

```
1 import nltk
2 import pickle
3 import argparse
4 from collections import Counter
5 from pycocotools.coco import COCO
6
7
8 class Vocabulary(object):
9     """Simple vocabulary wrapper."""
10     def __init__(self):
11         self.word2idx = {}
12         self.idx2word = {}
13         self.idx = 0
14
15     def add_word(self, word):
16         if not word in self.word2idx:
17             self.word2idx[word] = self.idx
18             self.idx2word[self.idx] = word
19             self.idx += 1
20
21     def __call__(self, word):
22         if not word in self.word2idx:
23             return self.word2idx['<unk>']
24         return self.word2idx[word]
25
26     def __len__(self):
27         return len(self.word2idx)
28
29 def build_vocab(json, threshold):
30     """Build a simple vocabulary wrapper."""
31     coco = COCO(json)
32     counter = Counter()
33     ids = coco.anns.keys()
34     for i, id in enumerate(ids):
35         caption = str(coco.anns[id]['caption'])
36         tokens = nltk.tokenize.word_tokenize(caption.lower())
37         counter.update(tokens)
38
39         if (i+1) % 1000 == 0:
40             print("[{}/{}] Tokenized the captions.".format(i+1, len(ids)))
41
42     # If the word frequency is less than 'threshold', then the word is
43     # discarded.
44     words = [word for word, cnt in counter.items() if cnt >= threshold]
45
46     # Create a vocab wrapper and add some special tokens.
47     vocab = Vocabulary()
48     vocab.add_word('<pad>')
49     vocab.add_word('<start>')
50     vocab.add_word('<end>')
51     vocab.add_word('<unk>')
52
53     # Add the words to the vocabulary.
54     for i, word in enumerate(words):
55         vocab.add_word(word)
56     return vocab
57
58 def main(args):
59     vocab = build_vocab(json=args.caption_path, threshold=args.threshold)
60     vocab_path = args.vocab_path
```

```
60     with open(vocab_path, 'wb') as f:
61         pickle.dump(vocab, f)
62     print("Total vocabulary size: {}".format(len(vocab)))
63     print("Saved the vocabulary wrapper to '{}'.format(vocab_path))
64
65
66 if __name__ == '__main__':
67     parser = argparse.ArgumentParser()
68     parser.add_argument('--caption_path', type=str,
69                          default='data/annotations/captions_train2014.json',
70                          help='path for train annotation file')
71     parser.add_argument('--vocab_path', type=str,
72                          default='./data/vocab.pkl',
73                          help='path for saving vocabulary wrapper')
74     parser.add_argument('--threshold', type=int, default=4,
75                          help='minimum word count threshold')
76     args = parser.parse_args()
77     main(args)
```