

```
1 import os
2 import random
3 from google.colab import drive
4 import json
```

```
1 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True)

```
1 files = []
2 for file in os.listdir('/content/drive/MyDrive/hindi_corpus/'):
3     if '.dat' in file:
4         files.append(file)
```

```
1 print("Total number of files used =", len(files))
```

Total number of files used = 1058

```
1 tagged_Sentences = []
2 curr = []
3 for f in files:
4     with open('/content/drive/MyDrive/hindi_corpus/' + f, 'r') as file:
5         for i in file.readlines():
6             i = i.split("\t")
7             if len(i) > 5:
8                 curr.append((i[1], i[4]))
9             else:
10                 tagged_Sentences.append(curr)
11                 curr = []
12 print("Total number of sentences =", len(tagged_Sentences))
```

☞ Total number of sentences = 18685

```
1 def creating_bigram_sents(current_sents):
2     bigram_sents = []
```

```

2 bigram_sents = []
3 for idx in range(len(current_sents) - 1):
4     bigram_sents.append([(current_sents[idx][0], current_sents[idx][1]), (current_sents[idx + 1][0], current_sents[idx + 1][1])])
5 return bigram_sents

```

```

1 def clean_sents(tagged_sentences):
2     final_sents = []
3     for sents in tagged_sentences:
4         new_sents = []
5         for word, tag in sents:
6             if tag == "SYM" or tag == '' or tag == 'PUNC':
7                 continue
8             new_sents.append((word, tag.lower()))
9         final_sents.extend(creating_bigram_sents(new_sents))
10    return final_sents

```

```

1 tagged_sentences = clean_sents(tagged_Sentences)
2 random.shuffle(tagged_sentences)
3 print("Total number of sentences after cleaning =", len(tagged_sentences))

```

Total number of sentences after cleaning = 331364

```
1 print(tagged_sentences[0])
```

[('आग', 'nn'), ('की', 'psp')]

```

1 def removing_tags(tagged_sentences):
2     return [word for word, tag in tagged_sentences]

```

```
1 removed = [removing_tags(i) for i in tagged_sentences]
```

```

1 tags = []
2 for i in tagged_sentences:
3     tags.extend([i[0][1], i[1][1]])

```

```
1 with open ('/content/drive/MyDrive/Cleaned_Sentences_Task/cleaned_sentences.txt', 'w+') as file:
```

```
2     json.dump(tagged_sentences, file)
```

```
1 with open ('/content/drive/MyDrive/Cleaned_Sentences_Task/sentences_without_tags.txt', 'w+') as file:  
2     json.dump(removed, file)
```

```
1 with open ('/content/drive/MyDrive/Cleaned_Sentences_Task/tags_original.txt', 'w+') as file:  
2     json.dump(tags, file)
```