# Predicting Stackoverflow tags
# 02807 Final project

November 29, 2016

## 1 Introduction

In this project a model predicting tags associating a given text will be constructed. The problem will be handled as a classification problem, hence a classifier will be trained.

The classifier will be trained on posts from Stackoverflow where associated tags have been given.

Initially only the top 20 frequent tags and posts containing these tags will be used for training and prediction.

Finally the model will be evaluated using all tags and posts.

During implementation and debugging a subset of the data will be used.

## 2 The data

The dataset consists of two XML files, one containing all possible tags and their corresponding counts, and one containing posts with *title*, *body*, *tags* and some meta data.

The total size of the files are approximately 49GB in uncompressed format.

## 3 Methods

In this section the different methods and steps in the process will be explained.

### 3.1 Preprocessing

The preprocessing step regards the transforming of questions in an `XML` file to processed questions in a `.csv` file. This also includes disregarding questions that does not have any of the top $N$ tags attached. At the same time two other `.csv`

files are created: One containing all unique words in the extracted questions, and one containing the unique tags used.

The processing of each questions contains the following steps (code can be found in appendix A.1)

1. Replace all links with *<link>*

2. Remove certain unwanted symbols

3. Remove suffix from words (e.g. *haven't → have*)

4. Remove line breaks

5. Replace digits with *<digit>*

6. Remove double whitespaces

7. Reduce words to their word stem (e.g. *lines → line*)

8. Lemmatize words (e.g. *better → good*)

Finally the unique words used as word dictionary were filtered by removing words occuring in more than 50% of the questions and words occuring in less than 0.1% of the questions. Also english stop words were removed making use of the `NLTK` library.

All these steps are used in order to reduce the dimensionality of the word space without really removing much information. Here the assumption is, that e.g. words like *better* and *good* kind of adds the same meaning to the sentence, and the same with e.g. two numbers.

## 3.2   Distributed file loading

Since the size of the final processed questions file is approximately 11GB, it will not be feasible to load into memory on most laptops. Therefore it will be necessary to load the file in smaller chunks.

The following code illustrates how the file `posts.csv` can be divided into byte-chunks. I.e. the following generator yields a list of tuples (`from_byte, size`) where `from_byte` is the index in the file in bytes and `size` is the size of the given chunk in bytes.

```
with open('posts.csv', 'rb') as f:
  while True:
    start = f.tell()
    f.seek(chunk_size, 1)
    s = f.readline()
    yield start, f.tell() - start
    if not s:    break
```

The `chunk_size` is a given minimum size of each chunk. The `f.readline()` makes sure the chunk ends at the end of a line.

A chunk of lines from the file can then be loaded using the following lines

```
1  # Seek to chunk start bytes
2  f.seek(from_bytes)
3
4  # Read end of chunk until end of line
5  chunk = f.read(size)
6
7  # Split in lines (Removing the last newline)
8  lines = chunk.rstrip('\n').split('\n')
```

## 3.3   K-means clustering

### 3.3.1   Serial

The regular serial in-memory version of K-means clustering algorithm is shown in algorithm 1.

---
**Algorithm 1** Serial K-means clustering algorithm

---
1:  **procedure** KMEANSCLUSTERING(X, K)
2:      # Initialize cluster centers
3:      **for** $k = 0$ to $K - 1$ **do**
4:          $\mu_k \leftarrow$ random point in X
5:      # Run iterations
6:      **while** $iter < max\_iter$ **do**
7:          # Update cluster means
8:          $\mu_{\text{old}} = \mu$
9:          **for** $k = 0$ to $K - 1$ **do**
10:              $C_k \leftarrow \{$Points in $X$ closest to $\mu_k\}$
11:              $\mu_k \leftarrow \frac{1}{|C_k|} \sum_{x_i \in C_k} x_i$
12:          # Check convergence criteria
13:          $norm \leftarrow \|\mu - \mu_{\text{old}}\|$
14:          **if** $norm < \epsilon$ **then**
15:              break

---

### 3.3.2   Distributed

The proposed distributed K-means clustering algorithm, which loads the data matrix $X$ in chunks, is shown in algorithm 2.

---

**Algorithm 2** Distributed K-means clustering algorithm

---

1: **procedure** KMEANSCLUSTERINGDISTRIBUTED(X, K)
2:     # Initialize cluster centers
3:     **for** $k = 0$ to $K - 1$ **do**
4:      $\mu_k \leftarrow$ random point in X
5:     # Run iterations
6:     **while** $iter < max\_iter$ **do**
7:      # Initialize shared cluster sums and cluster point counts.
8:      **for** $k = 0$ to $K - 1$ **do**
9:       $C\text{sum}_k \leftarrow \mathbf{0}$
10:      $C\text{count}_k \leftarrow 0$
11:     # Process each chunk in a distributed manner
12:     **for all** Chunks $X_{\text{chunk}}$ in X **do**
13:      **for** $k = 0$ to $K - 1$ **do**
14:       $C_k \leftarrow \{\text{Points in } X_{\text{chunk}} \text{ closest to } \mu_k\}$
15:       $C\text{sum}_k \leftarrow C\text{sum}_k + \sum_{x_i \in C_k} x_i$
16:       $C\text{count}_k \leftarrow C\text{count}_k + |C_k|$
17:     # Gather results and update cluster means
18:     $\mu_{\text{old}} = \mu$
19:     **for** $k = 0$ to $K - 1$ **do**
20:      $\mu_k \leftarrow \frac{C\text{sum}_k}{C\text{count}_k}$
21:     # Check convergence criteria
22:     $norm \leftarrow \|\mu - \mu_{\text{old}}\|$
23:     **if** $norm < \epsilon$ **then**
24:      break

---

### 3.3.3 Implementation

Simplified implementation of distributed K-means (only a single iteration is shown) see full code in appendix A.2:

```python
cluster_sums   = {k: np.zeros((1, word_count)) for k in
    range(0, K)}
cluster_counts = {k: 0 for k in range(0, K)}


for chunk in chunks:

  # Load chunk lines to sparse matrix
  X = chunk_to_sparse_mat(chunk)

  # Get closest cluster indices
  max_idx = sparse_matrix_to_cluster_indices(X, mu)

  # Assign points to clusters
  mu_subs = collections.defaultdict(list)
  for i, k in enumerate(max_idx):
    mu_subs[k].append(X[i].toarray())

  # Compute sub-means
  for k in range(0, K):
```

4

```
19      mu_sub = mu_subs[k]
20      if len(mu_sub) == 0:     continue
21      cluster_sums[k] += np.asarray(mu_sub).mean(axis=0)
22      cluster_counts[k] += 1
23
24  # Save old means
25  mu_old = np.array(mu, copy=True)
26
27  # Update means
28  for k in range(0, K):
29     count = cluster_counts[k]
30     if count == 0:   continue
31     mu[k] = cluster_sums[k] / cluster_counts[k]
32
33  # Check convergence criteria
34  mu_norm = np.linalg.norm(mu - mu_old)
35
36  if mu_norm < epsilon:
37     print('Converged after %d iterations' % (iteration+1))
38     break
```

## 3.4   Parallel processing

# 4   Results

# 5   Discussion

# 6   Conclusion

# A   Code snippets

## A.1   Preprocess text

```python
import re
import Stemmer
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()

# Precompile regular expressions
reg_links  = re.compile(r'http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@
    .&+]|[!*\(\),]|(?:%[0-9a-fA-F][0-9a-fA-F]))+')
re_digits  = re.compile(r'\b\d+\b')
re_spaces  = re.compile(r'\s{2,}')

reg_symbols = re.compile(r'[^A-Za-z0-9(),!?\'\'\']')
reg_symb_1 = re.compile(r',')
reg_symb_2 = re.compile(r'!')
reg_symb_3 = re.compile(r'\(')
reg_symb_4 = re.compile(r'\)')
reg_symb_5 = re.compile(r'\?')
reg_symb_6 = re.compile(r'\'')

reg_suf_1 = re.compile(r'\'s')
reg_suf_2 = re.compile(r'\'ve')
reg_suf_3 = re.compile(r'n\'t')
reg_suf_4 = re.compile(r'\'re')
reg_suf_5 = re.compile(r'\'d')
reg_suf_6 = re.compile(r'\'ll')

stemmer = Stemmer.Stemmer('english')
word_to_stem = {}
def stem_word(word):
    if not word in word_to_stem:
        word_to_stem[word] = stemmer.stemWord(word)
    return word_to_stem[word]

word_to_lemma = {}
def lemmatize_word(word):
    if not word in word_to_lemma:
        word_to_lemma[word] = lemmatizer.lemmatize(word)
    return word_to_lemma[word]

def clean_string(text):
  # Replace links with link identifier
  text = reg_links.sub('<link>', text)

  # Remove certain symbols
  text = reg_symbols.sub(' ', text)

  # Remove suffix from words
  text = reg_suf_1.sub(' ', text)
  text = reg_suf_2.sub(' ', text)
```

```
49    text = reg_suf_3.sub(' ', text)
50    text = reg_suf_4.sub(' ', text)
51    text = reg_suf_5.sub(' ', text)
52    text = reg_suf_6.sub(' ', text)
53
54    # Remove "'" from string
55    text = reg_symb_6.sub('', text)
56
57    # Replace breaks with spaces
58    text = text.replace('<br />', ' ')
59    text = text.replace('\r\n', ' ')
60    text = text.replace('\r', ' ')
61    text = text.replace('\n', ' ')
62
63    # Pad symbols with spaces on both sides
64    text = reg_symb_1.sub(' , ', text)
65    text = reg_symb_2.sub(' ! ', text)
66    text = reg_symb_3.sub(' ( ', text)
67    text = reg_symb_4.sub(' ) ', text)
68    text = reg_symb_5.sub(' ? ', text)
69
70    # Replace digits with 'DIGIT'
71    text = re_digits.sub('<DIGIT>', text)
72
73    # Remove double whitespaces
74    text = re_spaces.sub(' ', text)
75    text = text.strip()
76
77    # Convert to lowercase
78    text = text.lower()
79
80    # Stem each word
81    text = ' '.join(stem_word(word) for word in text.split(' '
        ))
82
83    # Lemmatize each word
84    text = ' '.join(lemmatize_word(word) for word in text.
        split(' '))
```

## A.2   Distributed K-means

```
1  import math
2  import collections
3  import numpy as np
4  import multiprocessing
5  import time
6
7  import helpers
8  import config
9
10 # Read tags
11 tags, tag2idx, tag_count = helpers.read_tags()
12
13 # Read words
```

```
14  words, word2idx, word_count = helpers.read_words()
15
16  # Clusters
17  K = tag_count
18
19  # Initialize cluster centers
20  mu = np.random.rand(K, word_count)
21
22  # Get chunks
23  chunk_reader = helpers.ChunkReader(post_filename=config.
        paths.TRAIN_DATA_IDX, chunk_size=config.data.CHUNK_SIZE)
        # TODO: Change
24  chunks = [chunk for chunk in chunk_reader]
25  chunk_count = len(chunks)
26
27  # Split chunks across processes
28  n = math.ceil(chunk_count / config.algorithm.PROCESS_COUNT)
29  chunks_split = []
30  for i in range(0, len(chunks), n):
31      chunks_split.append(chunks[i:i+n])
32
33  # Initialize shared variable manager
34  manager = multiprocessing.Manager()
35  lock = multiprocessing.Lock()
36
37  # Define function to run in parallel
38  def process_chunks(chunks, word_count, K, mu, cluster_sums,
        cluster_counts, lock):
39      for chunk in chunks:
40
41          # Convert to sparse matrix
42          X, _ = helpers.chunk_to_sparse_mat(chunk, word_count)
43
44          if X is None:    continue
45
46          # Get closest cluster indices
47          max_idx = helpers.sparse_matrix_to_cluster_indices(X, mu
        )
48
49          mu_subs = collections.defaultdict(list)
50          for i, k in enumerate(max_idx):
51              mu_subs[k].append(X[i].toarray())
52
53          # Compute sub-means
54          for k in range(0, K):
55              mu_sub = mu_subs[k]
56              if len(mu_sub) == 0:     continue
57
58              with lock:
59                  cluster_sums[k] = cluster_sums[k] + np.asarray(
        mu_sub, dtype=np.float32).mean(axis=0)
60                  cluster_counts[k] += 1
61
62
```

```python
for iteration in range(0, config.algorithm.MAX_ITER):
    start = time.time()

    cluster_sums = manager.dict({k: np.zeros((1, word_count))
      for k in range(0, K)})
    cluster_counts = manager.dict({k: 0 for k in range(0, K)})

    # Init processes
    processes = []
    for i, chunk_list in enumerate(chunks_split):
     p = multiprocessing.Process(target=process_chunks,
     kwargs={
        'chunks': chunk_list,
        'word_count': word_count,
        'K': K,
        'mu': mu,
        'cluster_sums': cluster_sums,
        'cluster_counts': cluster_counts,
        'lock': lock
     })
     processes.append(p)

    # Start processes
    for p in processes:
      p.start()

    #print('Started %d processes' % (len(processes)))

    # Wait for processes to finish
    for p in processes:
      p.join()

    # Save old means
    mu_old = np.array(mu, copy=True)

    # Update means
    for k in range(0, K):
      count = cluster_counts[k]
      if count == 0:  continue
      mu[k] = cluster_sums[k] / cluster_counts[k]

    # Check convergence criteria
    mu_norm = np.linalg.norm(mu - mu_old)

    print('Iteration %d took: %.4fs' % (iteration + 1, time.
      time() - start))

    if mu_norm < config.algorithm.EPSILON:
      print('Converged after %d iterations' % (iteration+1))
      break


# Determine cluster tags
cluster_tag_counts = {k: {tag: 0 for tag in range(0, K)} for
```

9

```
           k in range(0, K)}
114  for chunk in chunks:
115
116     # Convert to sparse matrix
117     X, tags = helpers.chunk_to_sparse_mat(chunk, word_count)
118
119     if X is None:    continue
120
121     # Get closest cluster indices
122     max_idx = helpers.sparse_matrix_to_cluster_indices(X, mu)
123
124     # Count cluster tags
125     for i, k in enumerate(max_idx):
126        for tag_idx in tags[i]:
127           cluster_tag_counts[k][tag_idx] += 1
128
129  # Assign tags to clusters
130  tags_labelled = []
131  cluster2tag = {}
132  for k, tag_counts in cluster_tag_counts.items():
133     tag_counts_sorted = sorted(tag_counts.items(), key=lambda
        x: x[1], reverse=True)
134     for tag, count in tag_counts_sorted:
135        if tag not in tags_labelled:
136           cluster2tag[k] = tag
137           tags_labelled.append(tag)
138           break
139
140  # Save cluster tags dict
141  config.data.save_cluster_tags(cluster_tags=cluster2tag)
142
143  # Save means
144  with open(config.paths.MU, 'wb') as f:
145     np.save(f, mu)
```

Include helper
functions and
config