

Deep Learning based Hotel Review Classification

국민대학교 정보보호연구실
정성민

2018-06-11

Information Security Lab. @2018



Outline

- Introduction & Members
- Dataset
- Related Algorithm
- Model Architecture
- Experiment Result
- Q&A
- Appendix

Introduction

호텔 리뷰가 주어지면 평점을 예측해주는 딥러닝 분류 모델을 만들자.

Really lovely
hotel!
Friendly and
helpful!!

Sentiment Analysis



Sheets were dirty.
The breakfast
tasted terrible!

Members

- Seongmin Jeong (Only One)
 - 정보보호연구실 석사과정
- Member's(=My) role
 - 데이터 수집 및 가공
 - Deep Learning 학습 모델 설계
 - 모델 검증 실험
 - hyper parameter 튜닝

May 6, 2018 – Jun 11, 2018

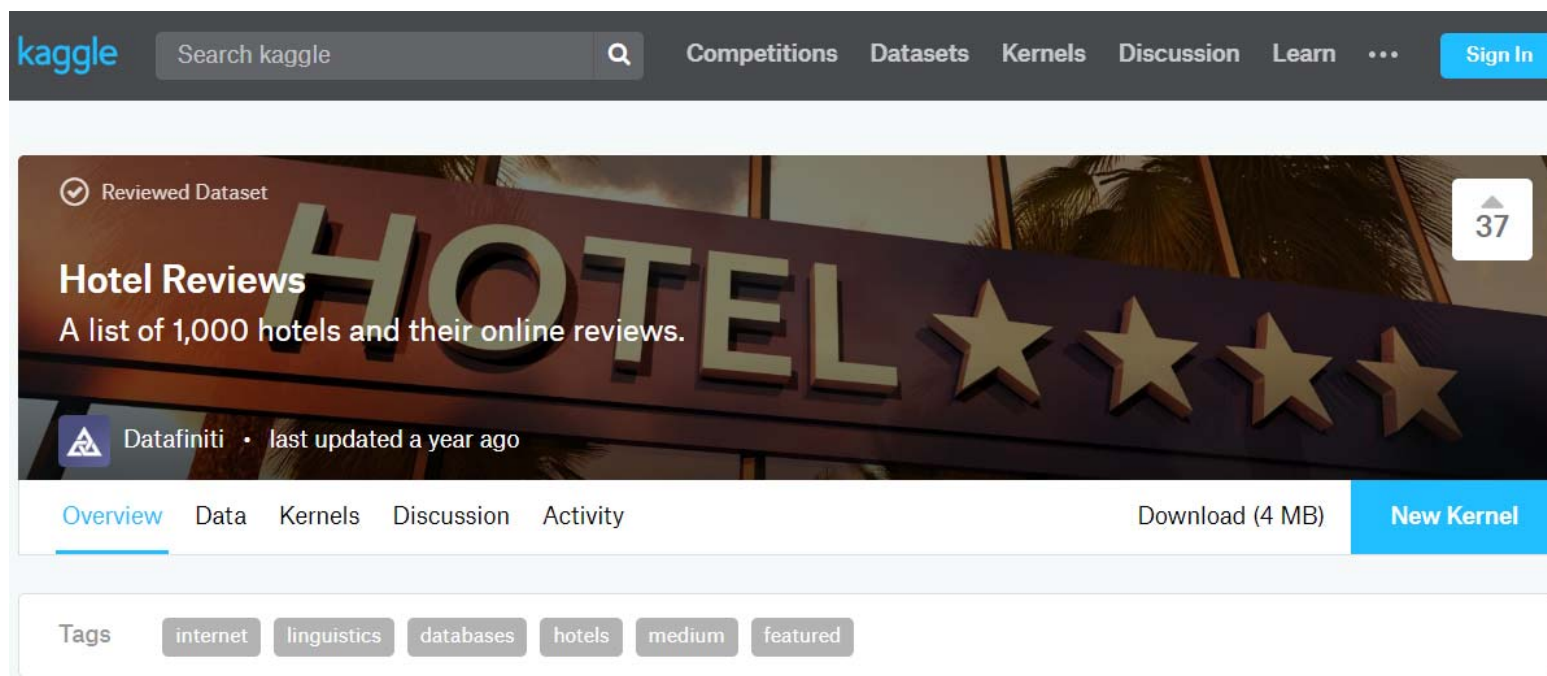
Contributions: Commits ▼

Contributions to master, excluding merge commits



Dataset

- Kaggle Datasets 에서 무료 제공 (Datafiniti 회사)
- 1,000 개의 Hotel에 대한 Review 데이터 (data: 약 36,000개, 19 columns)



출처: <https://www.kaggle.com/datafiniti/hotel-reviews>

Dataset

- Kaggle Datasets 에서 무료 제공 (Datafiniti 회사)
- 1,000 개의 Hotel에 대한 Review 데이터 (data: 약 36,000개, 19 columns)

Preview (first 100 rows) Column Metadata Column Metrics

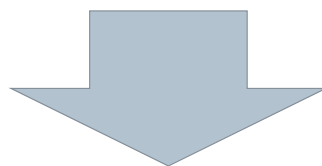
	reviews.rating	reviews.text	reviews.title
리뷰 평점, 리뷰 텍스트 두 정보만 사용	4	Pleasant 10 min walk along the sea front to the Water Bus. restaurants etc. Hotel was comfortable breakfast was good - quite a variety. Room aircon didn't work very well. Take mosquito repellant!	Good location away from the crowds
	5	Really lovely hotel. Stayed on the very top floor and were surprised by a Jacuzzi bath we didn't know we were getting! Staff were friendly and helpful and the included breakfast was great! Great location and great value for money. Didn't want to leave!	Great hotel with Jacuzzi bath!
	5	Ett mycket bra hotell. Det som drog ner betyget var att vi fick ett rum under taksarna dr det endast var full sthjd i 80 av rummets yta.	Lugnt 🍷🍷ge

출처: <https://www.kaggle.com/datafiniti/hotel-reviews>

Related Algorithm

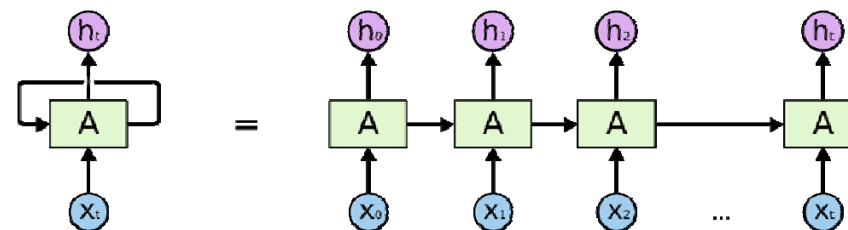
- Data: Text (length: dynamic)
- Type: English, Spanish

리뷰 데이터를 이해한 후 평점을 예측해야 한다.

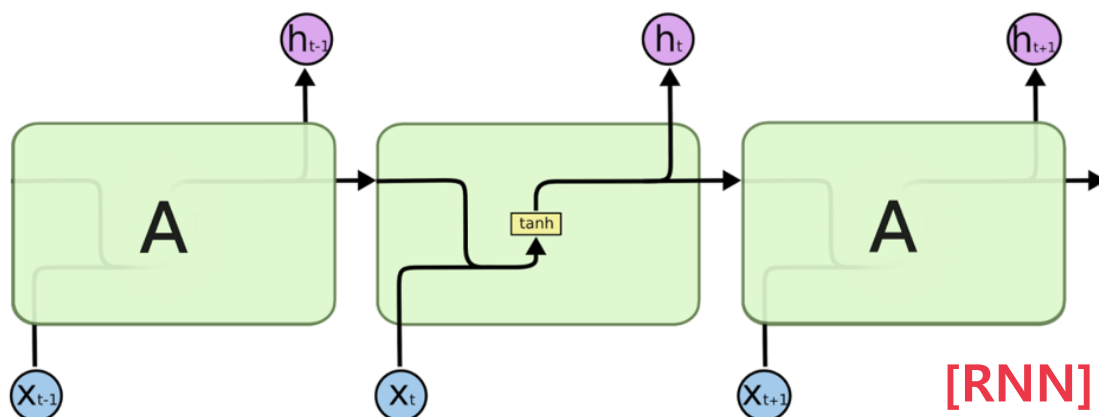


- Model: Many to one
- Type: Classification

LSTM 기반의 Many-to-one 학습 모델!



Related Algorithm



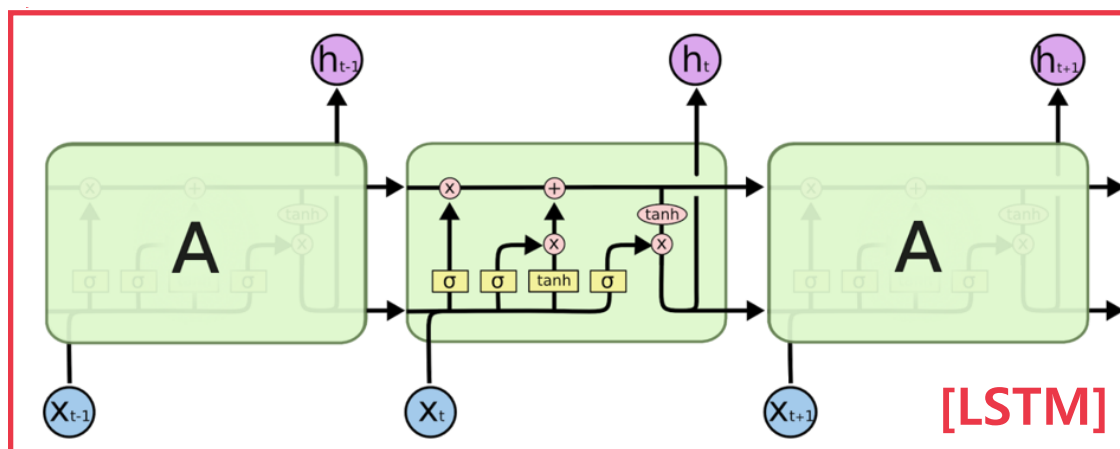
Vanishing Gradient Problem

출처: <https://ratsgo.github.io/natural%20language%20processing/2017/03/09/rnnlstm/>

2018-06-11

Information Security Lab. @2018

Related Algorithm



+ TensorFlow

출처: <https://ratsgo.github.io/natural%20language%20processing/2017/03/09/rnnlstm/>

2018-06-11

Information Security Lab. ©2018

Model Architecture

- LSTM Layer #
 - Single (1 Layer) → Stacked (e.g. 3 layer)
- Data Dimension
 - char-level : 44차원
 - word-level : 300차원
- Loss Function
 - Softmax cross entropy (with Adam Optimizer)

```
def inference_LSTM(x, x_seq_length, max_seq_length, data_dim, hidden_size, num_rnn_layers,
                  output_size, num_words, prob=1.0, train_flag=False):
    # x_one_hot = tf.one_hot(x, data_dim, name='x_one_hot') # (batch size, seq_length, data_dim)
    word_embedding = tf.get_variable('word_embedding', [num_words, data_dim])
    x_embedding_vector = tf.nn.embedding_lookup(word_embedding, x) # [input word #, embedding size]

    cell = rnn.BasicLSTMCell(num_units=hidden_size, state_is_tuple=True,
                             activation=tf.tanh)
    cells = [cell for _ in range(num_rnn_layers)]
    cells = rnn.MultiRNNCell(cells)
    # cell = rnn.DropoutWrapper(cell=cell, output_keep_prob=prob)

    outputs, state = tf.nn.dynamic_rnn(cells, x_embedding_vector,
                                       sequence_length=x_seq_length,
                                       dtype=tf.float32) # state: usually use final values

    # outputs : (?, ?, 300)
    # state : (3, ?, 300) ?

    y_reshape = tf.reshape(outputs, [-1, max_seq_length * hidden_size]) # 3278
    # y_reshape = tf.reshape(state, [-1, 6*hidden_size]) # why 6?

    y_fc = layers.fully_connected(inputs=y_reshape,
                                   num_outputs=output_size,
                                   activation_fn=None)

    return y_fc
```

LSTM inference 소스코드

출처: https://github.com/bonopi07/2018-1_advML_project/blob/master/src/network.py

Experiment Result (Accuracy)

- K-fold Cross Validation (K =5)
 - 모델 정확도 및 데이터셋 검증
- LSTM Layer #
 - Stacked (e.g. 3 layer)
- Data Dimension
 - word-level : 300차원 (word embedding vector)
- Loss Function
 - Softmax cross entropy (with Adam Optimizer)

실험 step	정확도 (accuracy)
1	89.17
2	88.42
3	91.03
4	84.98
5	86.11
평균	87.94

Contact us

Address:
7-613, Information Security Lab,
Kookmin University

Phone Number:
(010) 6667 - 5204

E-Mail
bonopi07@kookmin.ac.kr

Github
<https://github.com/bonopi07>



Any Questions?

Thank you for listening!

시연 장면

```

src [D:\working_board\2018-1_advML_project\src] - ...main.py [src] - PyCharm
File Edit View Navigate Code Refactor Run Tools VCS Window Help
src
├── main.py
├── model
│   ├── _processing.py
│   ├── config.ini
│   ├── data.py
│   ├── main.py
│   ├── model.py
│   └── network.py
└── External Libraries

main.py
43 accuracy_list = list()
44 data_idx = np.arange(len(total_data))
45 cv = KFold(n_splits=k_fold_value, shuffle=True, random_state=0)
46
47 for exp_idx, (train_idx, eval_idx) in enumerate(cv.split(data_idx)):
48     ...
49     separate data
50     ...
51     print('separate data {}'.format(exp_idx))
52     train_data = data.DataFeeder(total_data[train_idx], max_seq_len=max_seq_length,
53                                 batch_size=batch_size, mode='train')
54     eval_data = data.DataFeeder(total_data[eval_idx], max_seq_len=max_seq_length,
55                                batch_size=batch_size, mode='evaluate')
56
57     print('load model')
58     model_dic = {
59         'dropout_prob': float(config.get('CLASSIFIER', 'DROPOUT_PROB')),
60         'epoch': int(config.get('CLASSIFIER', 'EPOCH')),
61         'gpu_num': int(config.get('CLASSIFIER', 'GPU_NUM')),
62         'hidden_size': int(config.get('CLASSIFIER', 'HIDDEN_SIZE')),
63         'learning_rate': float(config.get('CLASSIFIER', 'LEARNING_RATE')),
64         'max_seq_length': max_seq_length,
65         'model_storage': config.get('CLASSIFIER', 'MODEL_STORAGE'),
66         'network': config.get('CLASSIFIER', 'NETWORK'),
67         'num_rnn_layers': int(config.get('CLASSIFIER', 'NUM_RNN_LAYERS')),
68         'num_words': number_of_words,
69         'output_size': int(config.get('CLASSIFIER', 'OUTPUT_SIZE'))
70     }
71
72     classifier = model.KISNet(model_num=step,
73                             train_data=train_data,
74                             eval_data=eval_data,
75                             model_dic=model_dic)
76
77     classifier.train()
78     accuracy = classifier.evaluate()
79     accuracy_list.append(accuracy)
80
81 run() for exp_idx, (train_idx, eval_idx) in enumerate(cv.split(data_idx)):
82     'dropout_prob'

```

```

Run main
[1850/6504] cost: 0.21 / acc: 0.94 / elapsed time: 1110.40
[1900/6504] cost: 0.12 / acc: 0.96 / elapsed time: 1143.10
[1950/6504] cost: 0.15 / acc: 0.95 / elapsed time: 1174.09
[2000/6504] cost: 0.22 / acc: 0.93 / elapsed time: 1202.64
[2050/6504] cost: 0.22 / acc: 0.92 / elapsed time: 1234.61
[2100/6504] cost: 0.15 / acc: 0.94 / elapsed time: 1265.06
[2150/6504] cost: 0.09 / acc: 0.96 / elapsed time: 1296.79

```

IDE and Plugin Updates: PyCharm is ready to update. (yesterday 오후 10:36) 59:14 CRLF UTF-8

```

main.py
58 accuracy_list = list()
59 data_idx = np.arange(len(total_data))
60 cv = KFold(n_splits=k_fold_value, shuffle=True, random_state=0)
61
62 for exp_idx, (train_idx, eval_idx) in enumerate(cv.split(data_idx)):
63     ...
64     separate data
65     ...
66     print('separate data {}'.format(exp_idx))
67     train_data = data.DataFeeder(total_data[train_idx], max_seq_len=max_seq_length,
68                                 batch_size=batch_size, mode='train')
69     eval_data = data.DataFeeder(total_data[eval_idx], max_seq_len=max_seq_length,
70                                batch_size=batch_size, mode='evaluate')
71
72     print('load model')
73     model_dic = {
74         'dropout_prob': float(config.get('CLASSIFIER', 'DROPOUT_PROB')),
75         'epoch': int(config.get('CLASSIFIER', 'EPOCH')),
76         'gpu_num': int(config.get('CLASSIFIER', 'GPU_NUM')),
77         'hidden_size': int(config.get('CLASSIFIER', 'HIDDEN_SIZE')),
78         'learning_rate': float(config.get('CLASSIFIER', 'LEARNING_RATE')),
79         'max_seq_length': max_seq_length,
80         'model_storage': config.get('CLASSIFIER', 'MODEL_STORAGE'),
81         'network': config.get('CLASSIFIER', 'NETWORK'),
82         'num_rnn_layers': int(config.get('CLASSIFIER', 'NUM_RNN_LAYERS')),
83         'num_words': number_of_words,
84         'output_size': int(config.get('CLASSIFIER', 'OUTPUT_SIZE'))
85     }
86
87     classifier = model.KISNet(model_num=step,
88                             train_data=train_data,
89                             eval_data=eval_data,
90                             model_dic=model_dic)
91
92     classifier.train()
93     accuracy = classifier.evaluate()
94     accuracy_list.append(accuracy)
95
96 run()

```

```

main main
[2000/6939] acc: 0.87 / elapsed time: 146.12
[3000/6939] acc: 0.91 / elapsed time: 205.91
[4000/6939] acc: 0.90 / elapsed time: 262.92
[5000/6939] acc: 0.88 / elapsed time: 320.30
[6000/6939] acc: 0.89 / elapsed time: 376.29
test time : 427.95388102531433
accuracy : 89.17380025940337
----evaluating finish----

```