# Notebook

February 21, 2019

```
In [1]: # import modules & set up logging
        import gensim, logging
        import smart_open, os
        logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.
        import datetime
        import pandas as pd
        import multiprocessing

        # fichier incltu dans le projet
        import save_notebook
```

D:\Outil\Anaconda\envs\majeure-ml-env\lib\site-packages\gensim\utils.py:1197: UserWarning: det
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

# 1 DÃľclaration donnÃľes

```
In [2]: now = str(datetime.datetime.now()).replace(" ","")
```

```
In [3]: word_embedding = "word2vec"
```

# 2 Prepare data

```
In [4]: filenames =  os.listdir("../wikipedia/data")
```

```
In [5]: #CrÃľÃľ un fichier ou chaque ligne continent tout un fichier
        # path="../wikipedia/data"
        path="../wikipedia/data/"
        with open('./data/wikipedia_informatic.txt', 'w+',encoding="utf8" ) as out_file:

            for fname in filenames:
        #         print(fname)
                if "ipynb_checkpoints" in fname:
                    continue
                try:
                    with open(path + fname, encoding="utf8") as in_file:
                        out_file.write(in_file.read().replace("\n",""))
```

```
            except:
                continue
```

In [6]: 
```python
# On lit et on tokenize le fichier
with open('./data/wikipedia_informatic.txt', 'r', encoding="utf8") as f:
    wiki_vocab = f.readlines()
wiki_vocab = [x.strip() for x in wiki_vocab]

wiki_vocab_tokenized = []
# for line in wiki_vocab:
#     print(gensim.utils.simple_preprocess(line))
# wiki_vocab_tokenized.append(gensim.utils.simple_preprocess(str(wiki_vocab)))
```

In [7]: 
```python
wiki_vocab_tokenized = gensim.utils.simple_preprocess(str(wiki_vocab))
```

## 3 Create model

In [8]: 
```python
# build vocabulary and train model
model = gensim.models.Word2Vec(
    [wiki_vocab_tokenized],
    size=150,
    seed=1234,
    window=10,
    min_count=2,
    workers=multiprocessing.cpu_count())

date_before_learning = datetime.datetime.now()
model.train([wiki_vocab_tokenized], total_examples=len(wiki_vocab_tokenized), epochs=3(
time_training = datetime.datetime.now() - date_before_learning
```

```
2019-02-21 17:14:51,773 : WARNING : consider setting layer size to a multiple of 4 for greater
2019-02-21 17:14:51,775 : INFO : collecting all words and their counts
2019-02-21 17:14:51,775 : INFO : PROGRESS: at sentence #0, processed 0 words, keeping 0 word ty
2019-02-21 17:14:53,206 : INFO : collected 134060 word types from a corpus of 5386246 raw words
2019-02-21 17:14:53,208 : INFO : Loading a fresh vocabulary
2019-02-21 17:14:53,600 : INFO : min_count=2 retains 62658 unique words (46% of original 134060
2019-02-21 17:14:53,601 : INFO : min_count=2 leaves 5314844 word corpus (98% of original 538624
2019-02-21 17:14:53,838 : INFO : deleting the raw counts dictionary of 134060 items
2019-02-21 17:14:53,842 : INFO : sample=0.001 downsamples 28 most-common words
2019-02-21 17:14:53,844 : INFO : downsampling leaves estimated 4091074 word corpus (77.0% of pr
2019-02-21 17:14:54,078 : INFO : estimated required memory for 62658 words and 150 dimensions:
2019-02-21 17:14:54,079 : INFO : resetting layer weights
2019-02-21 17:14:55,254 : INFO : training model with 4 workers on 62658 vocabulary and 150 feat
2019-02-21 17:14:55,263 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,264 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,264 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,308 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,310 : INFO : EPOCH - 1 : training on 5386246 raw words (10000 effective wor
```

```
2019-02-21 17:14:55,318 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,320 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,321 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,354 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,355 : INFO : EPOCH - 2 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,364 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,366 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,367 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,397 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,398 : INFO : EPOCH - 3 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,409 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,410 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,439 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,441 : INFO : EPOCH - 4 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,448 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,450 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,451 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,480 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,480 : INFO : EPOCH - 5 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,481 : INFO : training on a 26931230 raw words (50000 effective words) took
2019-02-21 17:14:55,481 : WARNING : under 10 jobs per worker: consider setting a smaller `batc
2019-02-21 17:14:55,483 : WARNING : Effective 'alpha' higher than previous training cycles
2019-02-21 17:14:55,484 : INFO : training model with 4 workers on 62658 vocabulary and 150 fea
2019-02-21 17:14:55,495 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,500 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,501 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,527 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,528 : INFO : EPOCH - 1 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,529 : WARNING : EPOCH - 1 : supplied example count (1) did not equal expec
2019-02-21 17:14:55,542 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,542 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,543 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,571 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,572 : INFO : EPOCH - 2 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,573 : WARNING : EPOCH - 2 : supplied example count (1) did not equal expec
2019-02-21 17:14:55,582 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,583 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,584 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,612 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,613 : INFO : EPOCH - 3 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,614 : WARNING : EPOCH - 3 : supplied example count (1) did not equal expec
2019-02-21 17:14:55,623 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,624 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,625 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,652 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,653 : INFO : EPOCH - 4 : training on 5386246 raw words (10000 effective wor
2019-02-21 17:14:55,654 : WARNING : EPOCH - 4 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:55,663 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,664 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,665 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,695 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,696 : INFO : EPOCH - 5 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:55,697 : WARNING : EPOCH - 5 : supplied example count (1) did not equal expect
2019-02-21 17:14:55,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,708 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,709 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,741 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,742 : INFO : EPOCH - 6 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:55,743 : WARNING : EPOCH - 6 : supplied example count (1) did not equal expect
2019-02-21 17:14:55,754 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,762 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,763 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,788 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,789 : INFO : EPOCH - 7 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:55,789 : WARNING : EPOCH - 7 : supplied example count (1) did not equal expect
2019-02-21 17:14:55,801 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,802 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,803 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,833 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,834 : INFO : EPOCH - 8 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:55,835 : WARNING : EPOCH - 8 : supplied example count (1) did not equal expect
2019-02-21 17:14:55,846 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,847 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,848 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,880 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,881 : INFO : EPOCH - 9 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:55,881 : WARNING : EPOCH - 9 : supplied example count (1) did not equal expect
2019-02-21 17:14:55,891 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,893 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,894 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,925 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,926 : INFO : EPOCH - 10 : training on 5386246 raw words (10000 effective we
2019-02-21 17:14:55,926 : WARNING : EPOCH - 10 : supplied example count (1) did not equal expec
2019-02-21 17:14:55,937 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,938 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,939 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:55,970 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:55,971 : INFO : EPOCH - 11 : training on 5386246 raw words (10000 effective we
2019-02-21 17:14:55,972 : WARNING : EPOCH - 11 : supplied example count (1) did not equal expec
2019-02-21 17:14:55,985 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:55,987 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:55,988 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,019 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,028 : INFO : EPOCH - 12 : training on 5386246 raw words (10000 effective we
2019-02-21 17:14:56,029 : WARNING : EPOCH - 12 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:56,041 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,043 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,043 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,075 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,076 : INFO : EPOCH - 13 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,077 : WARNING : EPOCH - 13 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,088 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,089 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,090 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,120 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,121 : INFO : EPOCH - 14 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,122 : WARNING : EPOCH - 14 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,132 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,134 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,134 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,165 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,166 : INFO : EPOCH - 15 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,166 : WARNING : EPOCH - 15 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,177 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,179 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,180 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,212 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,213 : INFO : EPOCH - 16 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,214 : WARNING : EPOCH - 16 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,225 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,226 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,227 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,257 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,259 : INFO : EPOCH - 17 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,260 : WARNING : EPOCH - 17 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,270 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,273 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,273 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,304 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,305 : INFO : EPOCH - 18 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,306 : WARNING : EPOCH - 18 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,316 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,318 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,319 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,349 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,350 : INFO : EPOCH - 19 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,351 : WARNING : EPOCH - 19 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,362 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,363 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,364 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,395 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,396 : INFO : EPOCH - 20 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,396 : WARNING : EPOCH - 20 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:56,406 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,408 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,409 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,439 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,440 : INFO : EPOCH - 21 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,441 : WARNING : EPOCH - 21 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,450 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,451 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,452 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,483 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,485 : INFO : EPOCH - 22 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,486 : WARNING : EPOCH - 22 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,500 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,502 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,533 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,535 : INFO : EPOCH - 23 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,535 : WARNING : EPOCH - 23 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,546 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,547 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,548 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,579 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,580 : INFO : EPOCH - 24 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,581 : WARNING : EPOCH - 24 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,592 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,593 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,593 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,624 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,625 : INFO : EPOCH - 25 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,626 : WARNING : EPOCH - 25 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,638 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,639 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,640 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,670 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,671 : INFO : EPOCH - 26 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,672 : WARNING : EPOCH - 26 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,683 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,684 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,685 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,715 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,716 : INFO : EPOCH - 27 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,717 : WARNING : EPOCH - 27 : supplied example count (1) did not equal expec
2019-02-21 17:14:56,756 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,764 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,765 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,765 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,766 : INFO : EPOCH - 28 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,768 : WARNING : EPOCH - 28 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:56,780 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,782 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,783 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,810 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,811 : INFO : EPOCH - 29 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,812 : WARNING : EPOCH - 29 : supplied example count (1) did not equal expe
2019-02-21 17:14:56,827 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,836 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,837 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,858 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,858 : INFO : EPOCH - 30 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,859 : WARNING : EPOCH - 30 : supplied example count (1) did not equal expe
2019-02-21 17:14:56,870 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,870 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,871 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,901 : INFO : EPOCH - 31 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,902 : WARNING : EPOCH - 31 : supplied example count (1) did not equal expe
2019-02-21 17:14:56,912 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,913 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,914 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,943 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,944 : INFO : EPOCH - 32 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,945 : WARNING : EPOCH - 32 : supplied example count (1) did not equal expe
2019-02-21 17:14:56,957 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:56,958 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:56,958 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:56,989 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:56,990 : INFO : EPOCH - 33 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:56,991 : WARNING : EPOCH - 33 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,001 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,002 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,003 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,033 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,034 : INFO : EPOCH - 34 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,035 : WARNING : EPOCH - 34 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,046 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,047 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,048 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,078 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,079 : INFO : EPOCH - 35 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,080 : WARNING : EPOCH - 35 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,090 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,092 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,093 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,126 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,127 : INFO : EPOCH - 36 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,127 : WARNING : EPOCH - 36 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:14:57,139 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,140 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,141 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,172 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,173 : INFO : EPOCH - 37 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,174 : WARNING : EPOCH - 37 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,189 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,190 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,192 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,222 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,223 : INFO : EPOCH - 38 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,224 : WARNING : EPOCH - 38 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,234 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,235 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,236 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,265 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,267 : INFO : EPOCH - 39 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,267 : WARNING : EPOCH - 39 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,278 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,279 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,280 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,310 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,311 : INFO : EPOCH - 40 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,312 : WARNING : EPOCH - 40 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,325 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,326 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,357 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,358 : INFO : EPOCH - 41 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,359 : WARNING : EPOCH - 41 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,370 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,383 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,384 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,403 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,404 : INFO : EPOCH - 42 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,404 : WARNING : EPOCH - 42 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,417 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,418 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,419 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,449 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,450 : INFO : EPOCH - 43 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,451 : WARNING : EPOCH - 43 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,462 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,463 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,463 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,493 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,494 : INFO : EPOCH - 44 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,495 : WARNING : EPOCH - 44 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:57,506 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,507 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,508 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,539 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,541 : INFO : EPOCH - 45 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,541 : WARNING : EPOCH - 45 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,552 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,553 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,554 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,583 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,584 : INFO : EPOCH - 46 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,585 : WARNING : EPOCH - 46 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,597 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,597 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,598 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,629 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,630 : INFO : EPOCH - 47 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,631 : WARNING : EPOCH - 47 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,644 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,653 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,654 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,677 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,678 : INFO : EPOCH - 48 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,679 : WARNING : EPOCH - 48 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,692 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,693 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,724 : INFO : EPOCH - 49 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,725 : WARNING : EPOCH - 49 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,737 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,738 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,739 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,770 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,771 : INFO : EPOCH - 50 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,772 : WARNING : EPOCH - 50 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,782 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,783 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,784 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,814 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,816 : INFO : EPOCH - 51 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,816 : WARNING : EPOCH - 51 : supplied example count (1) did not equal expe
2019-02-21 17:14:57,826 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,827 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,828 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,857 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,859 : INFO : EPOCH - 52 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,859 : WARNING : EPOCH - 52 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:14:57,870 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,871 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,871 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,900 : INFO : EPOCH - 53 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,900 : WARNING : EPOCH - 53 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,911 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,912 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,912 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,940 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,941 : INFO : EPOCH - 54 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,942 : WARNING : EPOCH - 54 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,951 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,953 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,954 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:57,981 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:57,982 : INFO : EPOCH - 55 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:57,983 : WARNING : EPOCH - 55 : supplied example count (1) did not equal expec
2019-02-21 17:14:57,992 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:57,993 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:57,994 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,022 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,023 : INFO : EPOCH - 56 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,023 : WARNING : EPOCH - 56 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,034 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,065 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,066 : INFO : EPOCH - 57 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,067 : WARNING : EPOCH - 57 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,077 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,078 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,079 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,109 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,110 : INFO : EPOCH - 58 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,111 : WARNING : EPOCH - 58 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,123 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,124 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,125 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,155 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,156 : INFO : EPOCH - 59 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,157 : WARNING : EPOCH - 59 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,171 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,172 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,173 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,207 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,208 : INFO : EPOCH - 60 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,209 : WARNING : EPOCH - 60 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:58,219 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,221 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,222 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,252 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,253 : INFO : EPOCH - 61 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,254 : WARNING : EPOCH - 61 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,262 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,263 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,264 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,292 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,293 : INFO : EPOCH - 62 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,293 : WARNING : EPOCH - 62 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,305 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,306 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,307 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,332 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,333 : INFO : EPOCH - 63 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,335 : WARNING : EPOCH - 63 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,343 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,344 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,345 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,372 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,374 : INFO : EPOCH - 64 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,374 : WARNING : EPOCH - 64 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,386 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,387 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,412 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,431 : INFO : EPOCH - 65 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,432 : WARNING : EPOCH - 65 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,445 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,447 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,448 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,474 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,476 : INFO : EPOCH - 66 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,476 : WARNING : EPOCH - 66 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,493 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,495 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,496 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,516 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,516 : INFO : EPOCH - 67 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,517 : WARNING : EPOCH - 67 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,526 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,527 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,527 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,555 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,556 : INFO : EPOCH - 68 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,557 : WARNING : EPOCH - 68 : supplied example count (1) did not equal expec
```

11

```
2019-02-21 17:14:58,566 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,567 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,568 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,596 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,597 : INFO : EPOCH - 69 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,598 : WARNING : EPOCH - 69 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,609 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,610 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,611 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,641 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,643 : INFO : EPOCH - 70 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,644 : WARNING : EPOCH - 70 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,655 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,656 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,657 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,687 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,688 : INFO : EPOCH - 71 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,689 : WARNING : EPOCH - 71 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,701 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,702 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,703 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,733 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,733 : INFO : EPOCH - 72 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,734 : WARNING : EPOCH - 72 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,745 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,749 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,752 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,780 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,782 : INFO : EPOCH - 73 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,782 : WARNING : EPOCH - 73 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,796 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,797 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,800 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,829 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,830 : INFO : EPOCH - 74 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,831 : WARNING : EPOCH - 74 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,842 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,843 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,844 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,874 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,876 : INFO : EPOCH - 75 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,876 : WARNING : EPOCH - 75 : supplied example count (1) did not equal expec
2019-02-21 17:14:58,887 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,888 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,888 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,918 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,919 : INFO : EPOCH - 76 : training on 5386246 raw words (10000 effective wc
2019-02-21 17:14:58,919 : WARNING : EPOCH - 76 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:58,931 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,933 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,933 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:58,963 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:58,965 : INFO : EPOCH - 77 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:58,966 : WARNING : EPOCH - 77 : supplied example count (1) did not equal expe
2019-02-21 17:14:58,979 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:58,980 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:58,981 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,012 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,013 : INFO : EPOCH - 78 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,014 : WARNING : EPOCH - 78 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,026 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,027 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,027 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,058 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,059 : INFO : EPOCH - 79 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,060 : WARNING : EPOCH - 79 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,072 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,073 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,074 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,105 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,106 : INFO : EPOCH - 80 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,107 : WARNING : EPOCH - 80 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,118 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,119 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,120 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,149 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,150 : INFO : EPOCH - 81 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,151 : WARNING : EPOCH - 81 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,160 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,162 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,163 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,194 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,195 : INFO : EPOCH - 82 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,196 : WARNING : EPOCH - 82 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,211 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,212 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,213 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,243 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,244 : INFO : EPOCH - 83 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,245 : WARNING : EPOCH - 83 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,257 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,258 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,259 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,289 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,290 : INFO : EPOCH - 84 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,291 : WARNING : EPOCH - 84 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:14:59,297 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,303 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,304 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,335 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,336 : INFO : EPOCH - 85 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,337 : WARNING : EPOCH - 85 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,349 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,350 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,351 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,382 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,383 : INFO : EPOCH - 86 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,384 : WARNING : EPOCH - 86 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,394 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,395 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,397 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,426 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,427 : INFO : EPOCH - 87 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,428 : WARNING : EPOCH - 87 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,440 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,442 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,443 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,472 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,473 : INFO : EPOCH - 88 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,474 : WARNING : EPOCH - 88 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,484 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,488 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,491 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,517 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,518 : INFO : EPOCH - 89 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,519 : WARNING : EPOCH - 89 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,530 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,531 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,533 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,563 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,564 : INFO : EPOCH - 90 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,565 : WARNING : EPOCH - 90 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,575 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,576 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,577 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,604 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,605 : INFO : EPOCH - 91 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,606 : WARNING : EPOCH - 91 : supplied example count (1) did not equal expec
2019-02-21 17:14:59,616 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,617 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,617 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,647 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,648 : INFO : EPOCH - 92 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,649 : WARNING : EPOCH - 92 : supplied example count (1) did not equal expec
```

```
2019-02-21 17:14:59,658 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,659 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,661 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,690 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,691 : INFO : EPOCH - 93 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,692 : WARNING : EPOCH - 93 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,703 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,704 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,705 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,734 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,735 : INFO : EPOCH - 94 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,736 : WARNING : EPOCH - 94 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,746 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,747 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,747 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,777 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,778 : INFO : EPOCH - 95 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,779 : WARNING : EPOCH - 95 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,791 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,793 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,793 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,824 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,825 : INFO : EPOCH - 96 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,826 : WARNING : EPOCH - 96 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,836 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,837 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,838 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,869 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,870 : INFO : EPOCH - 97 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,870 : WARNING : EPOCH - 97 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,881 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,883 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,883 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,913 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,914 : INFO : EPOCH - 98 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,915 : WARNING : EPOCH - 98 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,926 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,927 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,928 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:14:59,959 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:14:59,961 : INFO : EPOCH - 99 : training on 5386246 raw words (10000 effective wo
2019-02-21 17:14:59,961 : WARNING : EPOCH - 99 : supplied example count (1) did not equal expe
2019-02-21 17:14:59,975 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:14:59,976 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:14:59,977 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,008 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,009 : INFO : EPOCH - 100 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,010 : WARNING : EPOCH - 100 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:00,024 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,025 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,026 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,055 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,056 : INFO : EPOCH - 101 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,057 : WARNING : EPOCH - 101 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,067 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,068 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,068 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,096 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,097 : INFO : EPOCH - 102 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,098 : WARNING : EPOCH - 102 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,110 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,111 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,112 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,142 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,143 : INFO : EPOCH - 103 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,145 : WARNING : EPOCH - 103 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,156 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,157 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,158 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,191 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,192 : INFO : EPOCH - 104 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,193 : WARNING : EPOCH - 104 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,207 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,208 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,210 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,240 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,241 : INFO : EPOCH - 105 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,241 : WARNING : EPOCH - 105 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,252 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,253 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,254 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,285 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,286 : INFO : EPOCH - 106 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,287 : WARNING : EPOCH - 106 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,297 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,299 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,300 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,331 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,332 : INFO : EPOCH - 107 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,332 : WARNING : EPOCH - 107 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,342 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,344 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,345 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,375 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,376 : INFO : EPOCH - 108 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,377 : WARNING : EPOCH - 108 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:00,387 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,389 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,389 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,419 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,420 : INFO : EPOCH - 109 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,420 : WARNING : EPOCH - 109 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,430 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,433 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,434 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,463 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,465 : INFO : EPOCH - 110 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,465 : WARNING : EPOCH - 110 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,477 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,482 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,483 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,511 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,512 : INFO : EPOCH - 111 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,515 : WARNING : EPOCH - 111 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,526 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,528 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,529 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,561 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,562 : INFO : EPOCH - 112 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,563 : WARNING : EPOCH - 112 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,575 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,576 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,577 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,608 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,609 : INFO : EPOCH - 113 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,609 : WARNING : EPOCH - 113 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,621 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,623 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,624 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,655 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,656 : INFO : EPOCH - 114 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,657 : WARNING : EPOCH - 114 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,671 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,672 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,673 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,701 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,702 : INFO : EPOCH - 115 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,702 : WARNING : EPOCH - 115 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,712 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,713 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,713 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,741 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,742 : INFO : EPOCH - 116 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:00,743 : WARNING : EPOCH - 116 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:00,752 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,753 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,753 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,782 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,783 : INFO : EPOCH - 117 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,784 : WARNING : EPOCH - 117 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,794 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,794 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,795 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,839 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,844 : INFO : EPOCH - 118 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,846 : WARNING : EPOCH - 118 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,859 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,860 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,861 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,887 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,888 : INFO : EPOCH - 119 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,889 : WARNING : EPOCH - 119 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,898 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,899 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,900 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,927 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,928 : INFO : EPOCH - 120 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,929 : WARNING : EPOCH - 120 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,939 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,941 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,941 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:00,968 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:00,969 : INFO : EPOCH - 121 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:00,970 : WARNING : EPOCH - 121 : supplied example count (1) did not equal expe
2019-02-21 17:15:00,980 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:00,981 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:00,982 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,009 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,010 : INFO : EPOCH - 122 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,011 : WARNING : EPOCH - 122 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,021 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,023 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,024 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,054 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,056 : INFO : EPOCH - 123 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,057 : WARNING : EPOCH - 123 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,067 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,067 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,068 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,096 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,097 : INFO : EPOCH - 124 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,098 : WARNING : EPOCH - 124 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:01,108 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,110 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,110 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,138 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,139 : INFO : EPOCH - 125 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,140 : WARNING : EPOCH - 125 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,149 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,150 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,151 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,178 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,179 : INFO : EPOCH - 126 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,180 : WARNING : EPOCH - 126 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,194 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,195 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,196 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,223 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,224 : INFO : EPOCH - 127 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,225 : WARNING : EPOCH - 127 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,236 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,237 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,238 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,264 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,265 : INFO : EPOCH - 128 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,265 : WARNING : EPOCH - 128 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,275 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,276 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,277 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,306 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,306 : INFO : EPOCH - 129 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,307 : WARNING : EPOCH - 129 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,317 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,327 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,328 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,344 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,345 : INFO : EPOCH - 130 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,346 : WARNING : EPOCH - 130 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,356 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,358 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,358 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,387 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,388 : INFO : EPOCH - 131 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,388 : WARNING : EPOCH - 131 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,399 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,401 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,402 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,433 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,434 : INFO : EPOCH - 132 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,434 : WARNING : EPOCH - 132 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:01,446 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,447 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,448 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,478 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,479 : INFO : EPOCH - 133 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,480 : WARNING : EPOCH - 133 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,492 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,493 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,494 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,524 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,525 : INFO : EPOCH - 134 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,526 : WARNING : EPOCH - 134 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,540 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,542 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,543 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,573 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,574 : INFO : EPOCH - 135 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,575 : WARNING : EPOCH - 135 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,587 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,588 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,589 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,620 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,621 : INFO : EPOCH - 136 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,622 : WARNING : EPOCH - 136 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,634 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,639 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,641 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,665 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,666 : INFO : EPOCH - 137 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,667 : WARNING : EPOCH - 137 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,679 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,679 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,680 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,707 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,708 : INFO : EPOCH - 138 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,709 : WARNING : EPOCH - 138 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,718 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,719 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,720 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,747 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,748 : INFO : EPOCH - 139 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,749 : WARNING : EPOCH - 139 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,758 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,760 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,760 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,789 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,790 : INFO : EPOCH - 140 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:01,791 : WARNING : EPOCH - 140 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:01,800 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,803 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,803 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,833 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,834 : INFO : EPOCH - 141 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,835 : WARNING : EPOCH - 141 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,850 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,852 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,852 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,883 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,888 : INFO : EPOCH - 142 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,888 : WARNING : EPOCH - 142 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,900 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,902 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,903 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,933 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,934 : INFO : EPOCH - 143 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,935 : WARNING : EPOCH - 143 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,947 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,948 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,949 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:01,980 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:01,981 : INFO : EPOCH - 144 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:01,982 : WARNING : EPOCH - 144 : supplied example count (1) did not equal expe
2019-02-21 17:15:01,994 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:01,995 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:01,996 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,025 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,026 : INFO : EPOCH - 145 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:02,027 : WARNING : EPOCH - 145 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,039 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,041 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,042 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,072 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,073 : INFO : EPOCH - 146 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:02,073 : WARNING : EPOCH - 146 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,084 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,086 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,086 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,117 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,118 : INFO : EPOCH - 147 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:02,118 : WARNING : EPOCH - 147 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,130 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,131 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,131 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,162 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,163 : INFO : EPOCH - 148 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:02,164 : WARNING : EPOCH - 148 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:02,177 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,179 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,179 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,208 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,209 : INFO : EPOCH - 149 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,210 : WARNING : EPOCH - 149 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,223 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,225 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,226 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,257 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,258 : INFO : EPOCH - 150 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,258 : WARNING : EPOCH - 150 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,269 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,271 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,272 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,302 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,303 : INFO : EPOCH - 151 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,304 : WARNING : EPOCH - 151 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,315 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,317 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,317 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,348 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,349 : INFO : EPOCH - 152 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,349 : WARNING : EPOCH - 152 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,362 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,363 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,365 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,394 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,395 : INFO : EPOCH - 153 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,396 : WARNING : EPOCH - 153 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,406 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,412 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,417 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,439 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,440 : INFO : EPOCH - 154 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,440 : WARNING : EPOCH - 154 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,451 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,452 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,453 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,483 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,484 : INFO : EPOCH - 155 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,485 : WARNING : EPOCH - 155 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,497 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,498 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,499 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,543 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,550 : INFO : EPOCH - 156 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:02,550 : WARNING : EPOCH - 156 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:02,559 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,561 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,562 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,588 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,589 : INFO : EPOCH - 157 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,590 : WARNING : EPOCH - 157 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,600 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,601 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,602 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,628 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,629 : INFO : EPOCH - 158 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,629 : WARNING : EPOCH - 158 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,639 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,640 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,641 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,671 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,672 : INFO : EPOCH - 159 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,672 : WARNING : EPOCH - 159 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,690 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,699 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,700 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,718 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,719 : INFO : EPOCH - 160 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,720 : WARNING : EPOCH - 160 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,729 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,730 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,731 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,759 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,760 : INFO : EPOCH - 161 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,761 : WARNING : EPOCH - 161 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,771 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,772 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,773 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,798 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,799 : INFO : EPOCH - 162 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,800 : WARNING : EPOCH - 162 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,809 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,810 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,811 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,837 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,837 : INFO : EPOCH - 163 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,838 : WARNING : EPOCH - 163 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,846 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,847 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,848 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,875 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,876 : INFO : EPOCH - 164 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,876 : WARNING : EPOCH - 164 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:02,888 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,890 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,890 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,920 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,923 : INFO : EPOCH - 165 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,924 : WARNING : EPOCH - 165 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,934 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,935 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,935 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:02,963 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:02,970 : INFO : EPOCH - 166 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:02,971 : WARNING : EPOCH - 166 : supplied example count (1) did not equal expe
2019-02-21 17:15:02,983 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:02,984 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:02,985 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,015 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,016 : INFO : EPOCH - 167 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,017 : WARNING : EPOCH - 167 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,029 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,030 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,031 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,061 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,062 : INFO : EPOCH - 168 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,062 : WARNING : EPOCH - 168 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,073 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,074 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,075 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,105 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,106 : INFO : EPOCH - 169 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,107 : WARNING : EPOCH - 169 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,118 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,120 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,120 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,150 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,151 : INFO : EPOCH - 170 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,152 : WARNING : EPOCH - 170 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,164 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,165 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,166 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,196 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,197 : INFO : EPOCH - 171 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,198 : WARNING : EPOCH - 171 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,214 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,215 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,216 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,246 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,253 : INFO : EPOCH - 172 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,254 : WARNING : EPOCH - 172 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:03,265 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,267 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,268 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,297 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,298 : INFO : EPOCH - 173 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,299 : WARNING : EPOCH - 173 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,308 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,309 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,310 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,338 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,339 : INFO : EPOCH - 174 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,340 : WARNING : EPOCH - 174 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,351 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,352 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,353 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,383 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,384 : INFO : EPOCH - 175 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,385 : WARNING : EPOCH - 175 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,396 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,397 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,399 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,428 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,429 : INFO : EPOCH - 176 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,429 : WARNING : EPOCH - 176 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,441 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,442 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,443 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,474 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,477 : INFO : EPOCH - 177 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,477 : WARNING : EPOCH - 177 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,489 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,490 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,491 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,521 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,522 : INFO : EPOCH - 178 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,523 : WARNING : EPOCH - 178 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,535 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,537 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,538 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,568 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,569 : INFO : EPOCH - 179 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,569 : WARNING : EPOCH - 179 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,581 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,614 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,615 : INFO : EPOCH - 180 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:03,616 : WARNING : EPOCH - 180 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:03,624 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,626 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,626 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,653 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,654 : INFO : EPOCH - 181 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,655 : WARNING : EPOCH - 181 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,664 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,665 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,667 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,694 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,695 : INFO : EPOCH - 182 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,695 : WARNING : EPOCH - 182 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,709 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,709 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,739 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,742 : INFO : EPOCH - 183 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,743 : WARNING : EPOCH - 183 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,753 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,754 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,755 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,783 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,784 : INFO : EPOCH - 184 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,784 : WARNING : EPOCH - 184 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,801 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,802 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,803 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,832 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,834 : INFO : EPOCH - 185 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,834 : WARNING : EPOCH - 185 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,843 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,845 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,846 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,875 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,876 : INFO : EPOCH - 186 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,877 : WARNING : EPOCH - 186 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,888 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,889 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,890 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,920 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,921 : INFO : EPOCH - 187 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,922 : WARNING : EPOCH - 187 : supplied example count (1) did not equal expe
2019-02-21 17:15:03,933 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,934 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,935 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:03,964 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:03,965 : INFO : EPOCH - 188 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:03,966 : WARNING : EPOCH - 188 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:03,976 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:03,978 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:03,978 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,007 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,008 : INFO : EPOCH - 189 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,009 : WARNING : EPOCH - 189 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,020 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,021 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,022 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,052 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,053 : INFO : EPOCH - 190 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,054 : WARNING : EPOCH - 190 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,066 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,074 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,075 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,099 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,100 : INFO : EPOCH - 191 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,101 : WARNING : EPOCH - 191 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,113 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,114 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,116 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,144 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,145 : INFO : EPOCH - 192 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,146 : WARNING : EPOCH - 192 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,158 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,159 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,160 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,188 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,189 : INFO : EPOCH - 193 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,190 : WARNING : EPOCH - 193 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,204 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,205 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,206 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,235 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,236 : INFO : EPOCH - 194 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,237 : WARNING : EPOCH - 194 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,248 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,250 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,251 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,281 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,283 : INFO : EPOCH - 195 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,283 : WARNING : EPOCH - 195 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,296 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,297 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,300 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,328 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,338 : INFO : EPOCH - 196 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:04,339 : WARNING : EPOCH - 196 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:04,349 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,351 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,352 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,382 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,383 : INFO : EPOCH - 197 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,383 : WARNING : EPOCH - 197 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,395 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,397 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,398 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,428 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,429 : INFO : EPOCH - 198 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,429 : WARNING : EPOCH - 198 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,441 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,442 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,443 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,472 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,473 : INFO : EPOCH - 199 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,473 : WARNING : EPOCH - 199 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,482 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,484 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,484 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,511 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,512 : INFO : EPOCH - 200 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,512 : WARNING : EPOCH - 200 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,522 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,524 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,524 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,552 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,553 : INFO : EPOCH - 201 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,554 : WARNING : EPOCH - 201 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,566 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,567 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,568 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,597 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,598 : INFO : EPOCH - 202 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,599 : WARNING : EPOCH - 202 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,610 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,612 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,613 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,643 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,644 : INFO : EPOCH - 203 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,645 : WARNING : EPOCH - 203 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,656 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,657 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,658 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,687 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,688 : INFO : EPOCH - 204 : training on 5386246 raw words (10000 effective v
2019-02-21 17:15:04,689 : WARNING : EPOCH - 204 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:04,698 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,700 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,701 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,726 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,727 : INFO : EPOCH - 205 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,728 : WARNING : EPOCH - 205 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,738 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,739 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,740 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,769 : INFO : EPOCH - 206 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,769 : WARNING : EPOCH - 206 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,781 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,783 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,808 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,809 : INFO : EPOCH - 207 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,810 : WARNING : EPOCH - 207 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,820 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,821 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,823 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,851 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,852 : INFO : EPOCH - 208 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,853 : WARNING : EPOCH - 208 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,868 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,869 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,870 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,900 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,901 : INFO : EPOCH - 209 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,903 : WARNING : EPOCH - 209 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,914 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,916 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,917 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,946 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,948 : INFO : EPOCH - 210 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,948 : WARNING : EPOCH - 210 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,956 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,958 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,959 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:04,985 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:04,986 : INFO : EPOCH - 211 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:04,986 : WARNING : EPOCH - 211 : supplied example count (1) did not equal expe
2019-02-21 17:15:04,996 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:04,996 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:04,997 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,023 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,024 : INFO : EPOCH - 212 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,025 : WARNING : EPOCH - 212 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:05,035 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,035 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,036 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,063 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,064 : INFO : EPOCH - 213 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,065 : WARNING : EPOCH - 213 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,076 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,077 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,078 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,112 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,114 : INFO : EPOCH - 214 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,116 : WARNING : EPOCH - 214 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,130 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,135 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,137 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,163 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,164 : INFO : EPOCH - 215 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,165 : WARNING : EPOCH - 215 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,180 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,181 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,182 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,210 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,211 : INFO : EPOCH - 216 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,212 : WARNING : EPOCH - 216 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,220 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,222 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,224 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,247 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,248 : INFO : EPOCH - 217 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,248 : WARNING : EPOCH - 217 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,259 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,259 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,260 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,286 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,287 : INFO : EPOCH - 218 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,288 : WARNING : EPOCH - 218 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,298 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,300 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,301 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,326 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,327 : INFO : EPOCH - 219 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,328 : WARNING : EPOCH - 219 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,338 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,339 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,340 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,366 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,367 : INFO : EPOCH - 220 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,368 : WARNING : EPOCH - 220 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:05,379 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,380 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,381 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,408 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,408 : INFO : EPOCH - 221 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,409 : WARNING : EPOCH - 221 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,419 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,431 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,433 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,449 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,449 : INFO : EPOCH - 222 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,450 : WARNING : EPOCH - 222 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,459 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,460 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,461 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,488 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,489 : INFO : EPOCH - 223 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,490 : WARNING : EPOCH - 223 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,499 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,501 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,528 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,529 : INFO : EPOCH - 224 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,530 : WARNING : EPOCH - 224 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,543 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,575 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,576 : INFO : EPOCH - 225 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,577 : WARNING : EPOCH - 225 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,589 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,591 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,592 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,620 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,621 : INFO : EPOCH - 226 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,622 : WARNING : EPOCH - 226 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,636 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,637 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,638 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,668 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,669 : INFO : EPOCH - 227 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,669 : WARNING : EPOCH - 227 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,681 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,683 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,684 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,713 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,714 : INFO : EPOCH - 228 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:05,715 : WARNING : EPOCH - 228 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:05,726 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,727 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,727 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,754 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,755 : INFO : EPOCH - 229 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,756 : WARNING : EPOCH - 229 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,765 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,766 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,767 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,794 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,795 : INFO : EPOCH - 230 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,796 : WARNING : EPOCH - 230 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,808 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,809 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,810 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,836 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,836 : INFO : EPOCH - 231 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,837 : WARNING : EPOCH - 231 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,846 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,848 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,849 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,873 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,874 : INFO : EPOCH - 232 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,875 : WARNING : EPOCH - 232 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,886 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,886 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,887 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,915 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,915 : INFO : EPOCH - 233 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,916 : WARNING : EPOCH - 233 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,931 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,932 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,933 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:05,963 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:05,963 : INFO : EPOCH - 234 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:05,964 : WARNING : EPOCH - 234 : supplied example count (1) did not equal expe
2019-02-21 17:15:05,973 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:05,974 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:05,975 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,000 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,001 : INFO : EPOCH - 235 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,002 : WARNING : EPOCH - 235 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,012 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,013 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,014 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,041 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,042 : INFO : EPOCH - 236 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,043 : WARNING : EPOCH - 236 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:06,055 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,056 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,057 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,087 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,088 : INFO : EPOCH - 237 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,089 : WARNING : EPOCH - 237 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,099 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,100 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,102 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,132 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,133 : INFO : EPOCH - 238 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,134 : WARNING : EPOCH - 238 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,144 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,145 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,146 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,176 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,178 : INFO : EPOCH - 239 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,179 : WARNING : EPOCH - 239 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,192 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,195 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,197 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,224 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,225 : INFO : EPOCH - 240 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,225 : WARNING : EPOCH - 240 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,239 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,241 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,242 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,268 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,269 : INFO : EPOCH - 241 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,271 : WARNING : EPOCH - 241 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,283 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,284 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,285 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,314 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,314 : INFO : EPOCH - 242 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,315 : WARNING : EPOCH - 242 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,324 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,337 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,337 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,353 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,353 : INFO : EPOCH - 243 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,354 : WARNING : EPOCH - 243 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,364 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,364 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,365 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,394 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,394 : INFO : EPOCH - 244 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:06,395 : WARNING : EPOCH - 244 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:06,407 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,408 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,409 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,436 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,437 : INFO : EPOCH - 245 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,438 : WARNING : EPOCH - 245 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,447 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,448 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,449 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,476 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,477 : INFO : EPOCH - 246 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,478 : WARNING : EPOCH - 246 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,494 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,495 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,496 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,526 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,527 : INFO : EPOCH - 247 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,528 : WARNING : EPOCH - 247 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,541 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,542 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,542 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,571 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,572 : INFO : EPOCH - 248 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,573 : WARNING : EPOCH - 248 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,583 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,585 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,587 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,611 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,612 : INFO : EPOCH - 249 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,613 : WARNING : EPOCH - 249 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,625 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,626 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,627 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,653 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,654 : INFO : EPOCH - 250 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,655 : WARNING : EPOCH - 250 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,665 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,666 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,667 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,693 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,694 : INFO : EPOCH - 251 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,695 : WARNING : EPOCH - 251 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,707 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,707 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,708 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,736 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,737 : INFO : EPOCH - 252 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,737 : WARNING : EPOCH - 252 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:06,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,748 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,749 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,775 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,776 : INFO : EPOCH - 253 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,776 : WARNING : EPOCH - 253 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,790 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,791 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,792 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,820 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,824 : INFO : EPOCH - 254 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,825 : WARNING : EPOCH - 254 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,835 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,836 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,837 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,867 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,868 : INFO : EPOCH - 255 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,869 : WARNING : EPOCH - 255 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,879 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,880 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,881 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,911 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,912 : INFO : EPOCH - 256 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,913 : WARNING : EPOCH - 256 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,924 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,925 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,926 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:06,955 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:06,956 : INFO : EPOCH - 257 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:06,957 : WARNING : EPOCH - 257 : supplied example count (1) did not equal expe
2019-02-21 17:15:06,968 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:06,969 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:06,970 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,000 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,002 : INFO : EPOCH - 258 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,002 : WARNING : EPOCH - 258 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,014 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,015 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,016 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,045 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,046 : INFO : EPOCH - 259 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,046 : WARNING : EPOCH - 259 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,059 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,061 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,061 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,091 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,097 : INFO : EPOCH - 260 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,098 : WARNING : EPOCH - 260 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:07,109 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,110 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,111 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,142 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,143 : INFO : EPOCH - 261 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,144 : WARNING : EPOCH - 261 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,155 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,156 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,157 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,187 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,189 : INFO : EPOCH - 262 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,190 : WARNING : EPOCH - 262 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,205 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,206 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,207 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,236 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,237 : INFO : EPOCH - 263 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,238 : WARNING : EPOCH - 263 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,248 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,250 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,250 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,280 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,281 : INFO : EPOCH - 264 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,282 : WARNING : EPOCH - 264 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,293 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,295 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,296 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,324 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,325 : INFO : EPOCH - 265 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,326 : WARNING : EPOCH - 265 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,335 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,337 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,337 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,367 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,368 : INFO : EPOCH - 266 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,369 : WARNING : EPOCH - 266 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,382 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,383 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,384 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,414 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,416 : INFO : EPOCH - 267 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,417 : WARNING : EPOCH - 267 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,428 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,430 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,430 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,460 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,461 : INFO : EPOCH - 268 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,462 : WARNING : EPOCH - 268 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:07,475 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,478 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,479 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,507 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,508 : INFO : EPOCH - 269 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,509 : WARNING : EPOCH - 269 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,520 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,521 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,524 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,549 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,550 : INFO : EPOCH - 270 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,551 : WARNING : EPOCH - 270 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,560 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,561 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,562 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,588 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,589 : INFO : EPOCH - 271 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,590 : WARNING : EPOCH - 271 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,600 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,601 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,602 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,631 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,632 : INFO : EPOCH - 272 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,633 : WARNING : EPOCH - 272 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,644 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,646 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,648 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,678 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,679 : INFO : EPOCH - 273 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,680 : WARNING : EPOCH - 273 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,691 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,696 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,699 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,723 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,723 : INFO : EPOCH - 274 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,724 : WARNING : EPOCH - 274 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,736 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,738 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,739 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,767 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,768 : INFO : EPOCH - 275 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,768 : WARNING : EPOCH - 275 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,779 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,780 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,781 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,811 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,812 : INFO : EPOCH - 276 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,813 : WARNING : EPOCH - 276 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:07,824 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,831 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,832 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,857 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,858 : INFO : EPOCH - 277 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,859 : WARNING : EPOCH - 277 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,874 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,875 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,876 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,909 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,910 : INFO : EPOCH - 278 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,911 : WARNING : EPOCH - 278 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,928 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,930 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,930 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:07,961 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:07,961 : INFO : EPOCH - 279 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:07,962 : WARNING : EPOCH - 279 : supplied example count (1) did not equal expe
2019-02-21 17:15:07,975 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:07,976 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:07,977 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,004 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,005 : INFO : EPOCH - 280 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:08,006 : WARNING : EPOCH - 280 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,015 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,017 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,018 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,044 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,045 : INFO : EPOCH - 281 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:08,045 : WARNING : EPOCH - 281 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,058 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,060 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,061 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,088 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,089 : INFO : EPOCH - 282 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:08,090 : WARNING : EPOCH - 282 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,100 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,102 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,103 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,128 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,129 : INFO : EPOCH - 283 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:08,129 : WARNING : EPOCH - 283 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,141 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,142 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,143 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,170 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,171 : INFO : EPOCH - 284 : training on 5386246 raw words (10000 effective w
2019-02-21 17:15:08,172 : WARNING : EPOCH - 284 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:08,191 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,193 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,195 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,222 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,224 : INFO : EPOCH - 285 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,224 : WARNING : EPOCH - 285 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,236 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,237 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,239 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,265 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,267 : INFO : EPOCH - 286 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,267 : WARNING : EPOCH - 286 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,278 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,279 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,280 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,311 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,312 : INFO : EPOCH - 287 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,312 : WARNING : EPOCH - 287 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,323 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,325 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,326 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,351 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,352 : INFO : EPOCH - 288 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,353 : WARNING : EPOCH - 288 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,363 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,364 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,365 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,391 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,393 : INFO : EPOCH - 289 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,393 : WARNING : EPOCH - 289 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,404 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,405 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,406 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,435 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,436 : INFO : EPOCH - 290 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,436 : WARNING : EPOCH - 290 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,452 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,453 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,454 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,481 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,487 : INFO : EPOCH - 291 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,488 : WARNING : EPOCH - 291 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,498 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,499 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,503 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,529 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,530 : INFO : EPOCH - 292 : training on 5386246 raw words (10000 effective
2019-02-21 17:15:08,531 : WARNING : EPOCH - 292 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:08,543 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,544 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,545 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,572 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,572 : INFO : EPOCH - 293 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,573 : WARNING : EPOCH - 293 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,581 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,582 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,583 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,608 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,609 : INFO : EPOCH - 294 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,610 : WARNING : EPOCH - 294 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,620 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,621 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,621 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,648 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,648 : INFO : EPOCH - 295 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,649 : WARNING : EPOCH - 295 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,659 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,660 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,661 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,687 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,689 : INFO : EPOCH - 296 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,689 : WARNING : EPOCH - 296 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,701 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,702 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,703 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,734 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,735 : INFO : EPOCH - 297 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,736 : WARNING : EPOCH - 297 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,747 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,748 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,760 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,779 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,780 : INFO : EPOCH - 298 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,781 : WARNING : EPOCH - 298 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,793 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,794 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,795 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,824 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,825 : INFO : EPOCH - 299 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,826 : WARNING : EPOCH - 299 : supplied example count (1) did not equal expe
2019-02-21 17:15:08,835 : INFO : worker thread finished; awaiting finish of 3 more threads
2019-02-21 17:15:08,839 : INFO : worker thread finished; awaiting finish of 2 more threads
2019-02-21 17:15:08,840 : INFO : worker thread finished; awaiting finish of 1 more threads
2019-02-21 17:15:08,867 : INFO : worker thread finished; awaiting finish of 0 more threads
2019-02-21 17:15:08,868 : INFO : EPOCH - 300 : training on 5386246 raw words (10000 effective u
2019-02-21 17:15:08,869 : WARNING : EPOCH - 300 : supplied example count (1) did not equal expe
```

```
2019-02-21 17:15:08,870 : INFO : training on a 1615873800 raw words (3000000 effective words) t
```

```
In [9]: model.save(str("model/" + now.replace(".","-").replace(":","-") + ".model"))
```

```
2019-02-21 17:15:08,881 : INFO : saving Word2Vec object under model/2019-02-2117-14-38-361318.
2019-02-21 17:15:08,882 : INFO : not storing attribute vectors_norm
2019-02-21 17:15:08,884 : INFO : not storing attribute cum_table
2019-02-21 17:15:10,897 : INFO : saved model/2019-02-2117-14-38-361318.model
```

## 4  Test

```
In [10]: result = model.wv.most_similar(positive="microsoft",topn=10)
```

```
2019-02-21 17:15:10,905 : INFO : precomputing L2-norms of word weight vectors
```

```
In [11]: print(result)
```

```
[('dos', 0.8956843614578247), ('systÃÍmes', 0.8914103507995605), ('assimilÃľe', 0.8331711292266
```

## 5  Save

```
In [12]: import save_notebook
```

```
In [13]: name_notebook_exported = save_notebook.save_notebook("word2vec_with_gensim.ipynb")
```

```
In [14]: def write_result(word_embedding, time_training,name_notebook_exported, fname ):
            if not os.path.isfile(fname):
                f=open(fname, "a+")
                f.write("Nombre de mots;Model de word embedding;temps d'apprentissage;Noteboo
                f.close
            f=open(fname, "a+")
            f.write("\n" + str( str(len(wiki_vocab_tokenized)) + ";" +word_embedding + ";"+ s
            f.close
```

```
In [15]: write_result(word_embedding, time_training,name_notebook_exported, "resultats.csv" )
```

```
In [16]: name_notebook_exported
```

```
Out[16]: 'resultats/word2vec_with_gensim2019-02-2117-15-11-106036.pdf'
```