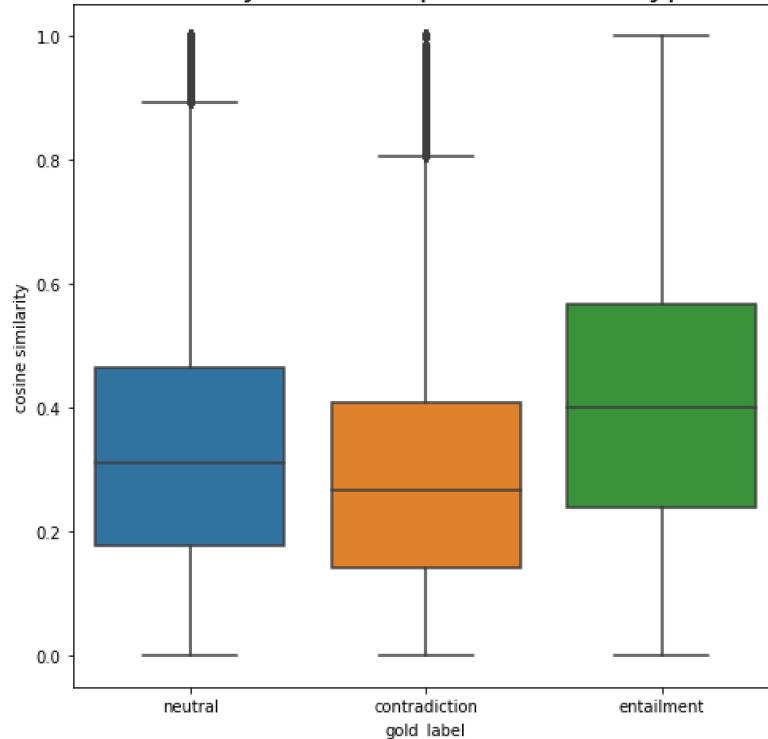
```
In [15]: 1 print(f"minimum cosine similarity: {train_df['cosine_similarity'].min()}\nma
2 {train_df['cosine_similarity'].max()}")
```

```
minimum cosine similarity: 0.0
maximum cosine similarity: 1.0000000000000002
```

In [16]:

```
1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['cosine_similarity'], a
3 ax.set_title('box plot of cosine similarity between premise and hypothesis f
4 ax.set_ylabel('cosine similarity')
5 plt.tight_layout()
6 plt.show()
```

box plot of cosine similarity between premise and hypothesis for each label



**cosine similarity does seem to have some predictive power. Particulary cosine similarities for contradiction and entailment do not overlap much**

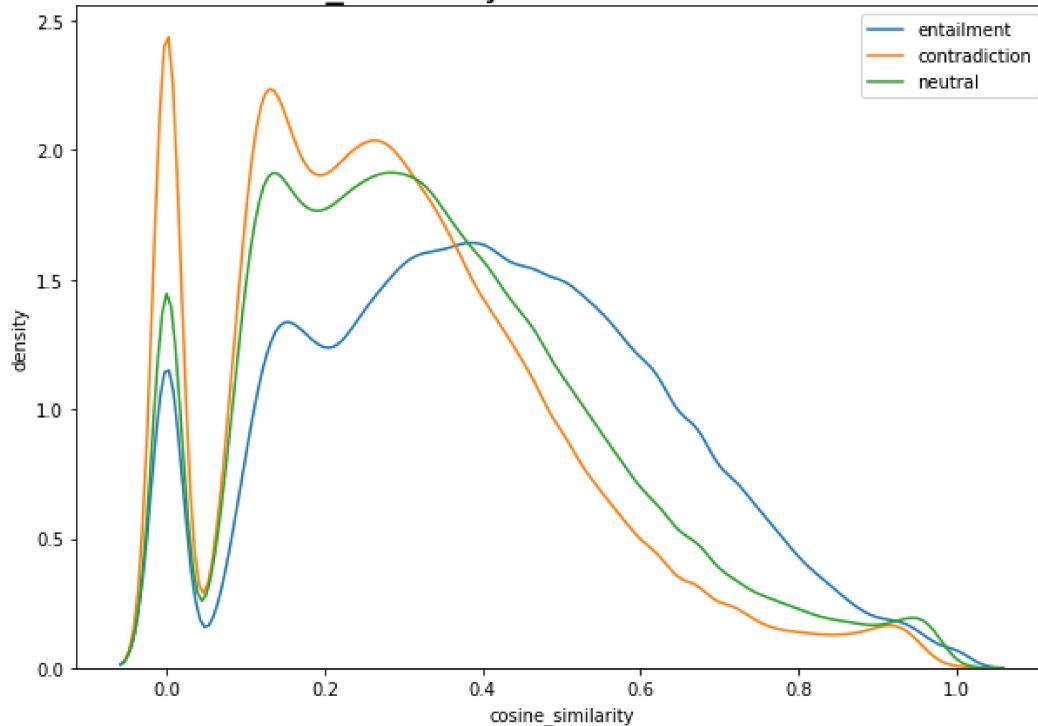
In [103]:

```

1 fig,ax = plt.subplots(1,1, figsize=(10,7))
2 sns.kdeplot(x = train_df[train_df['gold_label'] == 'entailment']['cosine_sim']
3 sns.kdeplot(x = train_df[train_df['gold_label'] == 'contradiction']['cosine_'
4 sns.kdeplot(x = train_df[train_df['gold_label'] == 'neutral']['cosine_simila'
5 ax.set_xlabel('cosine_similarity')
6 ax.set_ylabel('density')
7 ax.set_title('distribution of cosine_similarity between sentence 1 and sent'
8 plt.legend()
9 plt.show()

```

distribution of cosine\_similarity between sentence 1 and sentence 2



**mostly cosine similarities for all the labels follow similar distribution. But towards the higher end of the distribution(after cosine similarity value of 0.4) more desnity is seen for entailment pairs. This should be expected because if hypothesis entails premise then they should be similar**

## validation cosine similarity

```
In [17]: 1 val_cosine_sim = list()
2 for i,row in val_df.iterrows():
3     text1 = row['sentence1']
4     text2 = row['sentence2']
5     vector1 = text_to_vector(text1)
6     vector2 = text_to_vector(text2)
7     cosine = get_cosine(vector1, vector2)
8     val_cosine_sim.append(cosine)
9
10
```

```
In [19]: 1 val_df['cosine_similarity'] = val_cosine_sim
```

## Test cosine similarity

```
In [20]: 1 test_cosine_sim = list()
2 for i,row in test_df.iterrows():
3     text1 = row['sentence1']
4     text2 = row['sentence2']
5     vector1 = text_to_vector(text1)
6     vector2 = text_to_vector(text2)
7     cosine = get_cosine(vector1, vector2)
8     test_cosine_sim.append(cosine)
9
10
```

```
In [21]: 1 test_df['cosine_similarity'] = test_cosine_sim
```

```
In [ ]:
```

## Cosine similarity between tfidf vectors

```
In [31]: 1 #https://stackoverflow.com/questions/8897593/how-to-compute-the-similarity-between-strings-in-python
2 stemmer = nltk.stem.porter.PorterStemmer()
3 remove_punctuation_map = dict((ord(char), None) for char in string.punctuation)
4
5 def stem_tokens(tokens):
6     return [stemmer.stem(item) for item in tokens]
7
8 '''remove punctuation, lowercase, stem'''
9 def normalize(text):
10     return stem_tokens(nltk.word_tokenize(text.lower()).translate(remove_punctuation_map))
11
12 vectorizer = TfidfVectorizer(tokenizer=normalize)
13
14 def cosine_sim(text1, text2):
15     tfidf = vectorizer.fit_transform([text1, text2])
16     return ((tfidf * tfidf.T).A)[0,1]
17
```

## tfidf cosine similarity for training data

```
In [32]: 1 tfidf_cos_sim = list()
2 for i,row in train_df.iterrows():
3     cos_sim = cosine_sim(row['sentence1'], row['sentence2'])
4     tfidf_cos_sim.append(cos_sim)
```

```
In [33]: 1 train_df['tfidf_cosine_sim'] = tfidf_cos_sim
```

## tfidf cosine similarity for validation data

```
In [34]: 1 tfidf_cos_sim = list()
2 for i,row in val_df.iterrows():
3     cos_sim = cosine_sim(row['sentence1'], row['sentence2'])
4     tfidf_cos_sim.append(cos_sim)
```

```
In [35]: 1 val_df['tfidf_cosine_sim'] = tfidf_cos_sim
```

## tfidf cosine similarity for test data

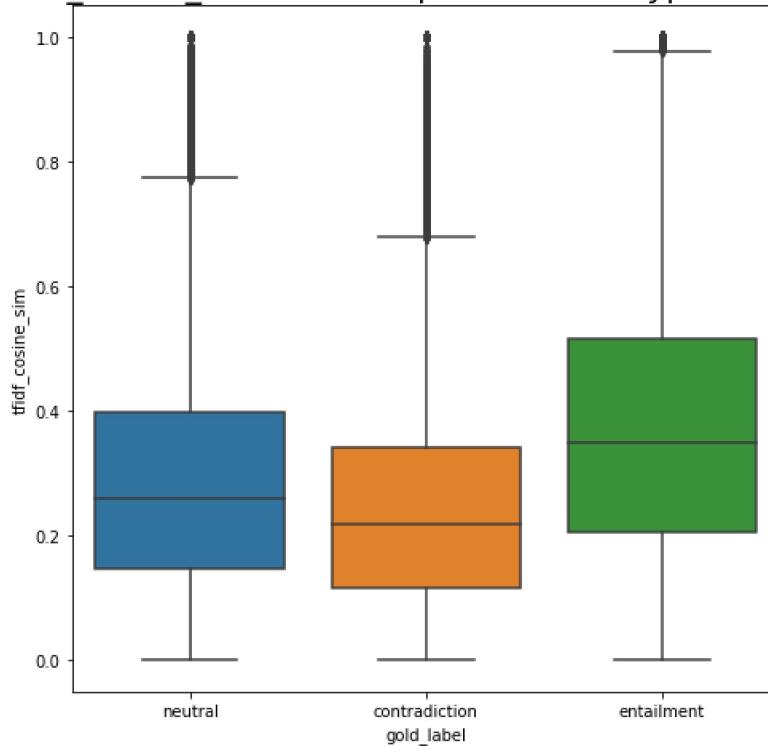
```
In [36]: 1 tfidf_cos_sim = list()
2 for i,row in test_df.iterrows():
3     cos_sim = cosine_sim(row['sentence1'], row['sentence2'])
4     tfidf_cos_sim.append(cos_sim)
```

```
In [37]: 1 test_df['tfidf_cosine_sim'] = tfidf_cos_sim
```

In [38]:

```
1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['tfidf_cosine_sim'], ax= ax)
3 ax.set_title('box plot of tfidf_cosine_sim between premise and hypothesis for each label')
4 ax.set_ylabel('tfidf_cosine_sim')
5 plt.tight_layout()
6 plt.show()
```

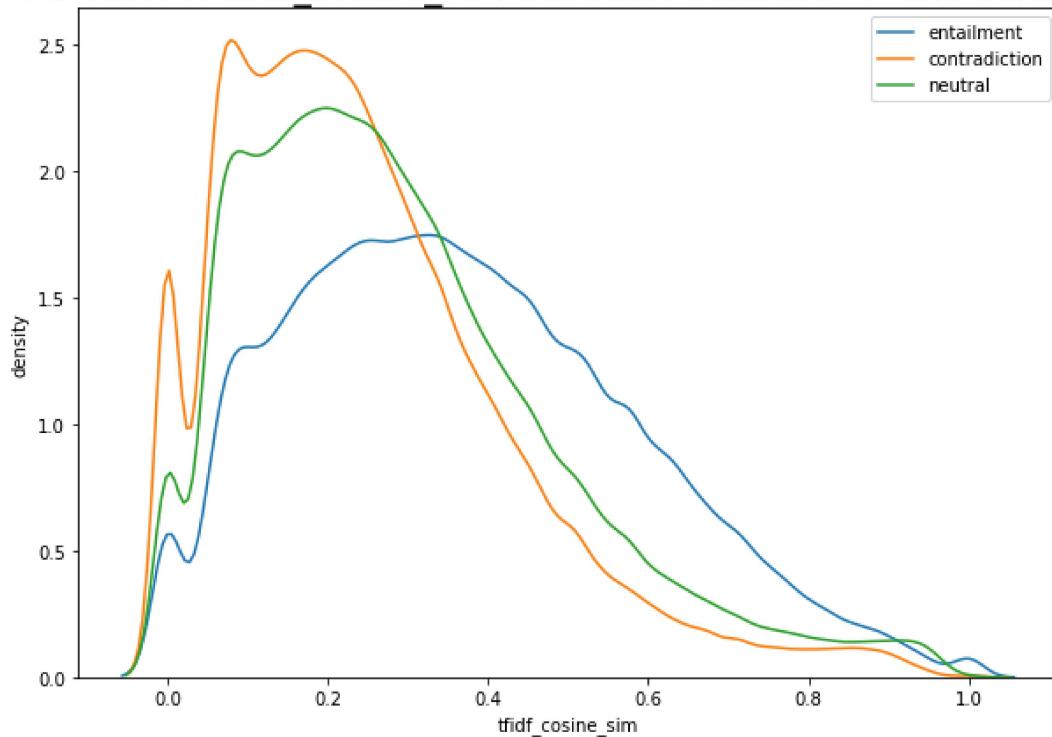
box plot of tfidf\_cosine\_sim between premise and hypothesis for each label



In [39]:

```
1 fig,ax = plt.subplots(1,1, figsize=(10,7))
2 sns.kdeplot(x = train_df[train_df['gold_label'] == 'entailment']['tfidf_cosine_sim'])
3 sns.kdeplot(x = train_df[train_df['gold_label'] == 'contradiction']['tfidf_cosine_sim'])
4 sns.kdeplot(x = train_df[train_df['gold_label'] == 'neutral']['tfidf_cosine_sim'])
5 ax.set_xlabel('tfidf_cosine_sim')
6 ax.set_ylabel('density')
7 ax.set_title('distribution of tfidf_cosine_sim between sentence 1 and sentence 2')
8 plt.legend()
9 plt.show()
```

distribution of tfidf\_cosine\_sim between sentence 1 and sentence 2



cosine similarity calculated on tfidf vectors of sentence1 and sentence2 seem to be good when it comes to separating target variable classes. This similarity is lowest for contradictory pair and highest for entailing pairs which is expected behaviour. There is a kind of a crossing region between 0.2 and 0.3 where kde for entailment and kde for contradiction intersect each other showing that for lower values of cosine similarity density of contradictory pairs is greater while for higher values of cosine similarity density of entailment pairs is greater

```
In [61]: 1 train_df.to_csv('train_df.csv')
2 val_df.to_csv('val_df.csv')
3 test_df.to_csv('test_df.csv')
```

```
In [174]: 1 train_df = pd.read_csv('train_df.csv', index_col = 0)
2 val_df = pd.read_csv('val_df.csv', index_col = 0)
3 test_df = pd.read_csv('test_df.csv', index_col = 0)
```

## Count of pos tags for sentence1 and sentence2

```
In [63]: 1 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\anike\AppData\Roaming\nltk_data...
[nltk_data]     Package wordnet is already up-to-date!
```

```
Out[63]: True
```

```
In [178]: 1 def pos_counter(text, s):
2     """
3         This function returns the count of occurrence of different parts of speech
4         argument s specifies which sentence are we processing. sentence1 or sentence2
5     """
6     pos_cols = ['CC', 'IN', 'JJ', 'JJR', 'JJS', 'NN', 'NNP', 'NNS', 'PRP', 'PRP$',
7     tokens = nltk.word_tokenize(text)
8     lemmatizer = WordNetLemmatizer()
9     lemmatized_tokens = [lemmatizer.lemmatize(token) for token in tokens]
10    tags = nltk.pos_tag(lemmatized_tokens)
11    counts = Counter(tag for word, tag in tags)
12    if s == 'sentence1':
13        count_dict = {s1_+k:v for k,v in dict(counts).items() if k in pos_
14    if s == 'sentence2':
15        count_dict = {s2_+k:v for k,v in dict(counts).items() if k in pos_
16
17    return count_dict
18
```

```
In [181]: 1 s = "dark matter written by blake crouch."
2 s_c = pos_counter(s, 'sentence2')
```

In [182]: 1 s\_c

Out[182]: {'s2\_JJ': 1, 's2\_NN': 3, 's2\_VBN': 1, 's2\_IN': 1}

## Meaning of each pos tag from nltk:

CC: conjunction, coordinating  
 CD: numeral, cardinal  
 DT: determiner  
 EX: existential there  
 IN: preposition or conjunction, subordinating  
 JJ: adjective or numeral, ordinal  
 JJR: adjective, comparative  
 JJS: adjective, superlative  
 LS: list item marker  
 MD: modal auxiliary  
 NN: noun, common, singular or mass  
 NNP: noun, proper, singular  
 NNS: noun, common, plural  
 PDT: pre-determiner  
 PRP: pronoun, personal  
 PRP\$: pronoun, possessive  
 RB: adverb  
 RBR: adverb, comparative  
 RBS: adverb, superlative  
 RP: particle  
 UH: interjection  
 VB: verb, base form  
 VBD: verb, past tense  
 VBG: verb, present participle or gerund  
 VBN: verb, past participle  
 VBP: verb, present tense, not 3rd person singular  
 VBZ: verb, present tense, 3rd person singular  
 WDT: WH-determiner  
 WP: WH-pronoun  
 WRB: Wh-adverb

In [80]: # Let us get all pos tags in sentence1

```

1 unique_tags = set()
2 for i, row in train_df.iterrows():
3     tag_dict = pos_counter(row['sentence1'])
4     for k, v in tag_dict.items():
5         unique_tags.add(k)
6
7

```

```
In [84]: 1 # let us get all pos tags in sentence2
          2 unique_tags2 = set()
          3 for i, row in train_df.iterrows():
          4     tag_dict = pos_counter(row['sentence2'])
          5     for k, v in tag_dict.items():
          6         unique_tags2.add(k)
          7
```

```
In [86]: 1 # combined unique pos tags in both sentences
          2 total_tags = unique_tags.union(unique_tags2)
```

```
In [87]: 1 len(total_tags)
```

Out[87]: 44

## Process applied to get count of pos tags for both the sentences

1. Create empty columns in dataframes for count of pos in two sentences.
2. Initialize all values first to 0. prefix 's1' for each tag represents it belongs to count of sentence1 and prefix 's2' represents it belongs to count of sentence2
3. consider these pos tags only: CC, IN, JJ, JJR, JJS, NN, NNP, NNS, PRP, PRP\$, RB, RBR, RBS, UH, VB, VBD, VBG, VBN, VBP, VBZ
4. now iterate through each row, for each row get dictionary of count of pos tags for sentence1 and sentence2
5. update count of each pos tag for both the sentences

## Generating pos tags count for train data

```
In [218]: 1 # adding empty columns to store count of each pos tag mentioned above
          2 train_df = train_df.reindex(columns=[*train_df.columns.tolist(), 's1_CC', 's1_IN',
          3                               's1_NN', 's1_NNP', 's1_NNS', 's1_PRP',
          4                               's1_UH', 's1_VB', 's1_VBD', 's1_VBG',
          5                               's2_JJ', 's2_JJR', 's2_JJS', 's2_NN',
          6                               's2_RBR', 's2_RBS', 's2_UH', 's2_VB', 's2_VBG'])
```

```
In [210]: 1 pos_cols = ['CC', 'IN', 'JJ', 'JJR', 'JJS', 'NN', 'NNP', 'NNS', 'PRP', 'PRP$']
```

```
In [219]: 1 #iterating through each row and storing count of each pos tag for both the s
2 for i,row in train_df.iterrows():
3     s1_pos = pos_counter(row['sentence1'], 'sentence1')
4     s2_pos = pos_counter(row['sentence2'], 'sentence2')
5     for k,v in s1_pos.items():
6         train_df.at[i, k] = v
7     for k,v in s2_pos.items():
8         train_df.at[i, k] = v
9
```

## Generating pos tags count for validation data

```
In [228]: 1 # adding empty columns to store count of each pos tag mentioned above
2 val_df = val_df.reindex(columns=[*val_df.columns.tolist(), 's1_CC', 's1_IN',
3                             's1_NN', 's1_NNP', 's1_NNS', 's1_PRP',
4                             's1_UH', 's1_VB', 's1_VBD', 's1_VBG',
5                             's2_JJ', 's2_JJR', 's2_JJS', 's2_NN',
6                             's2_RBR', 's2_RBS', 's2_UH', 's2_VB', 's2_VBG'])
```

```
In [229]: 1 #iterating through each row and storing count of each pos tag for both the s
2 for i,row in val_df.iterrows():
3     s1_pos = pos_counter(row['sentence1'], 'sentence1')
4     s2_pos = pos_counter(row['sentence2'], 'sentence2')
5     for k,v in s1_pos.items():
6         val_df.at[i, k] = v
7     for k,v in s2_pos.items():
8         val_df.at[i, k] = v
9
```

## Generating pos tags count for test data

```
In [230]: 1 # adding empty columns to store count of each pos tag mentioned above
2 test_df = test_df.reindex(columns=[*test_df.columns.tolist(), 's1_CC', 's1_I',
3                               's1_NN', 's1_NNP', 's1_NNS', 's1_PRP',
4                               's1_UH', 's1_VB', 's1_VBD', 's1_VBG',
5                               's2_JJ', 's2_JJR', 's2_JJS', 's2_NN',
6                               's2_RBR', 's2_RBS', 's2_UH', 's2_VB', 's2_VBG'])
```

```
In [231]: 1 #iterating through each row and storing count of each pos tag for both the s
2 for i,row in test_df.iterrows():
3     s1_pos = pos_counter(row['sentence1'], 'sentence1')
4     s2_pos = pos_counter(row['sentence2'], 'sentence2')
5     for k,v in s1_pos.items():
6         test_df.at[i, k] = v
7     for k,v in s2_pos.items():
8         test_df.at[i, k] = v
9
```

```
In [235]: 1 train_df.to_csv('train_df.csv')
2 val_df.to_csv('val_df.csv')
3 test_df.to_csv('test_df.csv')
```

```
In [236]: 1 train_df = pd.read_csv('train_df.csv', index_col = 0)
2 val_df = pd.read_csv('val_df.csv', index_col = 0)
3 test_df = pd.read_csv('test_df.csv', index_col = 0)
```

## Visualizing some important pos tags

### NN:Noun

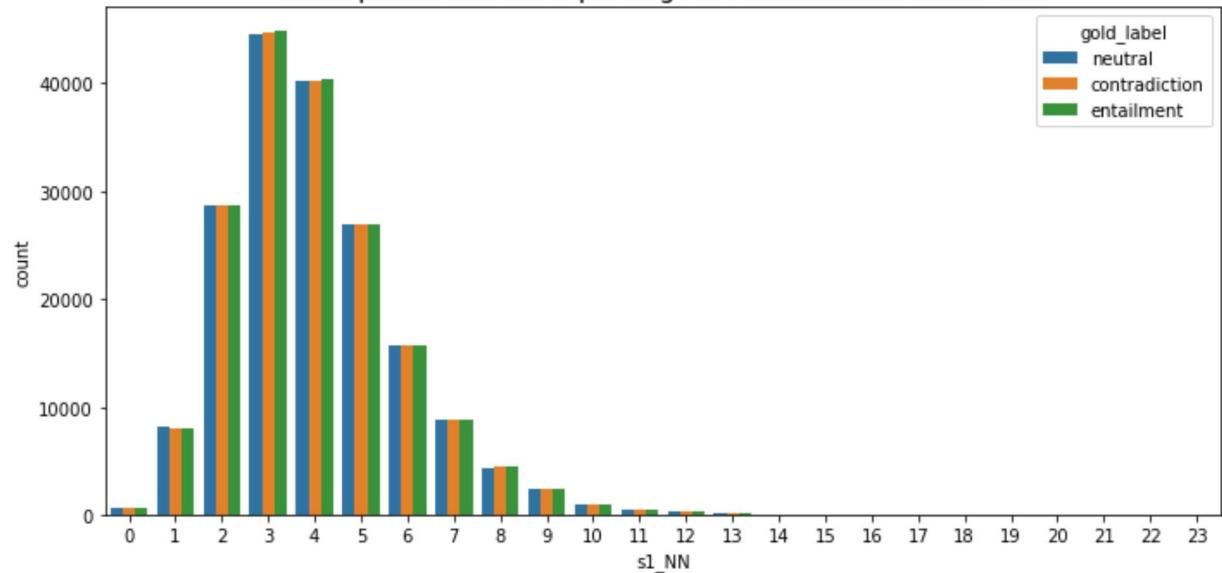
In [241]:

```

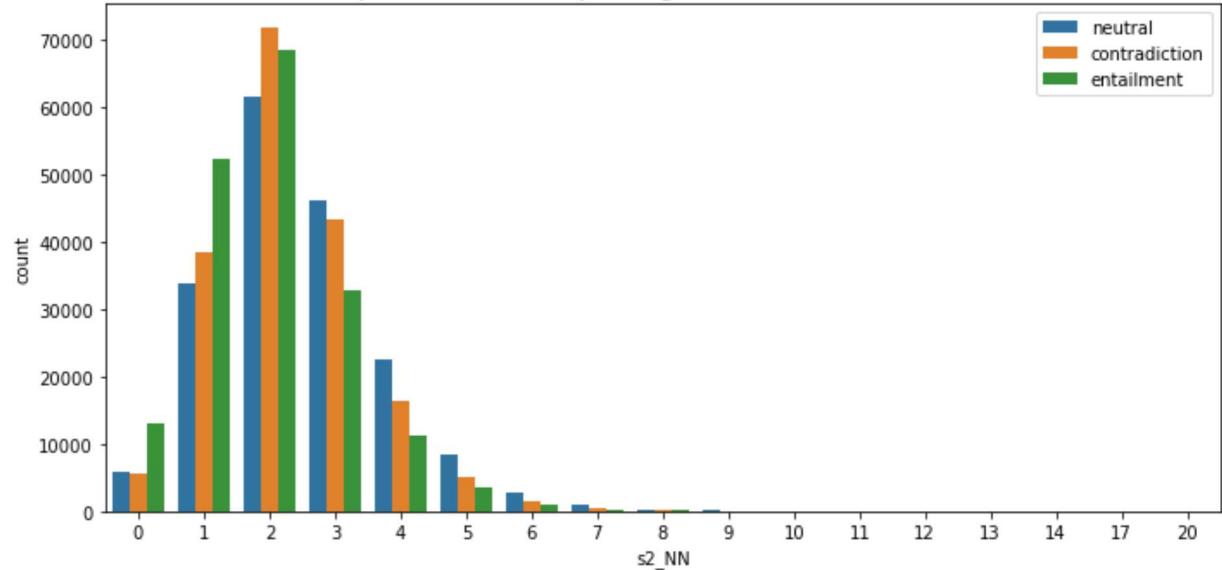
1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_NN'], hue = train_df['gold_label'], ax = ax[0]
3 sns.countplot(x = train_df['s2_NN'], hue = train_df['gold_label'], ax = ax[1]
4 ax[0].set_title('count plot of count NN pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count NN pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```

count plot of count NN pos tag for sentence1 for each label



count plot of count NN pos tag for sentence2 for each label



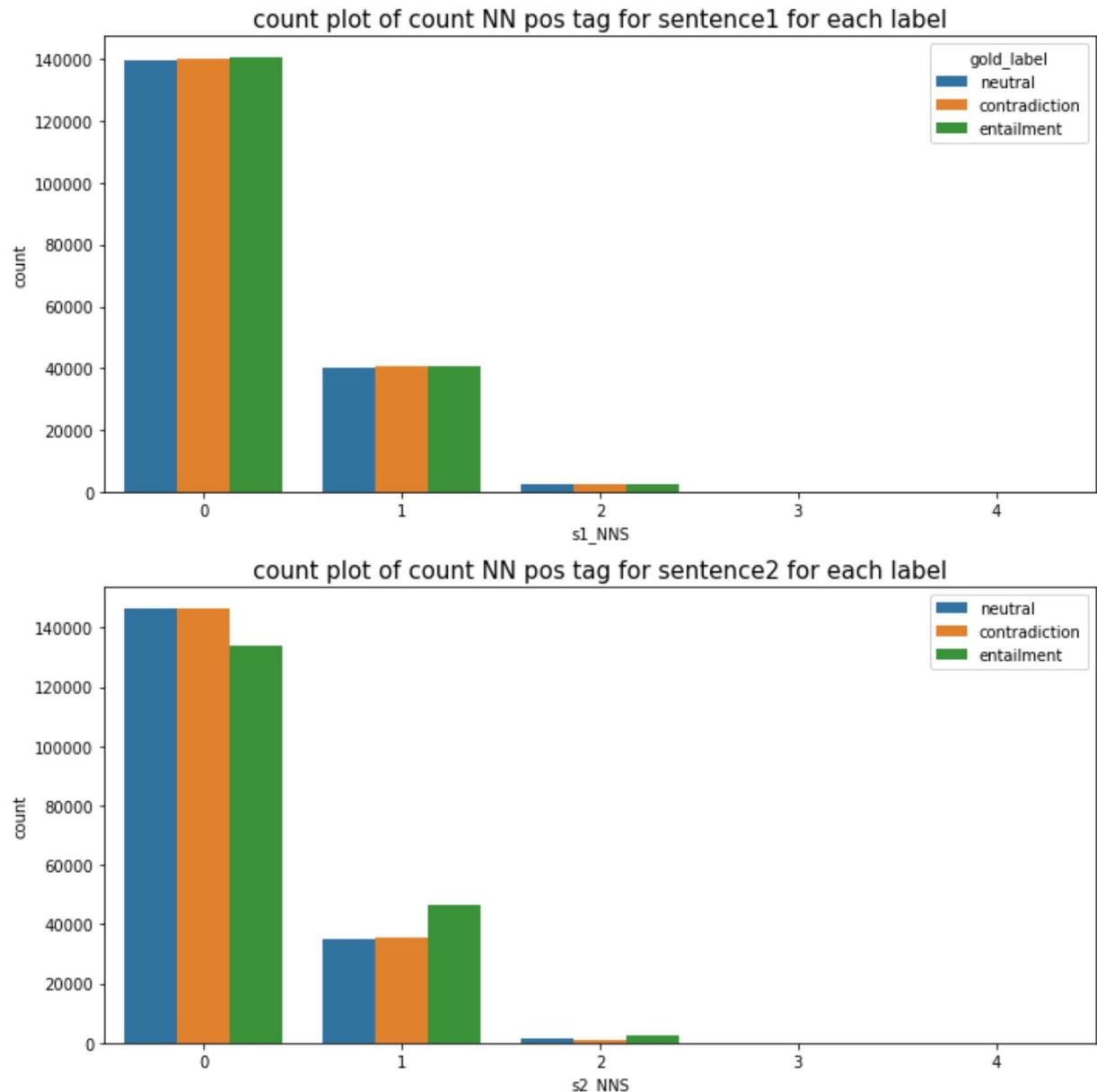
NN: noun, common, singular or mass NNP: noun, proper, singular NNS: noun, common, plural

In [249]:

```

1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_NNS'], hue = train_df['gold_label'], ax = ax[0])
3 sns.countplot(x = train_df['s2_NNS'], hue = train_df['gold_label'], ax = ax[1])
4 ax[0].set_title('count plot of count NN pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count NN pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```



**Let us add all noun pos tags and see the behaviour**

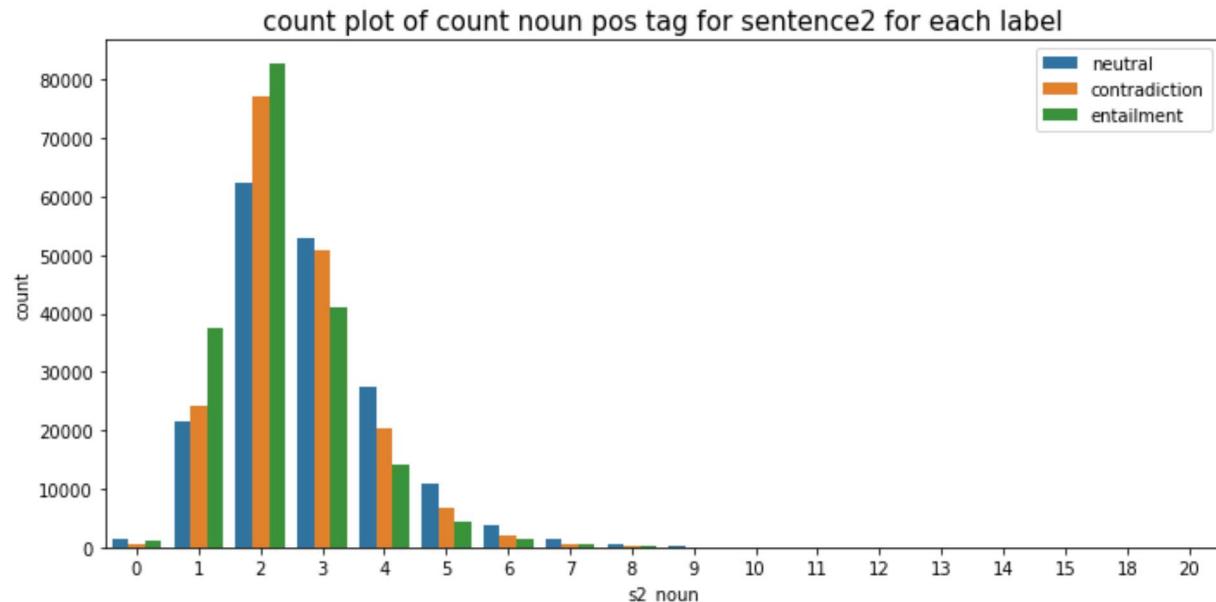
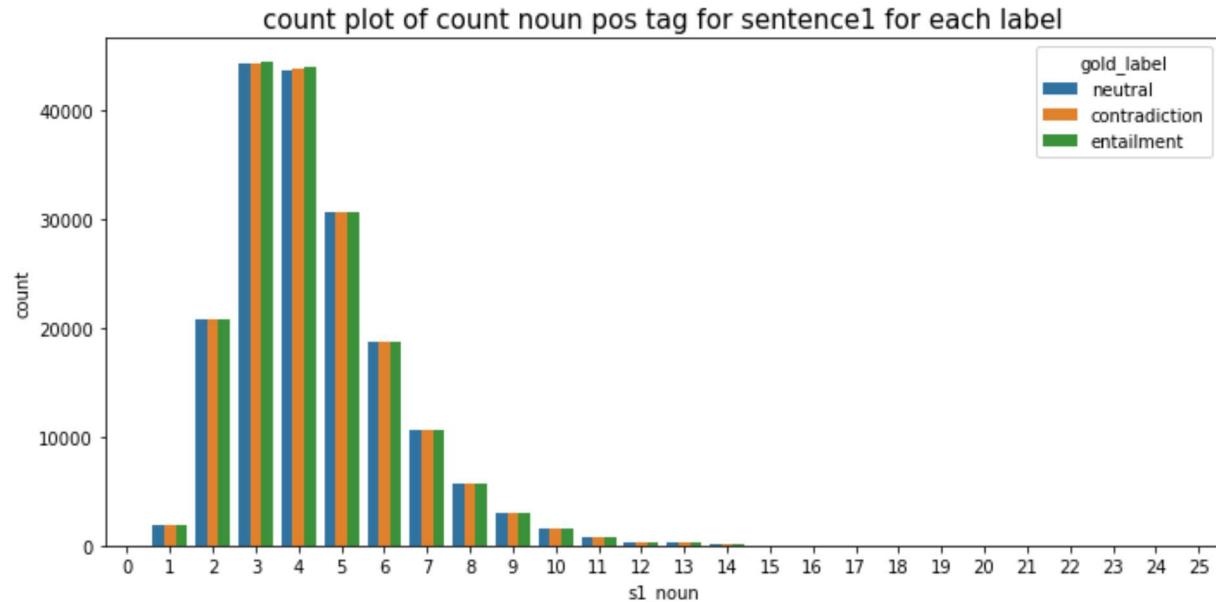
```
In [250]: 1 train_df['s1_noun'] = train_df['s1_NN'] + train_df['s1_NNP'] + train_df['s1_NNS']  
  
In [251]: 1 train_df['s2_noun'] = train_df['s2_NN'] + train_df['s2_NNP'] + train_df['s2_NNS']  
  
In [258]: 1 val_df['s1_noun'] = val_df['s1_NN'] + val_df['s1_NNP'] + val_df['s1_NNS']  
2 val_df['s2_noun'] = val_df['s2_NN'] + val_df['s2_NNP'] + val_df['s2_NNS']  
3 test_df['s1_noun'] = test_df['s1_NN'] + test_df['s1_NNP'] + test_df['s1_NNS']  
4 test_df['s2_noun'] = test_df['s2_NN'] + test_df['s2_NNP'] + test_df['s2_NNS']
```

In [252]:

```

1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_noun'], hue = train_df['gold_label'], ax = ax)
3 sns.countplot(x = train_df['s2_noun'], hue = train_df['gold_label'], ax = ax)
4 ax[0].set_title('count plot of count noun pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count noun pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```



In [291]:

```
1 train_df['ratio_noun'] = (train_df['s1_noun']/(train_df['s2_noun']+2))
```

In [296]:

```

1 val_df['ratio_noun'] = (val_df['s1_noun']/(val_df['s2_noun']+2))
2 test_df['ratio_noun'] = (test_df['s1_noun']/(test_df['s2_noun']+2))

```

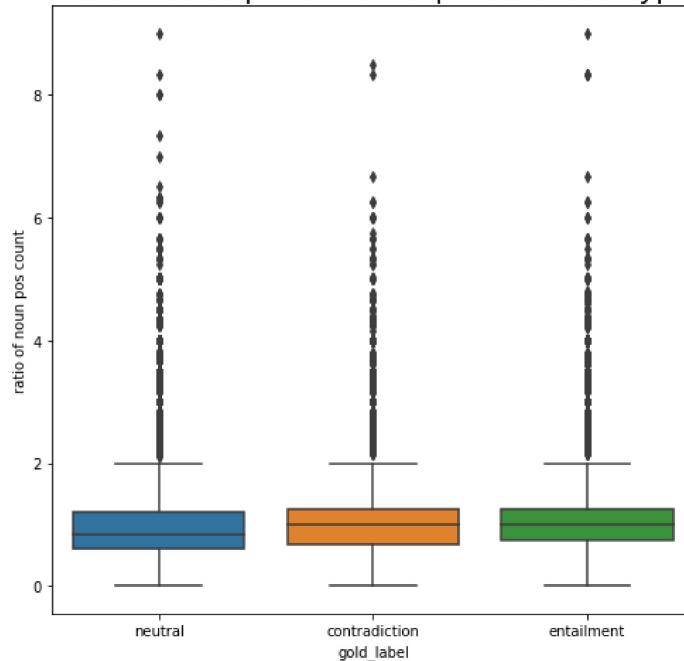
In [292]:

```

1 fig,ax = plt.subplots(1,1, figsize=(7,7), sharex=True)
2 sns.boxplot(x = train_df['gold_label'], y = train_df['ratio_noun'], ax = ax)
3 ax.set_title('box plot of ratio of count of noun pos between premise and hyp')
4 ax.set_ylabel('ratio of noun pos count')
5 plt.tight_layout()
6 plt.show()

```

box plot of ratio of count of noun pos between premise and hypothesis for each label



In [264]:

```

1 val_df['diff_noun'] = (val_df['s1_noun'] - val_df['s2_noun']).abs()
2 test_df['diff_noun'] = (test_df['s1_noun'] - test_df['s2_noun']).abs()

```

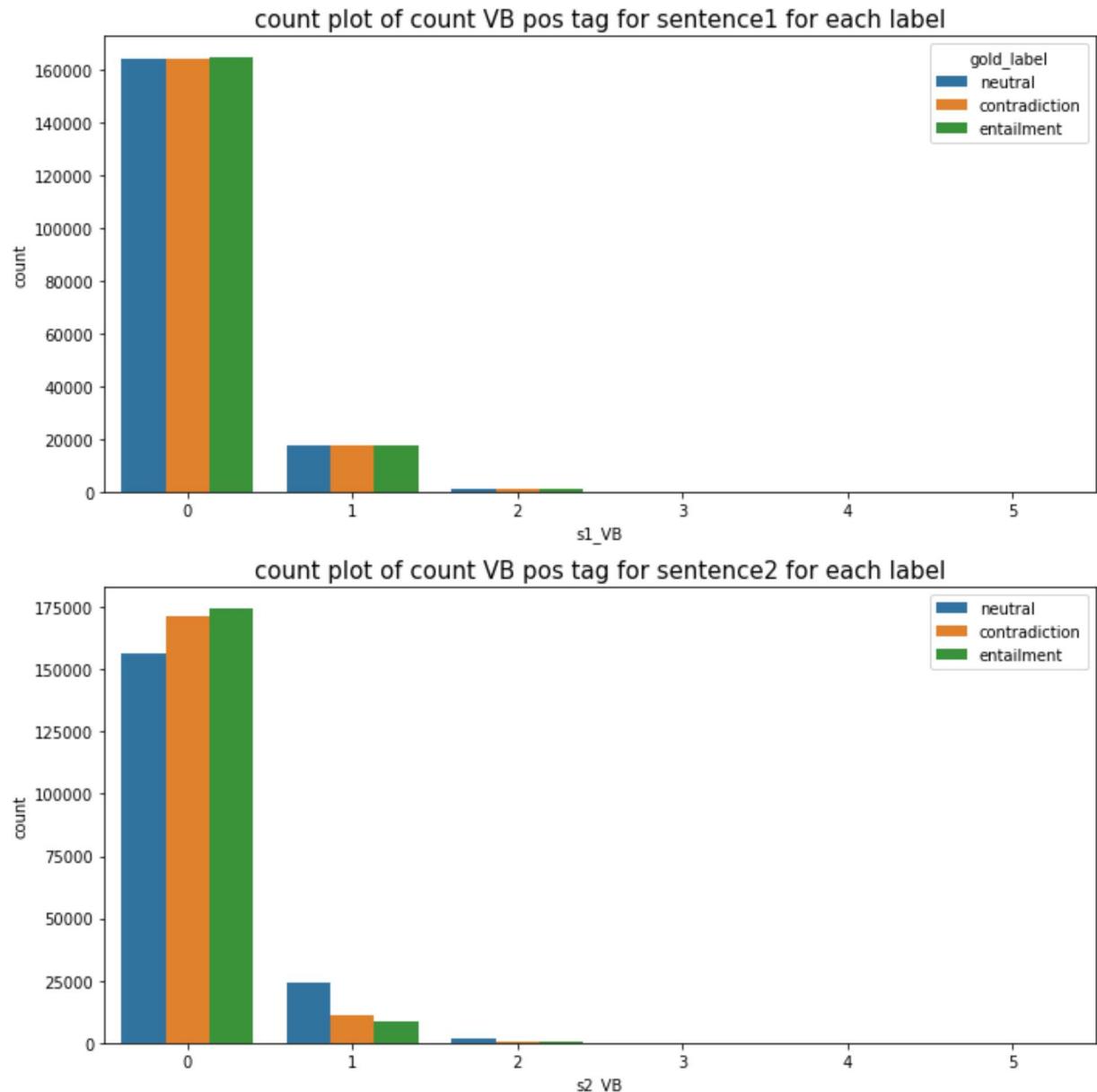
## VB: Verb

In [242]:

```

1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_VB'], hue = train_df['gold_label'], ax = ax[0]
3 sns.countplot(x = train_df['s2_VB'], hue = train_df['gold_label'], ax = ax[1]
4 ax[0].set_title('count plot of count VB pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count VB pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```



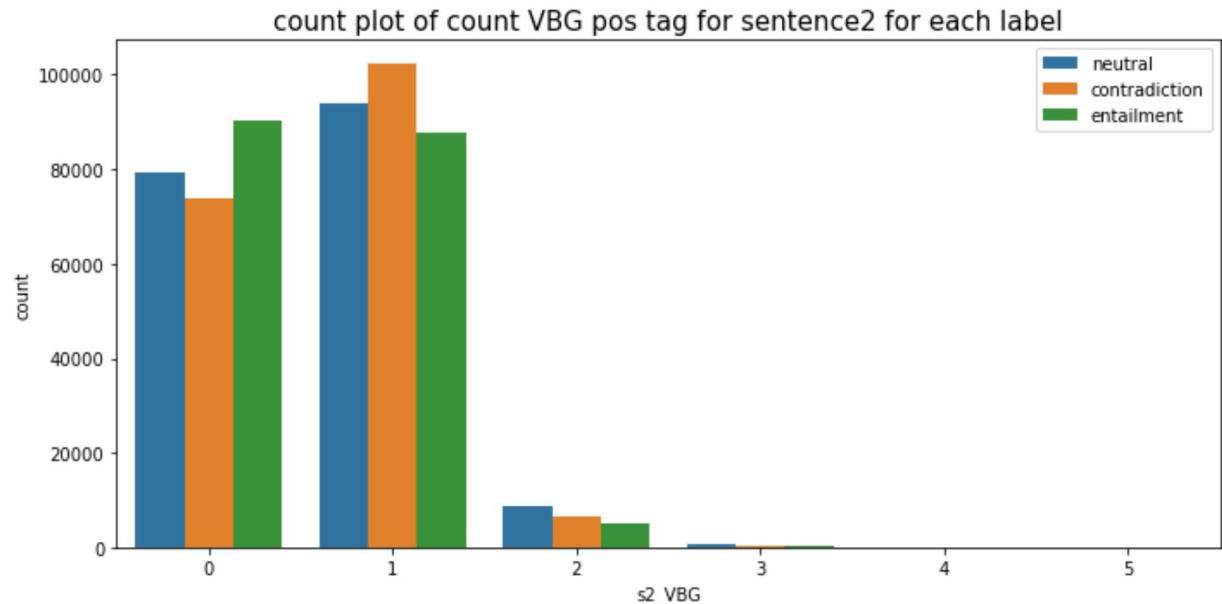
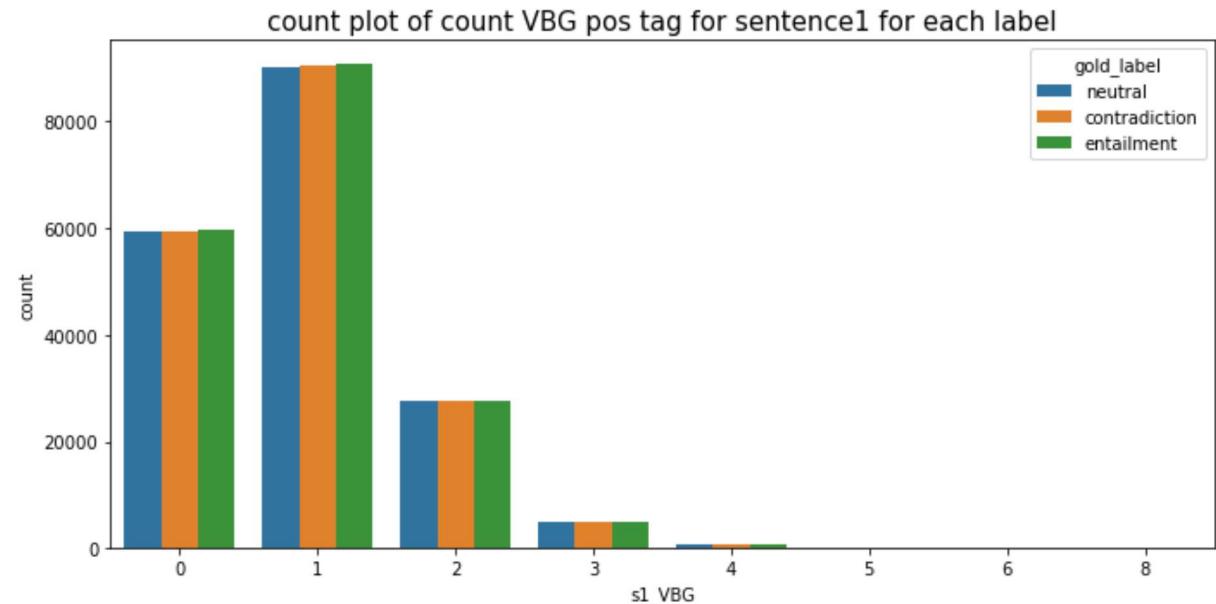
VB: verb, base form  
VBD: verb, past tense  
VBG: verb, present participle or gerund  
VBN: verb, past participle  
VBP: verb, present tense, not 3rd person singular  
VBZ: verb, present tense, 3rd person singular

In [245]:

```

1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_VBG'], hue = train_df['gold_label'], ax = ax[0])
3 sns.countplot(x = train_df['s2_VBG'], hue = train_df['gold_label'], ax = ax[1])
4 ax[0].set_title('count plot of count VBG pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count VBG pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```



In [255]:

```

1 train_df['s1_verb'] = train_df['s1_VB'] + train_df['s1_VBD'] + train_df['s1_VBN']
2 train_df['s2_verb'] = train_df['s2_VB'] + train_df['s2_VBD'] + train_df['s2_VBN']
3

```

In [259]:

```

1 val_df['s1_verb'] = val_df['s1_VB'] + val_df['s1_VBD'] + val_df['s1_VBG'] +
2 val_df['s2_verb'] = val_df['s2_VB'] + val_df['s2_VBD'] + val_df['s2_VBG'] +
3 test_df['s1_verb'] = test_df['s1_VB'] + test_df['s1_VBD'] + test_df['s1_VBG'] +
4 test_df['s2_verb'] = test_df['s2_VB'] + test_df['s2_VBD'] + test_df['s2_VBG']
5

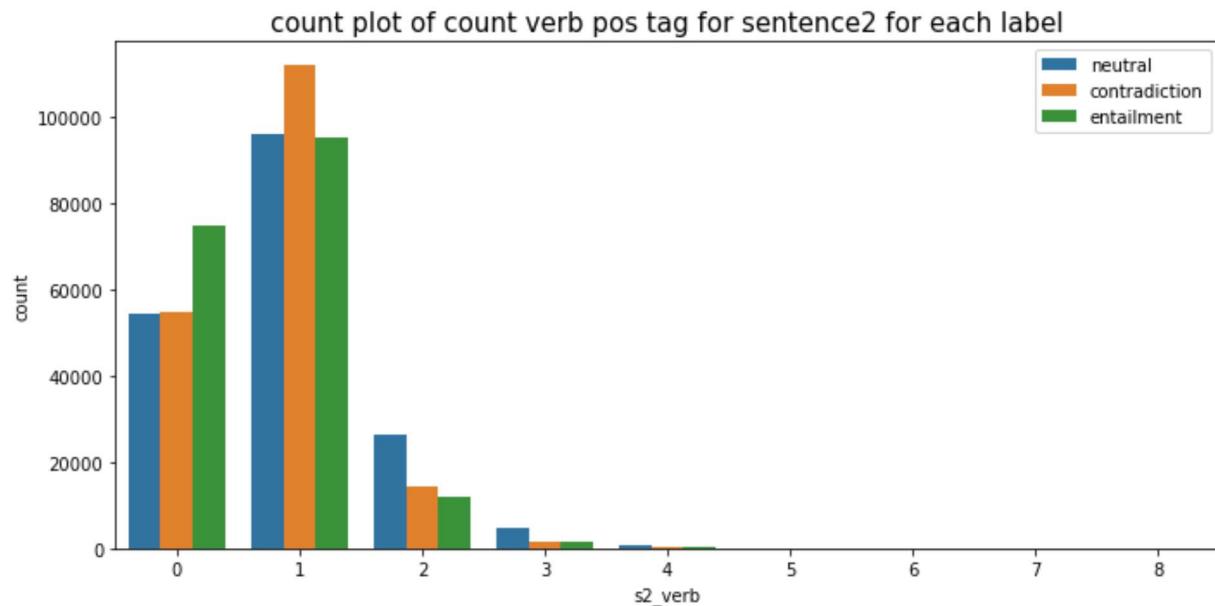
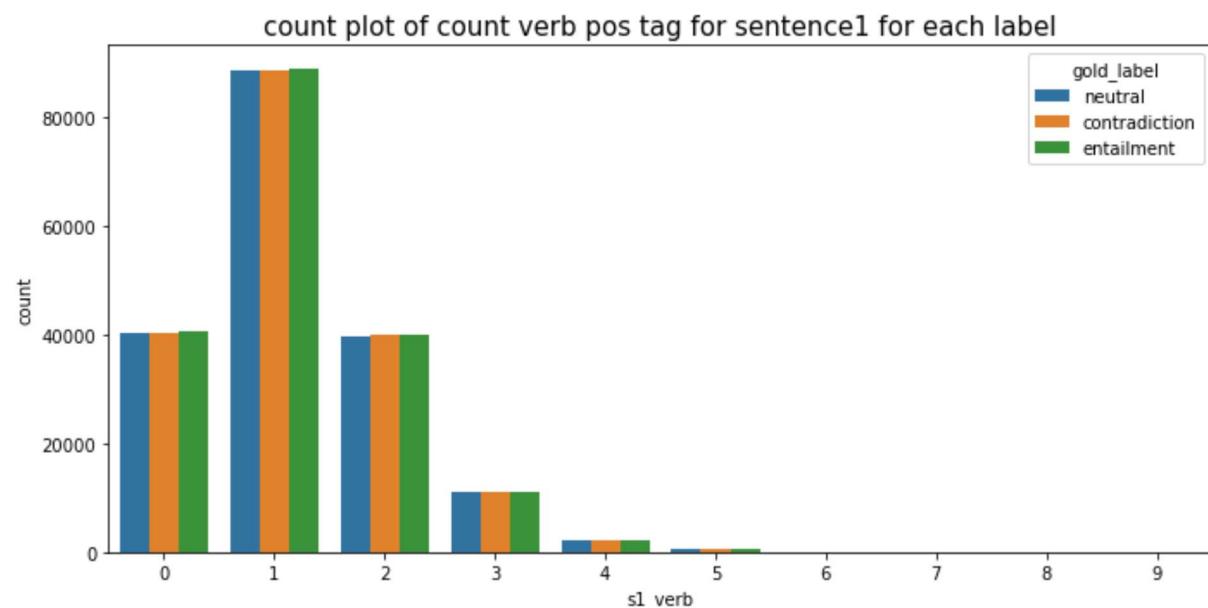
```

In [256]:

```

1 fig,ax = plt.subplots(2,1, figsize=(10,10))
2 sns.countplot(x = train_df['s1_verb'], hue = train_df['gold_label'], ax = ax)
3 sns.countplot(x = train_df['s2_verb'], hue = train_df['gold_label'], ax = ax)
4 ax[0].set_title('count plot of count verb pos tag for sentence1 for each label')
5 ax[0].set_ylabel('count')
6 ax[1].set_title('count plot of count verb pos tag for sentence2 for each label')
7 ax[1].set_ylabel('count')
8 plt.tight_layout()
9 plt.legend(loc = "upper right")
10 plt.show()

```



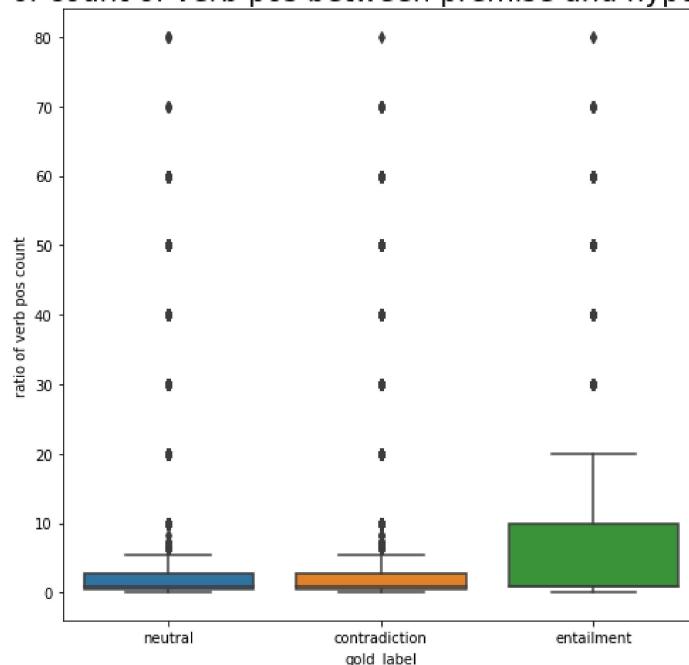
We can keep only summed up noun and verb pos tag counts and remove rest

```
In [266]: 1 train_df.to_csv('train_df.csv')
2 val_df.to_csv('val_df.csv')
3 test_df.to_csv('test_df.csv')
4
5
```

```
In [267]: 1 train_df = pd.read_csv('train_df.csv', index_col = 0)
2 val_df = pd.read_csv('val_df.csv', index_col = 0)
3 test_df = pd.read_csv('test_df.csv', index_col = 0)
```

```
In [290]: 1 train_df['ratio_verb'] = (train_df['s1_verb']/(train_df['s2_verb']+1e-1))
2
3 fig,ax = plt.subplots(1,1, figsize=(7,7), sharex=True)
4 sns.boxplot(x = train_df['gold_label'], y = train_df['ratio_verb'], ax = ax)
5 ax.set_title('box plot of ratio of count of verb pos between premise and hypothesis for each label')
6 ax.set_ylabel('ratio of verb pos count')
7 plt.tight_layout()
8 plt.show()
```

box plot of ratio of count of verb pos between premise and hypothesis for each label



```
In [297]: 1 val_df['ratio_verb'] = (val_df['s1_verb']/(val_df['s2_verb']+1e-1))
```

```
In [298]: 1 test_df['ratio_verb'] = (test_df['s1_verb']/(test_df['s2_verb']+1e-1))
```

Ratios of noun and verb pos count seem to be adding some discriminative power

In [305]: 1 train\_df.columns

```
Out[305]: Index(['sentence1', 'sentence2', 'gold_label', 'sentence1_preprocessed',
       'sentence2_preprocessed', 'sentence1_num_words', 'sentence2_num_words',
       'bleu_score', 'log_bleu_score', 'sentence1_polarity',
       'sentence2_polarity', 'sentence1_subjectivity',
       'sentence2_subjectivity', 'sum_polarity', 'diff_subjectivity',
       'cosine_similarity', 'tfidf_cosine_sim', 'diff_num_words', 's1_CC',
       's1_IN', 's1_JJ', 's1_JJR', 's1_JJS', 's1_NN', 's1_NNP', 's1_NNS',
       's1_PRP', 's1_PRP$', 's1_RB', 's1_RBR', 's1_RBS', 's1_UH', 's1_VB',
       's1_VBD', 's1_VBG', 's1_VBN', 's1_VBP', 's1_VBZ', 's2_CC', 's2_IN',
       's2_JJ', 's2_JJR', 's2_JJS', 's2_NN', 's2_NNP', 's2_NNS', 's2_PRP',
       's2_PRP$', 's2_RB', 's2_RBR', 's2_RBS', 's2_UH', 's2_VB', 's2_VBD',
       's2_VBG', 's2_VBN', 's2_VBP', 's2_VBZ', 's1_noun', 's2_noun', 's1_verb',
       's2_verb', 'ratio_noun', 'ratio_verb'],
      dtype='object')
```

## JJ: Adjective

In [306]: 1 # Let us add all adjective pos tags  
2 train\_df['s1\_adjective'] = train\_df['s1\_JJ'] + train\_df['s1\_JJR'] + train\_df['s1\_JJS']  
3 train\_df['s2\_adjective'] = train\_df['s2\_JJ'] + train\_df['s2\_JJR'] + train\_df['s2\_JJS']

In [320]: 1 val\_df['s1\_adjective'] = val\_df['s1\_JJ'] + val\_df['s1\_JJR'] + val\_df['s1\_JJS']  
2 val\_df['s2\_adjective'] = val\_df['s2\_JJ'] + val\_df['s2\_JJR'] + val\_df['s2\_JJS']  
3 test\_df['s1\_adjective'] = test\_df['s1\_JJ'] + test\_df['s1\_JJR'] + test\_df['s1\_JJS']  
4 test\_df['s2\_adjective'] = test\_df['s2\_JJ'] + test\_df['s2\_JJR'] + test\_df['s2\_JJS']

In [317]: 1 train\_df['diff\_adjective'] = (train\_df['s1\_adjective']-train\_df['s2\_adjective'])

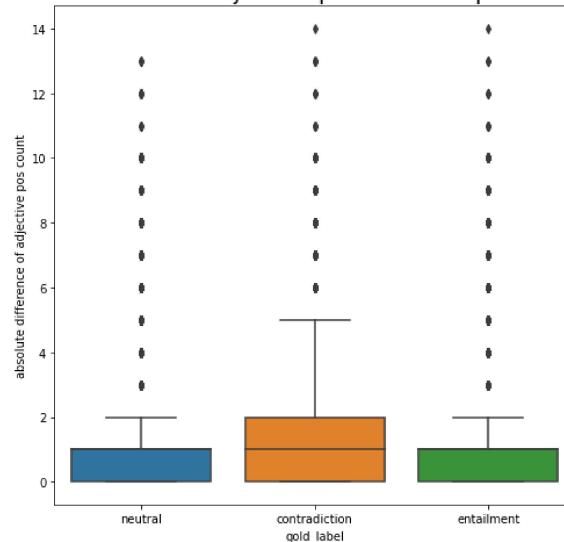
In [318]:

```

1 fig,ax = plt.subplots(1,1, figsize=(7,7), sharex=True)
2 sns.boxplot(x = train_df['gold_label'], y = train_df['diff_adjective'], ax =
3 ax.set_title('box plot of absolute difference of count of adjective pos betw
4 ax.set_ylabel('absolute difference of adjective pos count')
5 plt.tight_layout()
6 plt.show()

```

box plot of absolute difference of count of adjective pos between premise and hypothesis for each label



In [321]:

```

1 val_df['diff_adjective'] = (val_df['s1_adjective']-val_df['s2_adjective']).a
2 test_df['diff_adjective'] = (test_df['s1_adjective']-test_df['s2_adjective'])

```

In [324]:

```
1 train_df.columns
```

Out[324]: Index(['sentence1', 'sentence2', 'gold\_label', 'sentence1\_preprocessed', 'sentence2\_preprocessed', 'sentence1\_num\_words', 'sentence2\_num\_words', 'bleu\_score', 'log\_bleu\_score', 'sentence1\_polarity', 'sentence2\_polarity', 'sentence1\_subjectivity', 'sentence2\_subjectivity', 'sum\_polarity', 'diff\_subjectivity', 'cosine\_similarity', 'tfidf\_cosine\_sim', 'diff\_num\_words', 's1\_CC', 's1\_IN', 's1\_JJ', 's1\_JJR', 's1\_JJS', 's1\_NN', 's1\_NNP', 's1\_NNS', 's1\_PRP', 's1\_PRP\$', 's1\_RB', 's1\_RBR', 's1\_RBS', 's1\_UH', 's1\_VB', 's1\_VBD', 's1\_VBG', 's1\_VBN', 's1\_VBP', 's1\_VBZ', 's2\_CC', 's2\_IN', 's2\_JJ', 's2\_JJR', 's2\_JJS', 's2\_NN', 's2\_NNP', 's2\_NNS', 's2\_PRP', 's2\_PRP\$', 's2\_RB', 's2\_RBR', 's2\_RBS', 's2\_UH', 's2\_VB', 's2\_VBD', 's2\_VBG', 's2\_VBN', 's2\_VBP', 's2\_VBZ', 's1\_noun', 's2\_noun', 's1\_verb', 's2\_verb', 'ratio\_noun', 'ratio\_verb', 's1\_adjective', 's2\_adjective', 'diff\_adjective'],  
dtype='object')

## Levenshtein Distance

how many transformations you need to perform on the string A to make it equal to string B

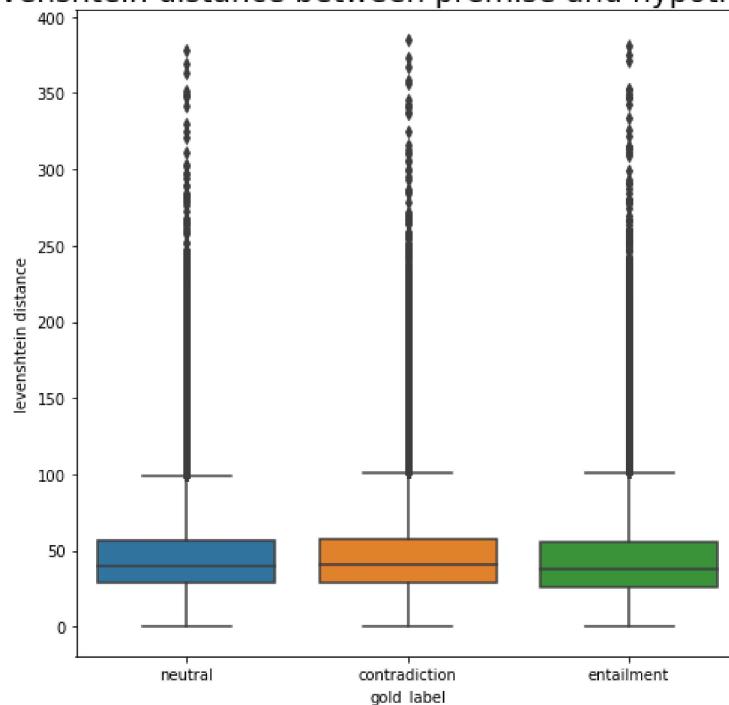
```
In [329]: 1 def get_levenshtein(s1, s2):
2     d = Levenshtein.distance(s1,s2)
3     return d
4
```

```
In [332]: 1 levelshtein_dist = list()
2 for i, row in train_df.iterrows():
3     d = get_levenshtein(row['sentence1'], row['sentence2'])
4     levelshtein_dist.append(d)
```

```
In [333]: 1 train_df['levenshtein_dist'] = levelshtein_dist
```

```
In [335]: 1 fig,ax = plt.subplots(1,1, figsize=(7,7))
2 sns.boxplot(x = train_df['gold_label'], y = train_df['levenshtein_dist'], ax=ax)
3 ax.set_title('box plot of levenshtein distance between premise and hypothesis')
4 ax.set_ylabel('levenshtein distance')
5 plt.tight_layout()
6 plt.show()
```

box plot of levenshtein distance between premise and hypothesis for each label



```
In [341]: 1 levelshtein_dist = list()
2 for i, row in val_df.iterrows():
3     d = get_levenshtein(row['sentence1'], row['sentence2'])
4     levelshtein_dist.append(d)
```

```
In [342]: 1 val_df['levenshtein_dist'] = levelshtein_dist
```

```
In [343]: 1 levelshtein_dist = list()
2 for i, row in test_df.iterrows():
3     d = get_levenshtein(row['sentence1'], row['sentence2'])
4     levelshtein_dist.append(d)
```

```
In [344]: 1 test_df['levenshtein_dist'] = levelshtein_dist
```

```
In [347]: 1 train_df.columns
```

```
Out[347]: Index(['sentence1', 'sentence2', 'gold_label', 'sentence1_preprocessed',
       'sentence2_preprocessed', 'sentence1_num_words', 'sentence2_num_words',
       'bleu_score', 'log_bleu_score', 'sentence1_polarity',
       'sentence2_polarity', 'sentence1_subjectivity',
       'sentence2_subjectivity', 'sum_polarity', 'diff_subjectivity',
       'cosine_similarity', 'tfidf_cosine_sim', 'diff_num_words', 's1_CC',
       's1_IN', 's1_JJ', 's1_JJR', 's1_JJS', 's1_NN', 's1_NNP', 's1_NNS',
       's1_PRP', 's1_PRP$', 's1_RB', 's1_RBR', 's1_RBS', 's1_UH', 's1_VB',
       's1_VBD', 's1_VBG', 's1_VBN', 's1_VBP', 's1_VBZ', 's2_CC', 's2_IN',
       's2_JJ', 's2_JJR', 's2_JJS', 's2_NN', 's2_NNP', 's2_NNS', 's2_PRP',
       's2_PRP$', 's2_RB', 's2_RBR', 's2_RBS', 's2_UH', 's2_VB', 's2_VBD',
       's2_VBG', 's2_VBN', 's2_VBP', 's2_VBZ', 's1_noun', 's2_noun', 's1_verb',
       's2_verb', 'ratio_noun', 'ratio_verb', 's1_adjective', 's2_adjective',
       'diff_adjective', 'levenshtein_dist'],
      dtype='object')
```

```
In [361]: 1 train_df.to_csv('train_df.csv')
2 val_df.to_csv('val_df.csv')
3 test_df.to_csv('test_df.csv')
```

```
In [362]: 1 train_df = pd.read_csv('train_df.csv', index_col = 0)
2 val_df = pd.read_csv('val_df.csv', index_col = 0)
3 test_df = pd.read_csv('test_df.csv', index_col = 0)
```

```
In [349]: 1 to_keep_x = ['sentence1_num_words', 'sentence2_num_words', 'log_bleu_score',
       'cosine_similarity', 'tfidf_cosine_sim', 'diff_num_words', 's1_CC',
       's1_IN', 's1_JJ', 's1_JJR', 's1_JJS', 's1_PRP', 's1_PRP$', 's1_RB',
       's2_CC', 's2_IN', 's2_JJ', 's2_JJR', 's2_JJS', 's2_NN', 's2_NNP',
       's2_NNS', 's2_PRP', 's2_PRP$', 's2_RB', 's2_RBR', 's2_RBS', 's2_UH',
       's2_VB', 's2_VBD', 's2_VBG', 's2_VBN', 's2_VBP', 's2_VBZ', 'ratio_noun',
       'ratio_verb', 's1_adjective', 's2_adjective', 'diff_adjective',
       'levenshtein_dist']
```

```
In [350]: 1 X_train = train_df[to_keep_x]
2 y_train = train_df['gold_label']
3 X_val = val_df[to_keep_x]
4 y_val = val_df['gold_label']
```

```
In [ ]: 1
```

In [353]:

```

1 # hyperparameter tuning
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import RandomizedSearchCV
4 # Number of trees in random forest
5 n_estimators = [100,200,500,600,800]
6 # Number of features to consider at every split
7 max_features = ['auto']
8 # Maximum number of Levels in tree
9 max_depth = [5,10,20,40,80]
10 # Minimum number of samples required to split a node
11 min_samples_split = [5, 10, 20, 40]
12 # Minimum number of samples required at each Leaf node
13 min_samples_leaf = [2, 4, 10]
14 # Method of selecting samples for training each tree
15 bootstrap = [True]
16 # Create the random grid
17 random_grid = {'n_estimators': n_estimators,
18                 'max_features': max_features,
19                 'max_depth': max_depth,
20                 'min_samples_split': min_samples_split,
21                 'min_samples_leaf': min_samples_leaf,
22                 'bootstrap': bootstrap}
23 # Use the random grid to search for best hyperparameters
24 # First create the base model to tune
25 rf = RandomForestClassifier()
26 # Random search of parameters, using 3 fold cross validation,
27 # search across 100 different combinations, and use all available cores
28 rf_random = RandomizedSearchCV(estimator = rf, param_distributions = random_
29 # Fit the random search model
30 rf_random.fit(X_train, y_train)

```

Out[353]: RandomizedSearchCV(cv=3, estimator=RandomForestClassifier(), n\_iter=50,

```

n_jobs=6,
param_distributions={'bootstrap': [True],
                     'max_depth': [5, 10, 20, 40, 80],
                     'max_features': ['auto'],
                     'min_samples_leaf': [2, 4, 10],
                     'min_samples_split': [5, 10, 20, 40],
                     'n_estimators': [100, 200, 500, 600,
                                     800]},
random_state=42)

```

In [354]: 1 rf\_random.best\_estimator\_

Out[354]: RandomForestClassifier(max\_depth=80, min\_samples\_leaf=2, min\_samples\_split=10, n\_estimators=800)

In [355]: 1 rf\_random.best\_score\_

Out[355]: 0.559056795693159

```
In [356]: 1 model = RandomForestClassifier(max_depth=80, min_samples_leaf=2, min_samples_
2                               n_estimators=800)
3 model.fit(X_train, y_train)
```

```
Out[356]: RandomForestClassifier(max_depth=80, min_samples_leaf=2, min_samples_split=10,
                                   n_estimators=800)
```

```
In [357]: 1 val_pred = model.predict(X_val)
```

```
In [359]: 1 from sklearn.metrics import accuracy_score
2 acc = accuracy_score(y_val, val_pred)
```

```
In [360]: 1 acc
```

```
Out[360]: 0.5722414143466775
```

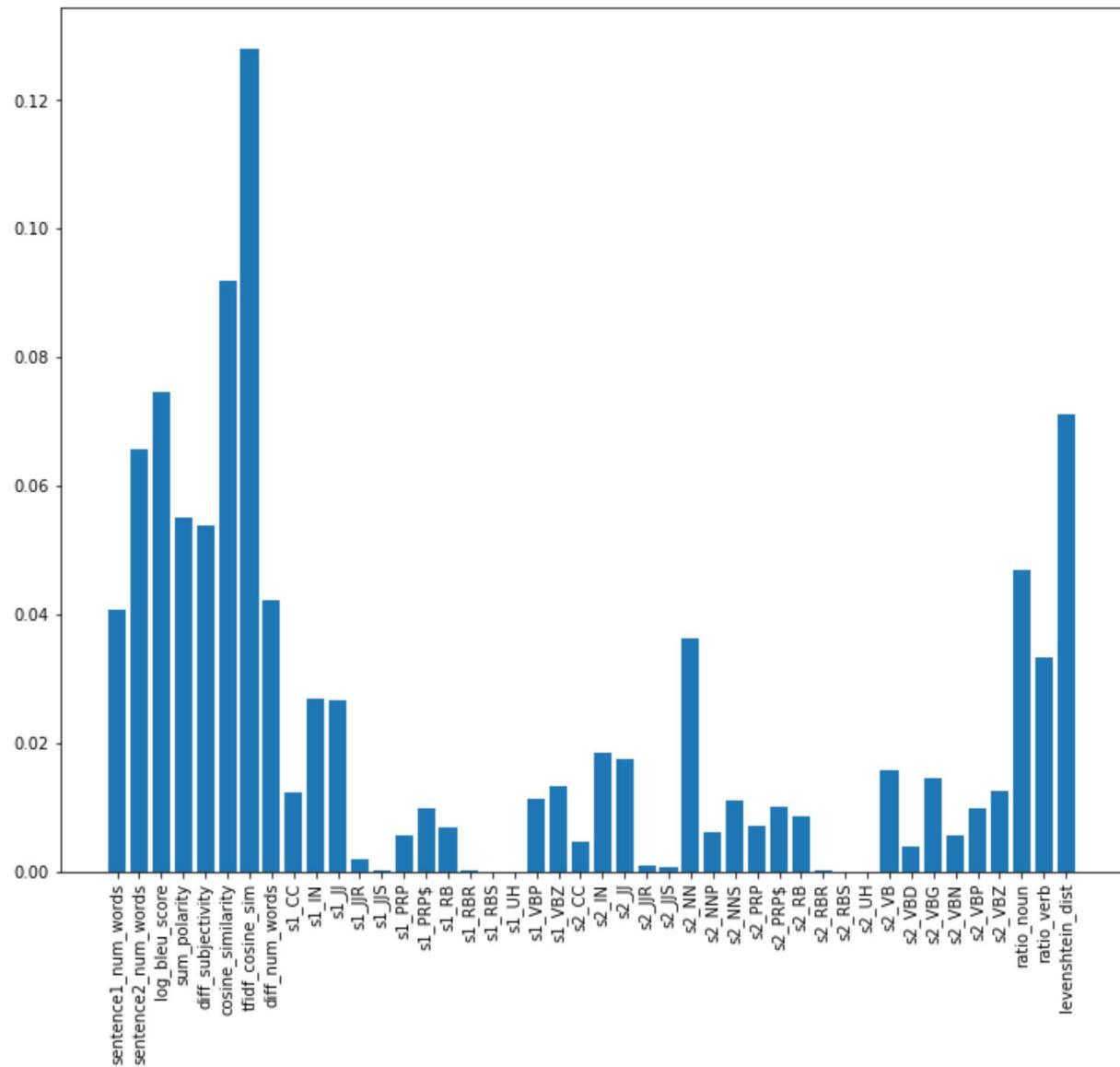
```
In [364]: 1 importance = model.feature_importances_
```

In [365]:

```

1 # Plotting feature importance
2 columns = X_train.columns
3 x_pos = np.arange(len(columns))
4 plt.figure(figsize = (12,10))
5 plt.bar(x_pos, importance)
6 plt.xticks(x_pos, columns, rotation='vertical')
7 plt.show()
8 #plt.plot(np.array(range(0,47)), importance, 'bo' )

```



In [ ]:

1