

```
In [1]: from nltk.tokenize import word_tokenize
        from nltk.text import Text
        from nltk.util import ngrams
        import pickle
```

```
In [2]: #function takes a filename as an argument
def preprocess(filename):
    with open(filename, 'r', encoding = 'utf8') as f:
        raw_text = f.read()
        raw_text = raw_text.replace('\n', ' ')

    tokens = word_tokenize(raw_text)

    bigrams = list(ngrams(tokens, 2))

    unigrams = list(ngrams(tokens, 1))

    bigram_dict = {b:bigrams.count(b) for b in set(bigrams)}
    unigram_dict = {u:unigrams.count(u) for u in set(unigrams)}

    #use the bigram list to create a bigram dictionary of bigrams and counts, ['token1
    countb = 1
    for element in bigram_dict.keys():
        print(element, '->', bigram_dict[element])
        countb += 1
        if countb > 5:
            break

    #use the unigram list to create a unigram dictionary of bigrams and counts, ['token
    countu = 1
    for element in unigram_dict.keys():
        print(element, '->', unigram_dict[element])
        countu += 1
        if countu > 5:
            break

    return unigram_dict, bigram_dict
```

```
In [3]: def main():
        #preprocess the text
        E_Uni, E_Bi = preprocess("LangId.train.English")
        F_Uni, F_Bi = preprocess("LangId.train.French")
        I_Uni, I_Bi = preprocess("LangId.train.Italian")

        #pickle the files
        pickle.dump(E_Uni, open('E_Uni.pickle', 'wb'))
        pickle.dump(E_Bi, open('E_Bi.pickle', 'wb'))

        pickle.dump(F_Uni, open('F_Uni.pickle', 'wb'))
        pickle.dump(F_Bi, open('F_Bi.pickle', 'wb'))

        pickle.dump(I_Uni, open('I_Uni.pickle', 'wb'))
        pickle.dump(I_Bi, open('I_Bi.pickle', 'wb'))
```

```
In [4]: if __name__ == "__main__":  
        main()
```

```
('approach', 'the') -> 1  
( 'child', '-' ) -> 1  
( 's', 'representative' ) -> 1  
( 'for', 'drafting' ) -> 1  
( 'know', '.' ) -> 1  
( 'decade', ) -> 2  
( '-', ) -> 819  
( 'Potential', ) -> 1  
( 'judgement', ) -> 7  
( 'modal', ) -> 1  
( 'en', 'vanter' ) -> 1  
( 'spécialisé', ',' ) -> 1  
( 'santé', 'doit' ) -> 1  
( 'pas', 'servir' ) -> 2  
( 'à', 'rien' ) -> 2  
( 'statistiques', ) -> 5  
( 'Indiens', ) -> 1  
( 'Bassin', ) -> 1  
( 'importe', ) -> 5  
( '-', ) -> 2359  
( 'tecnica', ',' ) -> 1  
( 'essere', 'esteso' ) -> 1  
( 'abbandono', ',' ) -> 1  
( 'scapito', 'dell' ) -> 1  
( 'potremmo', 'fornire' ) -> 1  
( 'ottimo', ) -> 4  
( 'partendo', ) -> 2  
( 'accolta', ) -> 1  
( 'Consentitemi', ) -> 2  
( 'riportare', ) -> 1
```

```
In [ ]:
```