# Deploying a Sentiment Analysis Model

**REVIEW**

**CODE REVIEW** 3

**HISTORY**

▶ website/index.html 1

▶ train/train.py 1

▼ serve/predict.py 1

```python
1   import argparse
2   import json
3   import os
4   import pickle
5   import sys
6   import sagemaker_containers
7   import pandas as pd
8   import numpy as np
9   import torch
10  import torch.nn as nn
11  import torch.optim as optim
12  import torch.utils.data
13
14  from model import LSTMClassifier
15
16  from utils import review_to_words, convert_and_pad
17
18  def model_fn(model_dir):
19      """Load the PyTorch model from the `model_dir` directory."""
20      print("Loading model.")
```

```python
22        # First, load the parameters used to create the model.
23        model_info = {}
24        model_info_path = os.path.join(model_dir, 'model_info.pth')
25        with open(model_info_path, 'rb') as f:
26            model_info = torch.load(f)
27
28        print("model_info: {}".format(model_info))
29
30        # Determine the device and construct the model.
31        device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
32        model = LSTMClassifier(model_info['embedding_dim'], model_info['hidden_dim'], mod
33
34        # Load the store model parameters.
35        model_path = os.path.join(model_dir, 'model.pth')
36        with open(model_path, 'rb') as f:
37            model.load_state_dict(torch.load(f))
38
39        # Load the saved word_dict.
40        word_dict_path = os.path.join(model_dir, 'word_dict.pkl')
41        with open(word_dict_path, 'rb') as f:
42            model.word_dict = pickle.load(f)
43
44        model.to(device).eval()
45
46        print("Done loading model.")
47        return model
48
49 def input_fn(serialized_input_data, content_type):
50        print('Deserializing the input data.')
51        if content_type == 'text/plain':
52            data = serialized_input_data.decode('utf-8')
53            return data
54        raise Exception('Requested unsupported ContentType in content_type: ' + content_t
55
56 def output_fn(prediction_output, accept):
57        print('Serializing the generated output.')
58        return str(prediction_output)
59
60 def predict_fn(input_data, model):
61        print('Inferring sentiment of input data.')
62
63        device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
64
65        if model.word_dict is None:
66            raise Exception('Model has not been loaded properly, no word_dict.')
67
68        # TODO: Process input_data so that it is ready to be sent to our model.
69        #       You should produce two variables:
70        #         data_X   - A sequence of length 500 which represents the converted revi
71        #         data_len - The length of the review
72
73        #data_X = None
74        #data_len = None
75
76        words = review_to_words(input_data)
77        data_X, data_len = convert_and_pad(model.word_dict, words)
78
79        # Using data_X and data_len we construct an appropriate input tensor. Remember
80        # that our model expects input data of the form 'len, review[500]'.
81        data_pack = np.hstack((data_len, data_X))
82        data_pack = data_pack.reshape(1, -1)
```

```
 83
 84      data = torch.from_numpy(data_pack)
 85      data = data.to(device)
 86
 87      # Make sure to put the model into evaluation mode
 88      model.eval()
 89
 90      # TODO: Compute the result of applying the model to the input data. The variable
 91      #       be a numpy array which contains a single integer which is either 1 or 0
 92
 93      with torch.no_grad():
 94          out = model.forward(data)
 95
 96      result = np.rint(out.numpy())
```

**AWESOME**

Expressions correct

```
 97
 98
 99      return result
100
```

RETURN TO PATH