

EEE-586 Term Project Report: Headline Generation from Distinctive Summaries

Burak Künkçü
21501843

Furkan Şahinuç
21400503

Selim Furkan Tekin
21501391

Abstract—In this paper, we demonstrate a LSTM based encoder-decoder network for the headline generation task. State-of-the-art BERT language model is used extract context vectors. Attention mechanism is used together with encoder-decoder network to capture relevancy information in input sequence and improve overall system performance. Several summaries are generated and according to their TF-IDF score, the best summary is used in headline generation instead of evaluating network on the whole content. Generated headlines are compared according to the BLEU metric.

I. INTRODUCTION

Seq2Seq encoder-decoder networks achieved state of the art performances for tasks involving transformation of one sequence into another such as translation tasks from a language to another etc. Such models often use recurrent network cells such as RNNs and LSTMs to construct encoder-decoder architectures which generates hidden vector representations and decomposes output sequences word-by-word in sequential time steps. In addition, it is possible to optimize such networks by adding attention mechanism which became highly popular in neural network training recently. Attention mechanism allows the network to look the input sequence selectively and learn alignments between input and output sequences.

In this term project, we have studied headline generation from distinctive summaries of articles. We have used LSTM recurrent cells to implement our encoder-decoder network that is used for sequence-to-sequence transformation. As a very special type of text summarization, headline generation is still a challenging task. It can be observed that, even though the nature of this task involves plenty of training data (lots of articles being published every day), only a very specific part of the published articles is essential for composing the article headline. Therefore, in order to generate a consistent and explanatory headline for an article, it is required to perform a successful selection of highlights in the context. In our work, we employ attention mechanism in order to selectively search the text for such highlights. We have used BERT, a state of the art language model, in order to generate contextualized word embeddings. Thus we can effectively exploit the context concept which is essential for the headline generation task. Also, instead of directly generating headlines from a whole article, first we generate several summaries of the article using pre-trained summarization methods. Then, we adaptively select the best summary for the article and use it as the input sequence for our headline generation model. In order to

choose best summary, we use a TF-IDF based method. We use BLEU, a frequently used evaluation metric used in machine translation, in order to compare our generated headlines with original labeled headlines.

We implement our model using *pytorch* framework. We train and evaluate our model on the all-the-news dataset (available at: <https://www.kaggle.com/snapcrack/all-the-news>)

II. RELATED WORK

There exist numerous studies on headline generation. These can be classified under three main approaches rule-based, statistics-based, and summarization-based. In this work we focus on summarization approach and think headlines as short summaries of texts. There are two main approaches to generating summary of a text, extraction-based and abstraction-based. Extraction based summaries are refined sentences formed by the words in text while abstraction based methods aims to generate summaries using paraphrasing and shortening.

One of the key example study of abstraction-based approach is [1] where they used neural encoder-decoder architecture. They used attention mechanism between first sentence of the article and given headlines. After the power of attention mechanism is shown, [2] implemented some side properties of language, such as abstract meaning, as an extension to attention based model. They used abstract meaning representation (AMR) parser to modify attention based model and improved the headline generation benchmarks compared with the baseline neural attention based model. Another extension to [1] made by [3]. They proposed conditional recurrent neural network (RNN) and designed a convolutional attention based encoder. Encoder computes scores over the words of input sentence considering position in the sentence and then these scores informs the decoder which part of the input sentence it should focus on and generate the next word. However, in [4] showed that attention based encoder-decoder approaches can be problematic due to unrelated filling of details when detail is missing in news. Another assumption that these [1], [2], [3], [4] works are making is, first sentences are the most informative. [5] showed that important information is distributed across several sentences and proposed multi sentence comparison model. A recent approach [6] proposed a model architecture where it takes multiple summaries as input to encoder and pass output states through control layer which modifies attention equation and feed resulting context

vector to decoder. These summaries can be extraction or abstraction based. In the control model they used Long-Short Term Memory (LSTM) units and suggest that this provides an hierarchy between summaries. However they did not mention LSTM functionality in detail.

III. PROBLEM DEFINITION

We can represent our dataset as $\mathcal{D} = \{\mathbf{X}_i, \mathbf{Y}_i\}_{i=1}^N$ where $|N|$ is the total number of news in our dataset. Let $\mathbf{X}_i = \{\mathbf{x}_j\}_{j=1}^T$ denote $|T|$ length of input text. Let $\mathbf{Y}_i = \{\mathbf{y}_k\}_{k=1}^{T'}$ represent the title of the news and $T < T'$. Where each \mathbf{x}_j and \mathbf{y}_k are one hot vectors of length $|V|$ where V is the vocabulary and $0 \leq j \leq T$, $0 \leq k \leq T'$. First we create candidate summaries $\{\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_K\}$ where $|K|$ is the number of candidate summaries. For each summary $\mathbf{S} = \{\mathbf{s}_l\}_{l=1}^L$ where $|L|$ is the length of the summary text there are synthesized words $\mathbf{s}_l \in V$. Our goal is to learn,

$$P(\mathbf{Y}|\mathbf{X}) = P(\mathbf{Y}|\{\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_K\}) \quad (1)$$

by maximising,

$$\hat{\mathbf{Y}} = \operatorname{argmax}_{\mathbf{Y}} P(\mathbf{Y}|\{\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_K\}) \quad (2)$$

where $\hat{\mathbf{Y}}$ represents predicted headline.

IV. PROPOSED METHODOLOGY

The methodology we followed in the main framework consists of encoder - decoder architecture with attention mechanism. As it is well known, such kind of architectures are special form of recurrent neural networks. Before examining how we implement this model, it would be better to clarify how the data is utilized and prepared for training.

A. Data & Pre-processing

In our project, the name of the used dataset is "All The News". It can be obtained from <https://www.kaggle.com/snapcrack/all-the-news>. This dataset contains news from several subjects such as politics, sport, music etc. There are approximately 140 thousands news and their respective titles. 10 % of the data was allocated as validation set, 10 % of the data was allocated as validation set and the remaining 80 % was used for training purposes.

First of all the dataset was subjected to a pre-processing. In pre-processing, all noisy characters and punctuation marks were removed. All letters are made lowercase and stop words are expelled from the dataset. Additionally, maximum word length is limited to 20 characters and $\langle \text{start} \rangle$ and $\langle \text{end} \rangle$ tokens were added in order to where the news start or finish.

Before summarization process it would be beneficial to mention some features of the news articles and titles. There are 142471 articles in our dataset. Every article has its unique title. However, 15 of the has no titles. This situation may stem from some errors while collecting the data. Since the number of articles without titles is very low compared to total size of dataset. The effect of removal of erroneous data is negligible.

Besides these general statistics, we also calculated average title and article length. Furthermore, it is important to know how many titles are extracted from the article sentences as a whole or paraphrased. How many words are common among titles and articles. These statistics can be observed in Table I.

TABLE I: Dataset Statistics

Average title length (in words)	10.2
Average article length (in words)	732.4
Average number of words in titles contained by article (per-cent)	73.2
Number of articles contains whole title	633

Here, we see that the number of titles which exactly come from articles are relatively low compared to total size of dataset. This situation shows that article titles use paraphrasing for important sentences. That's why learning the contexts of words while generating titles have significance.

After pre-processing, next operation is that summarizing the news. As it can be understood from the title of the project, our main objective is to generate and predict headlines from the summaries of the news. The reason why we used summaries instead of raw news is that feeding the network only by raw text would be too costly in terms of computation. Furthermore, whole news articles contain not only important parts of the news but also relatively redundant words too. Therefore, generating articles from the significant parts of the news would lead to obtain more informative titles. Nevertheless, "All The News" dataset does not contain summaries of news. For this reason, we create summaries for each news. While summarizing the news we utilized "LexRankSummarizer", "LsaSummarizer", "SumBasicSummarizer" libraries. Nonetheless, we need to choose the best reflecting summary among these three candidate. In order to choose best summary, we used a TF-IDF based method. As it can be remembered from the lectures. TF-IDF is a method that indicates a word's importance and specificity for document where it belongs to. Term frequency of a word respective to a document is calculated as number of the word in that document divided by total word count in all documents.

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_k} \quad (3)$$

In equation 3, $n_{i,j}$ denotes the number of occurrences word i in document j . Number in the denominator denotes the total number of words in all documents. Inverse document frequency of a word is formulated as equation 4.

$$idf_i = \log \frac{|D|}{|d : t_i \in d|} \quad (4)$$

Here D represents the set of all documents, cardinality of D is divided by the number of documents in which the term t_i takes place at least one time. Then, logarithm of the resulting value is taken. TF-IDF score of word for a document is calculated by multiplication of $tf_{i,j}$ and idf_i . It is reasonable to infer that if a word occurs many times in a document, that word can be an important word for that document. If

this stated word does not take place in other documents so many times, one can understand that stated word is a specific word for the first documents. Such kind of words carries a lot of information compared to other words in the documents. Naturally, high occurrence number does not enlighten whole problem. For example, stop words such as "the, and, of" take place so many times in every document, so these words does not give any information about the documents. Therefore, not only term frequency but also inverse document frequency should be taken into consideration.

In our case, we have chosen best summary among three candidates with the help TF-IDF like method. We treated union of three candidates as a document because all three candidates are related to one news article. Then, every TF-IDF scores of words were calculated. For every summary candidate, TF-IDF scores of words are summed and divided by total number of words in summary candidate. Here, normalization process is critical. Since the number of words in summaries are not the same, summaries consisting of more words would result in having high total TF-IDF score. That's why total TF-IDF scores should be normalized by total number words in summaries. In these calculations, <start> and <end> tokens were not included.

B. Proposed Model

Until this point, followed methodology is related to preparation of dataset to initialize training. After making the dataset ready to train. next step is to map words into embeddings. The most common and traditional pre-trained word embeddings are GloVe [7] and word2vec [8]. Although these tools are very powerful, they are not state of the art embeddings anymore. Language models such as ELMO [9], BERT [10] or ELECTRA [11] implement transformer models and generate contextualized word embeddings. In title generation task, since context concept is essential property, we decided to used one of these embedding models. However, those architectures have some effects in terms of model size and computation power. That's why we used "easy-bert" model for mapping words into embeddings [12]. This model directly gives the regular embeddings calculated with BERT parameters for each token or sequence. In fact, original BERT model does not give individual embeddings for words. It provides a pre-trained language model for fine tuning of the various language task. On the other hand traditional word embeddings do not take context into consideration in sufficient level. easy-bert incorporates the useful sides of traditional and contextualized models [12]. Therefore, we could obtain contextualized models for far cheaper way.

We have implemented a recurrent encoder-decoder neural network architecture with attention mechanism which is inspired by the Neural Machine Translation system idea that models the probability $p(y|x)$ of translating a source context (x_1, \dots, x_k) to a target context (y_1, \dots, y_m) . Our encoder-decoder neural network architecture consists of two parts: encoder es-

timates a representation z for each source context and decoder generates the headline word-by-word with log-probabilities:

$$\text{logp}(y_i|y_0, \dots, y_{i-1}, z) \quad (5)$$

Therefore, the modelled translation probability can be decomposed as:

$$\text{logp}(y|x) = \sum_i^m \text{logp}(y_i|y_0, \dots, y_{i-1}, z) \quad (6)$$

It is possible to use different recurrent architectures such as RNN, GRU or LSTM for generating the hidden representation z . Then, it is possible to decode each output word with probability:

$$p(y_i|y_0, \dots, y_{i-1}, z) = \text{softmax}(f(h_i)) \quad (7)$$

where h_i is the current hidden state representation of the used recurrent architecture and $f(\cdot)$ is a function that maps the hidden state representation to a vocabulary-sized vector. In our work, we have used LSTM as our recurrent architecture and all hidden states are computed throughout the LSTM unit. So far, the model only employs the encoder-decoder architecture. In order to employ the attention mechanism, we add an additional hidden state representation as follows

$$\tilde{h}_i = \tanh(W_c h_i^{\text{concat}}) \quad (8)$$

where $h_i^{\text{concat}} = [c_i; h_i]$ and c_i is the context vector that captures the relevant information in source context for the attention mechanism to help generating the target headline. W_c is a linear operation that maps the concatenated representation back to the original hidden state dimensions. Then, it is possible to write this additional representation into equation (7):

$$p(y_i|y_0, \dots, y_{i-1}, z) = \text{softmax}(W_f \tilde{h}_i) \quad (9)$$

where W_f is a matrix operator representation for the function $f(\cdot)$. Then in order to find context vectors c_i we use the global attention approach where all encoder hidden states are used while calculating c_i . Current hidden states h_i are then compared with source hidden states h_j^{source} in order to derive an alignment vector a_i with size equals to the source context length:

$$(a_i)_j = \frac{\exp(\text{score}(h_i, h_j^{\text{source}}))}{\exp(\sum_{j'} \text{score}(h_i, h_{j'}^{\text{source}}))} \quad (10)$$

where $\text{score}(\cdot, \cdot)$ is a context based function defined as:

$$\text{score}(h_i, h_j^{\text{source}}) = h_i^T W_a h_j^{\text{source}} \quad (11)$$

Then, the context-vector c_i is computed as a weighted sum of all encoder hidden states where the weights are the related entries of the alignment vector a_i . Figure 1 illustrates the attention mechanism that is being explained. After combining the encoder-decoder architecture with the attention mechanism, we have used an input feeding approach where the output vectors that composes the generated headline are concatenated with the source context sequence and feed back to the encoder-decoder architecture in the following time steps. Figure 2

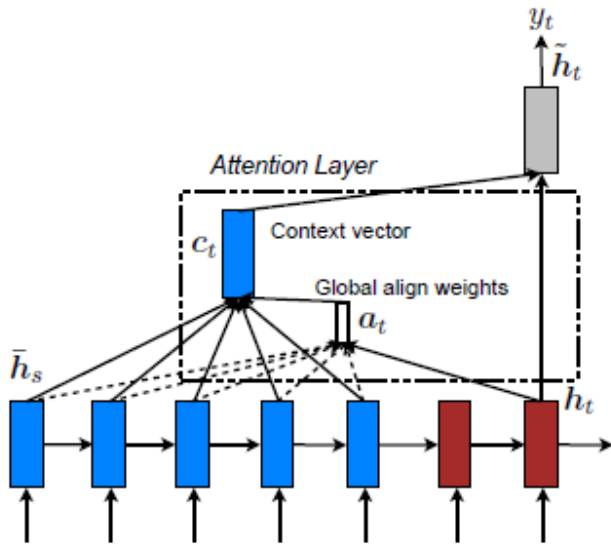


Fig. 1: Attention Mechanism

illustrates the entire model architecture together with input feeding architecture.

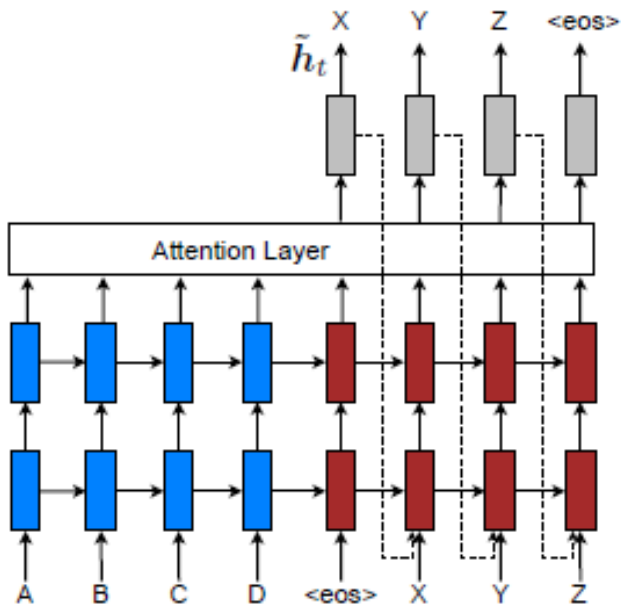


Fig. 2: Proposed Model

C. Evaluation

In order to compare our generated titles with original labeled ones we need a similarity metric. In such kind of applications BLEU metric is used frequently [13]. This metric is developed for machine translation tasks first but used in many different NLP tasks. BLEU is language independent metric and coincides highly with human evaluation. BLEU mainly, looks for matched n-grams between reference and candidate text.



Fig. 3: Training Loss

TABLE II: Sample Generated Headlines

Original Title	georgia man shoots another times hot pocket snapped
Predicted Title	georgia man pleads
Original Title	london mayor boris johnson backs brexit boosting campaign
Predicted Title	brexit britain theresa thatcher
Original Title	conservatives embrace principled populism
Predicted Title	trump populism trumpisms
Original Title	know every sex city outfit
Predicted Title	city outfit
Original Title	solid jobs created april unemployment falls
Predicted Title	unemployment falls

Here, n-gram concept is the same as the one we learned in lectures. However, comparison is positional independent. In addition reference text may comprise of multiple sentences or sentence blocks. BLEU score of a candidate with respect to the reference lies between 0.0 and 1.0. As it can be expected, value 1.0 represents complete match (it means that candidate is the same as reference) while value 0.0 indicates no correlation between candidate and reference [13]. In our project, original labels were used as references. Scores of candidates (our generated titles) is used for evaluation of success of our model. Obtained results will be demonstrated in next section.

V. RESULTS

Our loss curve is shown in Figure 3. We have observed that as the training process increases our model tends to over-fit. However we prevent it by using early stop and implementing regularisation techniques such as drop-out layer. In our training we followed teacher forcing technique. For evaluation our predictions are recurrently produced as shown in Figure 2. From the experiments made on the test set we have achieved **0.35** BLEU-1 score. Table II shows the original title and our predictions.

VI. DISCUSSION & CONCLUSION

In conclusion, we have implemented headline (or title) generator for news article. Since news articles are relatively long sequences. We used neural models with keeping memory.

For this purpose, we utilized RNNs with attention mechanism. Actually, RNNs are not successful in storing long term memory. In summarization tasks, since the input coming through model are long paragraphs, long term memory is a required feature for such kind of tasks. In this project, our main objective is to choose best informative text from summarizers and feed to our model which can keep long term memory. After this process, we would obtain best title.

Naturally, there are some drawbacks while generating headlines from scratch. Some possible, errors stem from similar news articles. Some keywords such as "Trump", "politics", "action" etc. take place a lot of times. Although news articles are gathered from different sources, actors of the news are same. This may make some words to have high number of occurrences. Therefore such words may be repeated at generated headlines. On the other hand, in such kind of tasks, implementing long sequences requires more complex models. We tried to overcome this issue by using summaries of the articles. Based on the summarization techniques performance of the generated headlines may change. As a future work, this problem may be examined. Overall, we have met our requirements and obtained reasonable results for this task as shown in "Results" section.

REFERENCES

- [1] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, (Lisbon, Portugal), pp. 379–389, Association for Computational Linguistics, Sept. 2015.
- [2] S. Takase, J. Suzuki, N. Okazaki, T. Hirao, and M. Nagata, "Neural headline generation on abstract meaning representation," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, (Austin, Texas), pp. 1054–1059, Association for Computational Linguistics, Nov. 2016.
- [3] S. Chopra, M. Auli, and A. M. Rush, "Abstractive sentence summarization with attentive recurrent neural networks," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, (San Diego, California), pp. 93–98, Association for Computational Linguistics, June 2016.
- [4] K. Lopyrev, "Generating news headlines with recurrent neural networks," *arXiv preprint arXiv:1512.01712*, 2015.
- [5] E. Alfonseca, D. Pighin, and G. Garrido, "Heady: News headline abstraction through event pattern clustering," in *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1243–1253, 2013.
- [6] J. Tan, X. Wan, and J. Xiao, "From neural sentence summarization to headline generation: A coarse-to-fine approach," in *IJCAI*, pp. 4109–4115, 2017.
- [7] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, (Doha, Qatar), pp. 1532–1543, Association for Computational Linguistics, Oct. 2014.
- [8] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems 26* (C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, eds.), pp. 3111–3119, Curran Associates, Inc., 2013.
- [9] M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, (New Orleans, Louisiana), pp. 2227–2237, Association for Computational Linguistics, June 2018.

- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, (Minneapolis, Minnesota), Association for Computational Linguistics, June 2019.
- [11] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "ELECTRA: Pre-training text encoders as discriminators rather than generators," in *ICLR*, 2020.
- [12] R. Rua, "easy-bert," Apr. 2019.
- [13] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, (Philadelphia, Pennsylvania, USA), pp. 311–318, Association for Computational Linguistics, July 2002.

VII. APPENDIX

A. Contributions

We all equally contributed to the project, here is the individual contribution of each partner:

- : Burak Künkçü: Dataset preparation, helping model implementation, model explanation in report
- Selim Furkan Tekin: Helping the implementation of the model. Training, testing and obtaining the results. Problem definition and related work.
- Furkan Şahinuç: Implementing TF-IDF algorithm on candidate summaries. Initialization of pre-trained vectors which are input to the model. First part of the methodology and discussion parts of the report.

```

1 import os
2 import random
3 import pandas as pd
4 import pickle as pkl
5
6 from collections import Counter
7 from transformers.preprocessing import Preprocess
8 from transformers.summarizer import Summarizer
9 from transformers.summary_selector import SummarySelector
10
11
12 class LoadData:
13     def __init__(self, dataset_path, **data_params):
14         self.summary_path = os.path.join(dataset_path, 'summary_set.pkl')
15         self.vocab_path = os.path.join(dataset_path, 'vocab.pkl')
16         self.all_sum_save_path = os.path.join(dataset_path, 'all_sum_set.pkl')
17         self.raw_data_path = os.path.join(dataset_path, 'all-the-news')
18
19         self.content_len = data_params.get('content_len', 50)
20         self.title_len = data_params.get('title_len', 15)
21         self.num_samples = data_params.get('num_samples', -1)
22         self.sentence_num = data_params.get('num_sentence', 3)
23         self.unk_threshold = data_params.get('unk_threshold', 10)
24
25         self.test_ratio = data_params.get('test_ratio', 0.1)
26         self.val_ratio = data_params.get('val_ratio', 0.1)
27         self.shuffle = data_params.get('shuffle', True)
28
29         if not os.path.isfile(self.vocab_path) or not os.path.isfile(self.summary_path):
30             contents, titles = self.__load_from_csv()
31             print('\nCreating data ...')
32             self.summaries, self.titles, self.word2int, self.int2word =
33                 ↪ self.__create_summary_set(contents, titles)
34
35             with open(self.vocab_path, 'wb') as f:
36                 pkl.dump([self.word2int, self.int2word], f)
37
38             with open(self.summary_path, 'wb') as f:
39                 pkl.dump([self.summaries, self.titles], f)
40
41         else:
42             print('\nLoading data from pickle ...')
43             with open(self.vocab_path, 'rb') as f:
44                 self.word2int, self.int2word = pkl.load(f)
45             with open(self.summary_path, 'rb') as f:
46                 self.summaries, self.titles = pkl.load(f)
47
48         # split test train validation
49         self.data_dict, self.label_dict = self.__split_data()

```

```

1  def __load_from_csv(self):
2      # read all articles 1, 2 and 3
3      file_paths = os.listdir(self.raw_data_path)
4      file_paths = [os.path.join(self.raw_data_path, file_name) for file_name in file_paths]
5
6      df_list = []
7      for file in file_paths:
8          df = pd.read_csv(file, index_col=0)
9          df_list.append(df)
10     articles = pd.concat(df_list, ignore_index=True)
11
12     contents = articles['content'].values[:self.num_samples]
13     titles = articles['title'].values[:self.num_samples]
14
15     return contents, titles
16
17 def __create_summary_set(self, in_contents, in_titles):
18     process_obj = Preprocess(content_len=self.content_len, title_len=self.title_len)
19
20     print("Creating summaries ...")
21     summarizer = Summarizer(num_sentence=self.sentence_num)
22     sum_collection, sum_titles = summarizer.transform(X=in_contents, y=in_titles)
23
24     summary_label = []
25     summary_content = []
26     for i in range(summarizer.num_summarizers):
27         titles = []
28         contents = []
29         for sum_list, title in zip(sum_collection, sum_titles):
30             contents.append(sum_list[i])
31             titles.append(title)
32
33         print('\nPreprocessing summary: {}'.format(i))
34         contents, titles = process_obj.transform(X=contents, y=titles)
35
36         summary_content.append(contents)
37         summary_label.append(titles)
38
39     with open(self.all_sum_save_path, 'wb') as f:
40         pickle.dump([summary_content, summary_label], f)
41
42     print('\nSelecting summaries ...')
43     selector = SummarySelector()
44     summary_content = selector.transform(summary_content)
45     summary_label = summary_label[0]
46
47     # create vocab
48     word2int, int2word = self.__create_vocab(summary_content, summary_label)
49
50     content_int = [[word2int[word] for word in content] for content in summary_content]
51     title_int = [[word2int[word] for word in title] for title in summary_label]
52
53     return content_int, title_int, word2int, int2word

```

```

1  def __create_vocab(self, contents, titles):
2      contents, titles = self.__rmv_less_frequent(contents, titles)
3      vocab = self.__get_counts(contents, titles)
4
5      word2int = {k: v for v, k in enumerate(vocab.keys())}
6      int2word = {v: k for k, v in word2int.items()}
7
8      return word2int, int2word
9
10 def __rmv_less_frequent(self, contents, titles):
11     vocab = self.__get_counts(contents, titles)
12
13     # replace the less frequent words with '<unk>'
14     for i in range(len(contents)):
15         for j in range(len(contents[i])):
16             word = contents[i][j]
17             if vocab[word] <= self.unk_threshold:
18                 contents[i][j] = '<unk>'
19         for k in range(len(titles[i])):
20             word = titles[i][k]
21             if vocab[word] <= self.unk_threshold:
22                 titles[i][k] = '<unk>'
23
24     return contents, titles
25
26 def __split_data(self):
27     dataset_length = len(self.summaries)
28
29     if self.shuffle:
30         zipped = list(zip(self.summaries, self.titles))
31         random.shuffle(zipped)
32         self.summaries, self.titles = zip(*zipped)
33         self.summaries = list(self.summaries)
34         self.titles = list(self.titles)
35
36     test_count = int(dataset_length * self.test_ratio)
37     val_count = int(dataset_length * self.val_ratio)
38
39     data_dict = dict()
40     data_dict['test'] = self.summaries[:test_count]
41     data_dict['validation'] = self.summaries[test_count:test_count + val_count]
42     data_dict['train'] = self.summaries[test_count + val_count:]
43
44     label_dict = dict()
45     label_dict['test'] = self.titles[:test_count]
46     label_dict['validation'] = self.titles[test_count:test_count + val_count]
47     label_dict['train'] = self.titles[test_count + val_count:]
48
49     return data_dict, label_dict
50
51 @staticmethod
52 def __get_counts(contents, titles):
53     # create_vocab
54     all_words = []
55     for c, t in zip(contents, titles):
56         all_words += c + t
57
58     vocab = Counter(all_words)
59     vocab = {k: v for k, v in sorted(vocab.items(), key=lambda x: x[1], reverse=True)}
60
61     return vocab

```



```

1 import os
2 import torch
3 import numpy as np
4 import torch.nn as nn
5
6 from transformers.word2vec import Word2VecTransformer
7
8
9 class Embedding(nn.Module):
10     def __init__(self, int2word, device, embed_name='glove'):
11         nn.Module.__init__(self)
12         self.vector_path = 'embedding/embed_{}.npy'.format(embed_name)
13         self.embed_name = embed_name
14
15         self.int2word = int2word
16         self.vocab_size = len(int2word)
17
18         if os.path.isfile(self.vector_path):
19             print('Loading saved embedding vectors')
20             self.weights = np.load(self.vector_path, allow_pickle=True)
21             self.weights = torch.from_numpy(self.weights)
22         else:
23             print('Creating embed tensor')
24             self.weights = self.create_embed_tensor()
25             np.save(self.vector_path, self.weights)
26
27         self.embed_dim = self.weights.shape[1]
28         self.embedding = nn.Embedding.from_pretrained(self.weights).to(device)
29
30     def forward(self, captions):
31         """
32         :param captions: (B, S)
33         :return: vectors: (B, S, Vector_dim)
34         """
35         vectors = self.embedding(captions)
36         return vectors
37
38     def create_embed_tensor(self):
39         transformer = Word2VecTransformer(self.embed_name)
40         vectors = []
41         for i in range(self.vocab_size):
42             print("\r{:.2f}%".format(i * 100 / self.vocab_size), flush=True, end='')
43             vec = transformer.transform(self.int2word[i])
44             vectors.append(vec)
45         embed = np.stack(vectors, axis=0)
46         return torch.from_numpy(embed)
47
48
49 if __name__ == '__main__':
50     import pickle
51
52     vocab = pickle.load(open('data/vocab.pkl', 'rb'))
53     word2int, int2word = vocab
54     Embedding(int2word)

```

```

1 import random
2 import torch
3 import torch.nn as nn
4
5 from models.embed import Embedding
6
7
8 class Encoder(nn.Module):
9     def __init__(self, vocab, hidden_dim, num_layers, device, embed_name):
10         super(Encoder, self).__init__()
11         self.num_layers = num_layers
12         self.hidden_dim = hidden_dim
13
14         self.embedding_layer = Embedding(vocab, device, embed_name=embed_name)
15         self.embed_dim = self.embedding_layer.embed_dim
16
17         self.rnn = nn.LSTM(input_size=self.embed_dim,
18                             hidden_size=self.hidden_dim,
19                             num_layers=num_layers,
20                             batch_first=True)
21
22     def forward(self, input_tensor):
23         """
24         :param input_tensor: (B,)
25         :param hidden: ((L, b, d), (L, b, d))
26         :return:
27         """
28         embed = self.embedding_layer.embedding(input_tensor).float()
29         output, hidden = self.rnn(embed)
30
31         return output, hidden
32
33
34 class Attention(nn.Module):
35     def __init__(self, enc_hidden, dec_hidden):
36         super(Attention, self).__init__()
37         self.enc_hidden = enc_hidden
38         self.dec_hidden = dec_hidden
39
40         self.w_in = nn.Linear(self.dec_hidden, self.enc_hidden, bias=False)
41         self.w_out = nn.Linear(self.dec_hidden+self.enc_hidden, self.dec_hidden)
42
43         self.softmax = nn.Softmax()
44         self.tanh = nn.Tanh()
45
46     def forward(self, dec_out, enc_out):
47         """
48         :param torch.Tensor dec_out: (batch, t_out, d_dec)
49         :param torch.Tensor enc_out: (batch, t_in, d_enc)
50         :return:
51         """
52         batch, t_out, d_dec = dec_out.shape
53         batch, t_in, d_enc = enc_out.shape
54
55         dec_out_ = dec_out.contiguous().view(batch*t_out, d_dec)
56         energy = self.w_in(dec_out_)
57         energy = energy.view(batch, t_out, d_enc)
58
59         # swap dimensions: (batch, d_enc, t_in)
60         enc_out_ = enc_out.transpose(1, 2)
61
62         # (batch, t_out, d_enc) x (batch, d_enc, t_in) -> (batch, t_out, t_in)
63         attn_energies = torch.bmm(energy, enc_out_)
64
65         alpha = self.softmax(attn_energies.view(batch*t_out, t_in))
66         alpha = alpha.view(batch, t_out, t_in)
67
68         # (batch, t_out, t_in) * (batch, t_in, d_enc) -> (batch, t_out, d_enc)
69         context = torch.bmm(alpha, enc_out)
70
71         #  $\hat{h}_t = \tanh(W [\text{context}, \text{dec\_out}])$ 
72         concat_c = torch.cat([context, dec_out], dim=2)
73         concat_c = concat_c.view(batch*t_out, d_enc + d_dec)
74         concat_c = self.w_out(concat_c)
75         attn_h = self.tanh(concat_c.view(batch, t_out, d_dec))

```

```

1  class Decoder(nn.Module):
2      def __init__(self, vocab, hidden_dim, enc_hidden, num_layers, dropout_prob, device,
3          ↪ embed_name):
4          super(Decoder, self).__init__()
5          self.hidden_dim = hidden_dim
6          self.enc_hidden = enc_hidden
7          self.num_layers = num_layers
8          self.drop_prob = dropout_prob
9
10         self.embedding_layer = Embedding(vocab, device, embed_name=embed_name)
11         self.embed_dim = self.embedding_layer.embed_dim
12         self.rnn = nn.LSTM(input_size=self.embed_dim,
13             hidden_size=self.hidden_dim,
14             num_layers=self.num_layers,
15             batch_first=True)
16
17         self.attention = Attention(dec_hidden=self.hidden_dim, enc_hidden=self.enc_hidden)
18         self.out_lin = nn.Linear(self.hidden_dim, self.embedding_layer.vocab_size)
19         self.dropout = nn.Dropout(self.drop_prob)
20
21     def forward(self, input_tensor, hidden, enc_out):
22         # Teacher-forcing,
23         input_tensor = input_tensor.unsqueeze(1)
24         embed = self.embedding_layer.embedding(input_tensor).float()
25         embed = self.dropout(embed)
26
27         dec_out, dec_hid = self.rnn(embed, hidden)
28         attn_out, alpha = self.attention(dec_out, enc_out)
29         output = self.out_lin(attn_out)
30
31         return output, hidden
32
33 class Seq2Seq(nn.Module):
34     def __init__(self, vocabs, device, **model_params):
35         super(Seq2Seq, self).__init__()
36         self.word2int, self.intword = vocabs
37         self.enc_hidden = model_params['encoder_hidden_dim']
38         self.dec_hidden = model_params['decoder_hidden_dim']
39         self.enc_num_layer = model_params['encoder_num_layer']
40         self.dec_num_layer = model_params['decoder_num_layer']
41         self.drop_prob = model_params['dropout_prob']
42         self.embed_name = model_params['embed_name']
43         self.device = device
44
45         self.encoder = Encoder(vocab=self.intword,
46             hidden_dim=self.enc_hidden,
47             num_layers=self.enc_num_layer,
48             device=self.device,
49             embed_name=self.embed_name)
50         self.decoder = Decoder(vocab=self.intword,
51             hidden_dim=self.dec_hidden,
52             enc_hidden=self.enc_hidden,
53             num_layers=self.dec_num_layer,
54             dropout_prob=self.drop_prob,
55             device=self.device,
56             embed_name=self.embed_name)
57
58     def forward(self, contents, titles, tf_ratio=0.0):
59         """
60
61         :param contents: (b, t')
62         :param titles: (b, t)
63         :param tf_ratio: teacher forcing ratio
64         :return:
65         """
66         batch, title_len = titles.shape
67
68         enc_out, hidden = self.encoder(contents)
69
70         # <start> is mapped to id=2
71         dec_inputs = torch.ones(batch, dtype=torch.long) * self.word2int['<start>']
72         dec_inputs = dec_inputs.to(self.device)
73         pred, hidden = self.decoder(dec_inputs, hidden, enc_out)
74

```

```

1 import re
2 import emoji
3 import string
4 import nltk
5
6 from nltk.tokenize import TweetTokenizer
7 from nltk.corpus import stopwords
8
9 eng_stopwords = set(stopwords.words('english'))
10
11
12 class Preprocess:
13     def __init__(self, content_len, title_len):
14         self.content_len = content_len
15         self.title_len = title_len
16         self.tokenizer = TweetTokenizer()
17
18         # remove url, long-mention, mention, hash-tag, numbers
19         self.re_list = [r'http\S+', r'\(@ ?[^\s].+)\', r'@ ?[^\s]+',
20                         r'# ?[^\s]+', '{}'.format(re.escape(string.punctuation))]
21
22         # replace '- The New York Times'
23         self.re_title = r'\-[a-z A-Z]+'
24
25         # before tokenizing, remove non textual emoji
26         self.rmv_emoji = lambda x: emoji.get_emoji_regexp().sub(r'', x)
27
28         # after tokenizing
29         self.rmv_stop = lambda x: [w for w in x if w not in eng_stopwords]
30         self.rmv_pun = lambda x: [w for w in x if w not in string.punctuation]
31         self.rmv_short_long = lambda x: [w for w in x if 20 >= len(w) >= 3]
32
33         self.rmv_list = [self.rmv_stop, self.rmv_pun, self.rmv_short_long]
34
35     def transform(self, X, y):
36         """
37         X is the content
38         y is the title
39
40         :param X: [str, str, ..., str]
41         :param y: [str, str, ..., str]
42         :return: [[token, ..., token], ...], [[token, ..., token], ...]
43         """
44         labels = []
45         clean_data = []
46         for count, (content, title) in enumerate(zip(X, y)):
47             print('\r{:.2f}%'.format(count * 100 / len(X)), flush=True, end='')
48             token_content = self._preprocess(str(content).lower(), mode='content')
49             title_content = self._preprocess(str(title).lower(), mode='title')
50             clean_data.append(token_content)
51             labels.append(title_content)
52
53         return clean_data, labels
54
55     def _preprocess(self, sentence, mode='content'):
56         """
57         :param sentence: string
58         :return: [token, ..., token]
59         """
60         if mode == 'title':
61             self.re_list.insert(0, self.re_title)
62
63         # remove url, long-mention, mention, hash-tag, non-textual emoji, punctuation
64         for re_op in self.re_list:
65             sentence = re.sub(re_op, '', sentence)
66         sentence = self.rmv_emoji(sentence)
67
68         # tokenize the sentence
69         tokens = self.tokenizer.tokenize(sentence)
70         for rmv_op in self.rmv_list:
71             tokens = rmv_op(tokens)
72
73         # lower every word
74         tokens = [token.lower() for token in tokens]
75

```

```

1 from sumy.parsers.plaintext import PlaintextParser
2 from sumy.nlp.tokenizers import Tokenizer
3 from sumy.summarizers.lex_rank import LexRankSummarizer
4 from sumy.summarizers.lsa import LsaSummarizer
5 from sumy.summarizers.sum_basic import SumBasicSummarizer
6
7
8 class Summarizer:
9     def __init__(self, num_sentence, trim_len=5000):
10         self.num_sentence = num_sentence
11         self.trim_len = trim_len
12         self.tokenizer = Tokenizer('english')
13
14         self.summarizers = [LexRankSummarizer(), LsaSummarizer(), SumBasicSummarizer()]
15         self.num_summarizers = len(self.summarizers)
16
17     def transform(self, X, y):
18         """
19         Create 5 summaries for each content. If content has less than
20         self.num_sentence it is not added to collection
21
22         :param X: [str, str, ...]
23         :param y: [str, str, ...]
24         :return: X=[[sum1, sum2, ..., sum5], ...], y=[str, str, ...]
25         """
26         labels = []
27         sum_collection = []
28         for count, (content, title) in enumerate(zip(X, y)):
29             print('\r{:.2f}%'.format(count * 100 / len(X)), flush=True, end='')
30
31             trimmed_content = content[:self.trim_len]
32
33             parser = PlaintextParser.from_string(trimmed_content, self.tokenizer)
34
35             sum_list = []
36             long_content = True
37             for summarizer in self.summarizers:
38                 summary = summarizer(parser.document, self.num_sentence)
39                 if len(summary) >= self.num_sentence:
40                     sum_txt = ' '.join(sentence._text for sentence in summary)
41                     sum_list.append(sum_txt)
42             else:
43                 long_content = False
44
45             if long_content:
46                 labels.append(y[count])
47                 sum_collection.append(sum_list)
48
49     return sum_collection, labels

```

```

1  import math
2
3
4  class SummarySelector:
5      def __init__(self):
6          self.start = 0
7          self.end = 1
8
9      def transform(self, summaries):
10
11         vocab = dict()
12         vocab['<start>'] = self.start
13         vocab['<end>'] = self.end
14         count = 2
15         for ind in range(len(summaries)):
16             for row in summaries[ind]:
17                 for word in row:
18                     if word not in vocab:
19                         vocab[word] = count
20                         count += 1
21
22         idf = dict()
23
24         total_word_count = 0
25         total_document_count = len(summaries[0])
26
27         candidate_corpus = []
28
29         for num in range(len(summaries[0])):
30             temp = summaries[0][num] + summaries[1][num] + summaries[2][num]
31             candidate_corpus.append(temp)
32
33         for word in list(vocab.keys()):
34             belonging_doc = 0
35
36             for doc in candidate_corpus:
37                 if word in doc:
38                     belonging_doc += 1
39
40             if belonging_doc == 0:
41                 continue
42
43             idf[word] = math.log(total_document_count / belonging_doc)
44
45         for candidate in candidate_corpus:
46             total_word_count += len(candidate) - 6 # start end x3
47
48         print(total_word_count)
49
50         chosen_summaries = []
51         for num in range(len(summaries[0])):
52             scores = []
53             candidates = []
54             candidates.append(summaries[0][num][:])
55             candidates.append(summaries[1][num][:])
56             candidates.append(summaries[2][num][:])
57
58             for candidate in candidates:
59                 score = 0
60                 for word in candidate:
61                     if word == '<start>' or word == '<end>':
62                         continue
63                     score += (candidate_corpus[num].count(word) / total_word_count) * idf[word]
64
65             scores.append(score / (len(candidate) - 2))
66
67             best_index = scores.index(max(scores))
68
69             chosen_summaries.append(candidates[best_index])
70
71         return chosen_summaries

```

```

1 import os
2 import pickle as pkl
3 import numpy as np
4 from gensim.models import KeyedVectors
5
6 np.random.seed(42)
7
8
9 class Word2VecTransformer:
10     def __init__(self, embed_name='glove'):
11         project_dir = os.path.dirname(os.path.dirname(os.path.abspath(__file__)))
12         self.embed_name = embed_name
13         if embed_name == 'glove':
14             self.vector_path = os.path.join(project_dir, 'embedding/glove_dict.pkl')
15             self.model = pkl.load(open(self.vector_path, 'rb'))
16             self.vector_size = len(next(iter(self.model.values())))
17
18         elif embed_name == 'bert':
19             from easybert import Bert
20             self.vector_path = 'https://tfhub.dev/google/bert_multi_cased_L-12_H-768_A-12/1'
21             self.vector_size = 768
22
23             print('Downloading BERT Model')
24             self.model = Bert(self.vector_path)
25
26         elif embed_name == 'fasttext':
27             self.model_path = os.path.join(project_dir, 'embedding/fasttext.pkl')
28             self.vector_path = os.path.join(project_dir, 'embedding/fasttext.vec')
29
30             if os.path.isfile(self.model_path):
31                 model_file = open(self.model_path, 'rb')
32                 self.model = pkl.load(model_file)
33             else:
34                 print('Loading Fasttext model')
35                 self.model = KeyedVectors.load_word2vec_format(self.vector_path, binary=False,
36                                     ↪ unicode_errors='replace')
37                 model_file = open(self.model_path, 'wb')
38                 pkl.dump(self.model, model_file)
39
40             self.vector_size = self.model.vector_size
41
42         self.special_tokens = ['<start>', '<end>', '<pad>', '<unk>']
43         self.special_tok_dict = {}
44         for token in self.special_tokens:
45             self.special_tok_dict[token] = np.random.rand(self.vector_size)
46
47     def transform(self, x):
48         """
49         returns bert vector of input word
50
51         :param str x: input word
52         :return:
53         """
54         if x in self.special_tokens:
55             r = self.special_tok_dict[x]
56         else:
57             if self.embed_name == 'bert':
58                 with self.model:
59                     r = self.model.embed(x)
60             elif self.embed_name == 'glove':
61                 if x in self.model:
62                     r = self.model[x]
63                 else:
64                     r = np.random.rand(self.vector_size)
65             elif self.embed_name == 'fasttext':
66                 if x in self.model.wv:
67                     r = self.model.wv[x]
68                 else:
69                     r = np.random.rand(self.vector_size)
70
71         return r

```

```

1  from dataset import HeadlineDataset
2  from torch.utils.data import DataLoader
3
4
5  class BatchGenerator:
6      def __init__(self, data_dict, label_dict, **params):
7          self.data_dict = data_dict
8          self.label_dict = label_dict
9
10         self.batch_size = params.get('batch_size', 16)
11         self.num_works = params.get('num_works', 4)
12         self.shuffle = params.get('shuffle', True)
13
14         self.dataset_dict, self.dataloader_dict = self.__create_data()
15
16     def generate(self, data_type):
17         """
18         :param data_type: can be 'test', 'train' and 'validation'
19         :return: img tensor, label numpy_array
20         """
21         selected_loader = self.dataloader_dict[data_type]
22         yield from selected_loader
23
24     def num_iter(self, data_type):
25         dataset = self.dataset_dict[data_type]
26         return len(dataset) // self.batch_size
27
28     def __create_data(self):
29
30         im_dataset = {}
31         for i in ['test', 'train', 'validation']:
32             im_dataset[i] = HeadlineDataset(articles=self.data_dict[i],
33                                           titles=self.label_dict[i])
34
35         im_loader = {}
36         for i in ['test', 'train', 'validation']:
37             im_loader[i] = DataLoader(im_dataset[i],
38                                     batch_size=self.batch_size,
39                                     shuffle=self.shuffle,
40                                     num_workers=self.num_works,
41                                     drop_last=True)
42
43         return im_dataset, im_loader

```



```
1 model_params = {
2     'encoder_hidden_dim': 512,
3     'decoder_hidden_dim': 512,
4     'encoder_num_layer': 1,
5     'decoder_num_layer': 1,
6     'dropout_prob': 0.3,
7     'embed_name': 'fasttext'
8 }
9
10 train_params = {
11     'num_epoch': 50,
12     'learn_rate': 0.0003,
13     'train_tf_ratio': 0.5,
14     'val_tf_ratio': 0.1,
15     'clip': 5,
16     'eval_every': 400
17 }
18
19 data_params = {
20     "content_len": 50,
21     "title_len": 15,
22     "num_samples": 100,
23     "num_sentence": 3,
24     "test_ratio": 0.1,
25     "val_ratio": 0.1,
26     "shuffle": True,
27     "unk_threshold": 15
28 }
29
30 batch_params = {
31     'batch_size': 128,
32     'num_works': 0,
33     'shuffle': True,
34 }
```

```

1 import os
2 import random
3 import pandas as pd
4 import pickle as pkl
5
6 from collections import Counter
7 from transformers.preprocessing import Preprocess
8 from transformers.summarizer import Summarizer
9 from transformers.summary_selector import SummarySelector
10
11
12 class LoadData:
13     def __init__(self, dataset_path, **data_params):
14         self.summary_path = os.path.join(dataset_path, 'summary_set.pkl')
15         self.vocab_path = os.path.join(dataset_path, 'vocab.pkl')
16         self.all_sum_save_path = os.path.join(dataset_path, 'all_sum_set.pkl')
17         self.raw_data_path = os.path.join(dataset_path, 'all-the-news')
18
19         self.content_len = data_params.get('content_len', 50)
20         self.title_len = data_params.get('title_len', 15)
21         self.num_samples = data_params.get('num_samples', -1)
22         self.sentence_num = data_params.get('num_sentence', 3)
23         self.unk_threshold = data_params.get('unk_threshold', 10)
24
25         self.test_ratio = data_params.get('test_ratio', 0.1)
26         self.val_ratio = data_params.get('val_ratio', 0.1)
27         self.shuffle = data_params.get('shuffle', True)
28
29         if not os.path.isfile(self.vocab_path) or not os.path.isfile(self.summary_path):
30             contents, titles = self.__load_from_csv()
31             print('\nCreating data ...')
32             self.summaries, self.titles, self.word2int, self.int2word =
33                 ↪ self.__create_summary_set(contents, titles)
34
35             with open(self.vocab_path, 'wb') as f:
36                 pkl.dump([self.word2int, self.int2word], f)
37
38             with open(self.summary_path, 'wb') as f:
39                 pkl.dump([self.summaries, self.titles], f)
40
41         else:
42             print('\nLoading data from pickle ...')
43             with open(self.vocab_path, 'rb') as f:
44                 self.word2int, self.int2word = pkl.load(f)
45             with open(self.summary_path, 'rb') as f:
46                 self.summaries, self.titles = pkl.load(f)
47
48         # split test train validation
49         self.data_dict, self.label_dict = self.__split_data()
50
51     def __load_from_csv(self):
52         # read all articles 1, 2 and 3
53         file_paths = os.listdir(self.raw_data_path)
54         file_paths = [os.path.join(self.raw_data_path, file_name) for file_name in file_paths]
55
56         df_list = []
57         for file in file_paths:
58             df = pd.read_csv(file, index_col=0)
59             df_list.append(df)
60         articles = pd.concat(df_list, ignore_index=True)
61
62         contents = articles['content'].values[:self.num_samples]
63         titles = articles['title'].values[:self.num_samples]
64
65         return contents, titles
66
67     def __create_summary_set(self, in_contents, in_titles):
68         process_obj = Preprocess(content_len=self.content_len, title_len=self.title_len)
69
70         print("Creating summaries ...")
71         summarizer = Summarizer(num_sentence=self.sentence_num)
72         sum_collection, sum_titles = summarizer.transform(X=in_contents, y=in_titles)
73
74         summary_label = []
75         summary_content = []

```

```

1 import sys
2 import pickle as pkl
3
4 from config import data_params, batch_params, model_params, train_params
5 from load_data import LoadData
6 from batch_generator import BatchGenerator
7 from trainer import train
8 from tester import test
9
10
11 def main(mode):
12     dataset_path = 'data'
13     data = LoadData(dataset_path=dataset_path, **data_params)
14
15     print('Creating Batch Generator...')
16     batch_gen = BatchGenerator(data_dict=data.data_dict,
17                               label_dict=data.label_dict,
18                               **batch_params)
19
20     if mode == 'train':
21         train(vocabs=[data.word2int, data.int2word],
22              batch_gen=batch_gen,
23              train_params=train_params,
24              model_params=model_params)
25
26     elif mode == 'test':
27         print('Loading model')
28         model_file = open('results/seq2seq.pkl', 'rb')
29         model = pkl.load(model_file)
30         print('Testing model...')
31         test(model, data.int2word, batch_gen)
32
33
34 if __name__ == '__main__':
35     if len(sys.argv) == 1:
36         print('no arguments given, default process of training has started')
37         run_mode = 'train'
38     else:
39         run_mode = sys.argv[1]
40     main(run_mode)

```

```

1 import torch
2 import numpy as np
3
4 from transformers.bleu import compute_bleu
5 from trainer import translate
6
7
8 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
9
10
11 def test(net, vocab, batch_gen):
12     net.eval()
13     net.to(device)
14
15     bleu_all_score = []
16     for idx, (content, title) in enumerate(batch_gen.generate('test')):
17
18         print('\rttest:{}'.format(idx), flush=True, end='')
19
20         content, title = content.to(device), title.to(device)
21         title_pre = net(content, title, tf_ratio=0.1)
22
23         referance_caps, translate_caps = translate(outputs=[title, title_pre],
24                                                    int2word=vocab,
25                                                    top_k=10,
26                                                    print_count=10)
27
28         print('\n')
29         for i in range(len(referance_caps)):
30             bleu_score_list = []
31             for j in range(1, 5):
32                 bleu, _, _ = compute_bleu(referance_caps[i], translate_caps[i], max_order=j)
33                 bleu_score_list.append(bleu)
34             bleu_all_score.append(bleu_score_list)
35
36     bleu_all_score = np.array(bleu_all_score)
37     bleu_score_arr = np.mean(bleu_all_score, axis=0)
38
39     for i in range(4):
40         print("BLUE {}: {:.2f}".format(i+1, bleu_score_arr[i]))

```

```

1 import torch
2 import pickle
3 import numpy as np
4 import collections
5 import torch.nn as nn
6 import torch.optim as optim
7 import torch.nn.functional as F
8
9 from models.seq2seq import Seq2Seq
10
11 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
12
13
14 def train(vocabs, batch_gen, train_params, model_params):
15     word2int, int2word = vocabs
16     num_epoch = train_params['num_epoch']
17     learn_rate = train_params['learn_rate']
18     clip = train_params['clip']
19     eval_every = train_params['eval_every']
20     train_tf_ratio = train_params['train_tf_ratio']
21     val_tf_ratio = train_params['val_tf_ratio']
22
23     net = Seq2Seq(vocabs=vocabs, device=device, **model_params).to(device)
24     net.train()
25
26     opt = optim.Adam(net.parameters(), lr=learn_rate)
27     weights = calc_class_weights(batch_gen.data_dict, batch_gen.label_dict)
28     criterion = nn.CrossEntropyLoss(weight=weights, ignore_index=word2int['<pad>'])
29
30     print('Training is starting ...')
31     train_loss_list = []
32     val_loss_list = []
33     for epoch in range(num_epoch):
34         running_loss = 0
35
36         for idx, (x_cap, y_cap) in enumerate(batch_gen.generate('train')):
37             print('\rtrain: {}/{}'.format(idx, batch_gen.num_iter('train')), flush=True, end='')
38             x_cap, y_cap = x_cap.to(device), y_cap.to(device)
39
40             opt.zero_grad()
41             output = net(x_cap, y_cap, train_tf_ratio)
42
43             loss = criterion(output.view(-1, output.size(2)), y_cap.view(-1).long())
44             loss.backward()
45
46             nn.utils.clip_grad_norm_(net.parameters(), clip)
47             opt.step()
48
49             running_loss += loss.item()
50
51             if (idx+1) % eval_every == 0:
52                 print('\n')
53                 val_loss = evaluate(net, word2int, batch_gen, weights, val_tf_ratio)
54                 print("\nEpoch: {}/{}...".format(epoch + 1, num_epoch),
55                       "Step: {}...".format(idx),
56                       "Loss: {:.4f}...".format(running_loss / idx),
57                       "Val Loss: {:.4f}\n".format(val_loss))
58
59             print('\nCreating sample captions')
60             sample(net, vocabs, generator=batch_gen.generate('validation'))
61             print('\n')
62
63             train_loss_list.append(running_loss / idx)
64             val_loss_list.append(val_loss)
65
66             loss_file = open('results/losses.pkl', 'wb')
67             model_file = open('results/seq2seq.pkl', 'wb')
68             pickle.dump([train_loss_list, val_loss_list], loss_file)
69             pickle.dump(net, model_file)
70
71             print('Training finished, saving the model')
72             model_file = open('seq2seq.pkl', 'wb')
73             pickle.dump(net, model_file)
74
75

```