

Sarah Ahmad (saz848)  
Vikram Kohli (vkp260)  
Junhao Li (jlo708)

## Checkpoint 5: Modeling with Neural Networks

**Classifier goal:** To predict whether the officers involved in a complaint will receive complaints in the future based on the text of summaries from past allegations.

To train a classifier to predict whether an officer will receive a complaint in the future based on the summary of an allegation, we created two sets of data for the classifier. The first contained the summary for each allegation with a summary as well as a label denoting whether the majority of active officers in the allegation had received a recent allegation, with a recent allegation being an allegation after 2014. The second dataset contained the combination of all allegations with summaries with the officers in the allegations and a label denoting whether the specific officer had received a recent allegation. Both datasets were constructed by exporting SQL queries on the CPDB to csv files.

After splitting this data into training and test datasets, we fed it into a modified version of a classifier neural network developed by OpenAI, and written in PyTorch. The OpenAI model uses as combination of transformers and unsupervised pre-training followed by discriminative fine-tuning on the particular task, and has been shown to achieve state of the art results on a variety of language tasks. Since the code for this model was written to be used for a multiple-choice task, we rewrote a substantial amount of the code in order to take our datasets as inputs, as well as to perform binary classification based on whether or not an officer will receive a complaint in the future given the summary of their previous complaint. (For reference, the unedited OpenAI classifier can be found at [https://github.com/huggingface/pytorch-openai-transformer-lm/blob/master/model\\_pytorch.py](https://github.com/huggingface/pytorch-openai-transformer-lm/blob/master/model_pytorch.py)); screenshots have been included at the bottom of this file that display the modifications made.

Following this modification, we split our data into a training set of size 1000 and a test set of size 80, and ran the network for 4 epochs. When testing in the case where allegations were repeated for each officer that had the same allegation placed against them, we found that our neural net had a maximum validation set accuracy of 78%, which is both relatively high and convincingly beat the ZeroR result of 70% after 4 epochs. However, in the case where allegation summaries were only considered once (and not for each officer who received the allegation), we were able to receive a test accuracy of only around 62%, which actually ended up being lower than the ZeroR result of 67% (unfortunately, this result was deleted, and so we're unable to

display the trace that produced it. Furthermore, even in the case where validation accuracy was high, the test accuracy for the second case only peaked at around 60.3% --- less than the performance of ZeroR, which achieves 72% accuracy. This is most likely not an issue with the model, and rather with the fact that our data set is rather small and doesn't include very many distinct summaries. Alternatively, this may point to the fact that the language used in the summaries of complaints against officers isn't a very good predictor of whether or not they are likely to have another allegation placed against them in the future. Determining which of these is the case would likely require a greater amount of summary data, as this would allow us to determine if the current model's failings are due to there just being insufficient information for training the neural net. Unfortunately, that vast majority of allegations listed do not have summaries and this seems consistent through the years of allegations so the amount of allegations without summaries does not appear to be changing.

Using the Pytorch OpenAI classifier as a basis, the words within the summaries were translated into vectors using their positions within sentences and pre-learned encodings. We had considered simpler features for encoding the summary such as mapping the words to integers or inputting an array of word frequencies but believed the encoded vectors to provide the best representation. However, experimentation with the form of the input features could certainly be explored further.

One issue we did appear to run into was that the training of the neural net was somewhat prohibitively slow; it took around 2 to 3 hours for the training to fully complete and for us to get our final results. This could be resolved by running on GPU, or by using a fewer number of examples / fine-tuning the code further.

# Results

## Queries for training/testing data

```
WITH data_repeatcomplaints as (  
  SELECT officer_id,  
         CASE WHEN SUM(CASE WHEN incident_date > '2013-12-31 23:59:59-06'  
                             THEN 1  
                             ELSE 0  
                           END) > 0  
         THEN 1  
         ELSE 0  
         END AS label  
  FROM data_allegation AS a  
  JOIN data_officer allegation AS oa ON a.id = oa.allegation_id  
  JOIN data_officer AS o ON oa.officer_id = o.id  
  WHERE active <> 'No'  
  GROUP BY officer_id  
  ORDER BY officer_id  
)  
SELECT allegation_id, oa.officer_id, summary, label  
FROM data_allegation AS a  
JOIN data_officer allegation AS oa ON a.id = oa.allegation_id  
JOIN data_repeatcomplaints AS rc ON rc.officer_id = oa.officer_id  
WHERE summary <> '';
```

```
WITH data_repeatcomplaints as (  
  SELECT officer_id,  
         CASE WHEN SUM(CASE WHEN incident_date > '2013-12-31 23:59:59-06'  
                             THEN 1  
                             ELSE 0  
                           END) > 0  
         THEN 1  
         ELSE 0  
         END AS recent_complaint  
  FROM data_allegation AS a  
  JOIN data_officer allegation AS oa ON a.id = oa.allegation_id  
  JOIN data_officer AS o ON oa.officer_id = o.id  
  WHERE active <> 'No'  
  GROUP BY officer_id  
  ORDER BY officer_id  
)  
SELECT allegation_id, summary,  
       CASE WHEN SUM(recent_complaint) > (COUNT(recent_complaint) / 2.0)  
       THEN 1  
       ELSE 0  
       END AS label  
FROM data_allegation AS a  
JOIN data_officer allegation AS oa ON a.id = oa.allegation_id  
JOIN data_repeatcomplaints AS rc ON rc.officer_id = oa.officer_id  
WHERE summary <> ''  
GROUP BY allegation_id, summary;
```

## Results from PyTorch neural net

Results after 4 epochs of training ---

*read as (epoch, num training examples, validation loss, test set loss, validation accuracy, test accuracy)*

```
(base) C:\Users\Vikram\Desktop\checkpoint_5>python train.py --dataset checkpoint5 --desc checkpoint5 --submit --analysis
Namespace(afn='gelu', analysis=True, attn_pdrop=0.1, b1=0.9, b2=0.999, bpe_path='model/vocab_40000.bpe', clf_pdrop=0.1,
data_dir='data/', dataset='checkpoint5', desc='checkpoint5', e=1e-08, embd_pdrop=0.1, encoder_path='model/encoder_bpe_40000.json', l2=0.01, lm_coef=0.5, log_dir='log/', lr=6.25e-05, lr_schedule='warmup_linear', lr_warmup=0.002, max_grad_norm=1, n_batch=8, n_ctx=512, n_embd=768, n_head=12, n_iter=5, n_layer=12, n_transfer=12, n_valid=199, opt='adam', resid_pdrop=0.1, save_dir='save/', seed=42, submission_dir='submission/', submit=True, vector_l2=False)
device cpu n_gpu 0
Encoding dataset...
C:\Users\Vikram\AppData\Local\Continuum\anaconda3\lib\site-packages\torch\nn\functional.py:52: UserWarning: size_average and reduce args will be deprecated, please use reduction='none' instead.
  warnings.warn(warning.format(ret))
Loading weights...
running epoch 0
Logging
1 100 22.284 27.479 67.84 60.30
running epoch 1
Logging
2 200 10.815 16.141 68.34 59.80
running epoch 2
Logging
3 300 23.680 37.422 68.34 59.80
running epoch 3
Logging
4 400 10.852 23.223 78.89 60.30
```

Screenshots of code for PyTorch neural net (our modifications can be found at

<https://github.com/JunhaoSLi/DataScienceF2018>):

[datasets.py](#) (written to pipe the CSV data into the network)

```

11
12 def _checkpoint5(path):
13     with open(path, encoding='utf-8') as f:
14         f = csv.reader(f)
15         statements = []
16         results = []
17         for i, line in enumerate(tqdm(list(f), ncols=80, leave=False)):
18             statements.append(line[1])
19             results.append(int(line[2]))
20     return statements, results
21
22 def checkpoint5(data_dir, n_train=800, n_valid=199):
23     statements, results = _checkpoint5(os.path.join(data_dir, 'training.csv'))
24     test_statements, _ = _checkpoint5(os.path.join(data_dir, 'testing.csv'))
25     tr_statements, va_statements, tr_results, va_results = train_test_split(statements, results, test_size=n_valid, random_state=seed)
26     train_statements = []
27     train_results = []
28     for statement, result in zip(tr_statements, tr_results):
29         train_statements.append(statement)
30         train_results.append(result)
31
32     val_statements = []
33     val_results = []
34     for statement, result in zip(va_statements, va_results):
35         val_statements.append(statement)
36         val_results.append(result)
37
38     train_results = np.asarray(train_results, dtype=np.int32)
39     val_results = np.asarray(val_results, dtype=np.int32)
40
41     return (train_statements, train_results), (val_statements, val_results), test_statements

```

train.py (creates the network, does training, displays results)

```

25 def transform_checkpoint5(statements):
26     n_batch = len(statements)
27     xmb = np.zeros((n_batch, n_ctx, 2), dtype=np.int32)
28     mmb = np.zeros((n_batch, n_ctx), dtype=np.float32)
29     start = encoder['_start_']
30     for i, statement in enumerate(statements):
31         if not isinstance(statement[0], int):
32             statement = sum(statement, [])
33         full = [start] + statement[:max_len] + [clf_token]
34         length = len(full)
35         xmb[i, :length, 0] = full
36         mmb[i, :length] = 1
37     xmb[:, :, 1] = np.arange(n_vocab + n_special, n_vocab + n_special + n_ctx)
38     return xmb, mmb
39
40
41 def iter_apply(Xs, Ms, Ys):
42     # fns = [lambda x: np.concatenate(x, 0), lambda x: float(np.sum(x))]
43     logits = []
44     cost = 0
45     with torch.no_grad():
46         dh_model.eval()
47         for xmb, mmb, ymb in iter_data(Xs, Ms, Ys, n_batch=n_batch_train, truncate=False, verbose=True):
48             n = len(xmb)
49             XMB = torch.tensor(xmb, dtype=torch.long).to(device)
50             YMB = torch.tensor(ymb, dtype=torch.long).to(device)
51             MMB = torch.tensor(mmb).to(device)
52             _, clf_logits = dh_model(XMB)
53             clf_logits *= n
54             clf_losses = compute_loss_fct(XMB, YMB, MMB, clf_logits, only_return_losses=True)
55             clf_losses *= n
56             logits.append(clf_logits.to("cpu").numpy())

```



```

62 def iter_predict(Xs, Ms):
63     logits = []
64     with torch.no_grad():
65         dh_model.eval()
66         for xmb, mmb in iter_data(Xs, Ms, n_batch=n_batch_train, truncate=False, verbose=True):
67             n = len(xmb)
68             XMB = torch.tensor(xmb, dtype=torch.long).to(device)
69             MMB = torch.tensor(mmb).to(device)
70             _, clf_logits = dh_model(XMB)
71             logits.append(clf_logits.to("cpu").numpy())
72     logits = np.concatenate(logits, 0)
73     return logits
74
75
76 def log(save_dir, desc):
77     global best_score
78     print("Logging")
79     tr_logits, tr_cost = iter_apply(trX[:n_valid], trM[:n_valid], trY[:n_valid])
80     va_logits, va_cost = iter_apply(vaX, vaM, vaY)
81     tr_cost = tr_cost / len(trY[:n_valid])
82     va_cost = va_cost / n_valid
83     tr_acc = accuracy_score(trY[:n_valid], np.argmax(tr_logits, 1)) * 100.
84     va_acc = accuracy_score(vaY, np.argmax(va_logits, 1)) * 100.
85     logger.log(n_epochs=n_epochs, n_updates=n_updates, tr_cost=tr_cost, va_cost=va_cost, tr_acc=tr_acc, va_acc=va_acc)
86     print('%d %d %.3f %.3f %.2f %.2f' % (n_epochs, n_updates, tr_cost, va_cost, tr_acc, va_acc))
87     if submit:
88         score = va_acc
89         if score > best_score:
90             best_score = score
91             path = os.path.join(save_dir, desc, 'best_params')
92             torch.save(dh_model.state_dict(), make_path(path))

```

```

94
95 def predict(dataset, submission_dir):
96     filename = filenames[dataset]
97     pred_fn = pred_fns[dataset]
98     label_decoder = label_decoders[dataset]
99     predictions = pred_fn(iter_predict(teX, teM))
100     if label_decoder is not None:
101         predictions = [label_decoder[prediction] for prediction in predictions]
102     path = os.path.join(submission_dir, filename)
103     os.makedirs(os.path.dirname(path), exist_ok=True)
104     with open(path, 'w') as f:
105         f.write('{}\t{}\n'.format('index', 'prediction'))
106         for i, prediction in enumerate(predictions):
107             f.write('{}\t{}\n'.format(i, prediction))
108
109
110 def run_epoch():
111     for xmb, mmb, ymb in iter_data(*shuffle(trX, trM, trYt, random_state=np.random),
112                                     n_batch=n_batch_train, truncate=True, verbose=True):
113         global n_updates
114         dh_model.train()
115         XMB = torch.tensor(xmb, dtype=torch.long).to(device)
116         YMB = torch.tensor(ymb, dtype=torch.long).to(device)
117         MMB = torch.tensor(mmb).to(device)
118         lm_logits, clf_logits = dh_model(XMB)
119         compute_loss_fct(XMB, YMB, MMB, clf_logits, lm_logits)
120         n_updates += 1
121         if n_updates in [1000, 2000, 4000, 8000, 16000, 32000] and n_epochs == 0:
122             log(save_dir, desc)
123
124
125 argmax = lambda x: np.argmax(x, 1)
126

```

```

138
139 if __name__ == '__main__':
140     parser = argparse.ArgumentParser()
141     parser.add_argument('--desc', type=str, help="Description")
142     parser.add_argument('--dataset', type=str)
143     parser.add_argument('--log_dir', type=str, default='log/')
144     parser.add_argument('--save_dir', type=str, default='save/')
145     parser.add_argument('--data_dir', type=str, default='data/')
146     parser.add_argument('--submission_dir', type=str, default='submission/')
147     parser.add_argument('--submit', action='store_true')
148     parser.add_argument('--analysis', action='store_true')
149     parser.add_argument('--seed', type=int, default=42)
150     parser.add_argument('--n_iter', type=int, default=5)
151     parser.add_argument('--n_batch', type=int, default=8)
152     parser.add_argument('--max_grad_norm', type=int, default=1)
153     parser.add_argument('--lr', type=float, default=6.25e-5)
154     parser.add_argument('--lr_warmup', type=float, default=0.002)
155     parser.add_argument('--n_ctx', type=int, default=512)
156     parser.add_argument('--n_embd', type=int, default=768)
157     parser.add_argument('--n_head', type=int, default=12)
158     parser.add_argument('--n_layer', type=int, default=12)
159     parser.add_argument('--embd_pdrop', type=float, default=0.1)
160     parser.add_argument('--attn_pdrop', type=float, default=0.1)
161     parser.add_argument('--resid_pdrop', type=float, default=0.1)
162     parser.add_argument('--clf_pdrop', type=float, default=0.1)
163     parser.add_argument('--l2', type=float, default=0.01)
164     parser.add_argument('--vector_l2', action='store_true')
165     parser.add_argument('--opt', type=str, default='adam')
166     parser.add_argument('--afn', type=str, default='gelu')
167     parser.add_argument('--lr_schedule', type=str, default='warmup_linear')
168     parser.add_argument('--encoder_path', type=str, default='model/encoder_bpe_40000.json')
169     parser.add_argument('--bpe_path', type=str, default='model/vocab_40000.bpe')

```

```

180 random.seed(args.seed)
181 np.random.seed(args.seed)
182 torch.manual_seed(args.seed)
183 torch.cuda.manual_seed_all(args.seed)
184
185 # Constants
186 submit = args.submit
187 dataset = args.dataset
188 n_ctx = args.n_ctx
189 save_dir = args.save_dir
190 desc = args.desc
191 data_dir = args.data_dir
192 log_dir = args.log_dir
193 submission_dir = args.submission_dir
194
195 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
196 n_gpu = torch.cuda.device_count()
197 print("device", device, "n_gpu", n_gpu)
198
199 logger = ResultLogger(path=os.path.join(log_dir, '{}.jsonl'.format(desc)), **args.__dict__)
200 text_encoder = TextEncoder(args.encoder_path, args.bpe_path)
201 encoder = text_encoder.encoder
202 n_vocab = len(text_encoder.encoder)
203
204 print("Encoding dataset...")
205 ((trX1, trY),
206  (vaX1, vaY),
207  teX1) = encode_dataset(*checkpoint5(data_dir, n_valid=args.n_valid),
208                          encoder=text_encoder)
209
210 #print(trX1)
211 encoder['start'] = len(encoder)

```

```

210     mp4_loader(loader)
211     encoder['_start_'] = len(encoder)
212     encoder['_delimiter_'] = len(encoder)
213     encoder['_classify_'] = len(encoder)
214     clf_token = encoder['_classify_']
215     n_special = 3
216     max_len = n_ctx // 2 - 2
217     n_ctx = min(max(
218         [len(x1[:max_len]) for x1 in trX1]
219         + [len(x1[:max_len]) for x1 in vaX1]
220         + [len(x1[:max_len]) for x1 in teX1]
221         ) + 2, n_ctx)
222
223     vocab = n_vocab + n_special + n_ctx
224     trX, trM = transform_checkpoint5(trX1)
225     vaX, vaM = transform_checkpoint5(vaX1)
226     if submit:
227         teX, teM = transform_checkpoint5(teX1)
228
229     n_train = len(trY)
230     n_valid = len(vaY)
231     n_batch_train = args.n_batch * max(n_gpu, 1)
232     n_updates_total = (n_train // n_batch_train) * args.n_iter
233
234     dh_model = DoubleHeadModel(args, clf_token, ('classification', 2), vocab, n_ctx)
235
236     criterion = nn.CrossEntropyLoss(reduce=False)
237     model_opt = OpenAIAdam(dh_model.parameters(),
238         lr=args.lr,
239         schedule=args.lr_schedule,
240         warmup=args.lr_warmup,
241         t_total=n_updates_total,

```



```

247         max_grad_norm=args.max_grad_norm)
248     compute_loss_fct = ClassificationLossCompute(criterion,
249                                                 criterion,
250                                                 args.lm_coef,
251                                                 model_opt)
252     load_openai_pretrained_model(dh_model.transformer, n_ctx=n_ctx, n_special=n_special)
253
254     dh_model.to(device)
255     dh_model = nn.DataParallel(dh_model)
256
257     n_updates = 0
258     n_epochs = 0
259     if dataset != 'stsb':
260         trYt = trY
261     if submit:
262         path = os.path.join(save_dir, desc, 'best_params')
263         torch.save(dh_model.state_dict(), make_path(path))
264     best_score = 0
265     for i in range(args.n_iter):
266         print("running epoch", i)
267         run_epoch()
268         n_epochs += 1
269         log(save_dir, desc)
270     if submit:
271         path = os.path.join(save_dir, desc, 'best_params')
272         dh_model.load_state_dict(torch.load(path))
273         predict(dataset, args.submission_dir)
274         if args.analysis:
275             checkpoints5_analysis(data_dir, os.path.join(args.submission_dir, 'Checkpoint5.tsv'),
276                                 os.path.join(log_dir, 'checkpoint5.jsonl'))
277

```

```

1  import argparse
2  import os
3  import random
4
5  import numpy as np
6  import torch
7  import torch.nn as nn
8  from sklearn.metrics import accuracy_score
9  from sklearn.utils import shuffle
10
11  from analysis import checkpoints5 as checkpoints5_analysis
12  from datasets import checkpoints5
13  from model_pytorch import DoubleHeadModel, load_openai_pretrained_model
14  from opt import OpenAIAdam
15  from text_utils import TextEncoder
16  from utils import (encode_dataset, iter_data,
17                    ResultLogger, make_path)
18  from loss import MultipleChoiceLossCompute, ClassificationLossCompute
19
20  ...
21  to run:
22  python train.py --dataset checkpoints5 --desc checkpoints5 --submit --analysis
23  ...
24
25  def transform_checkpoints5(statements):
26      n_batch = len(statements)
27      xmb = np.zeros((n_batch, n_ctx, 2), dtype=np.int32)
28      mmb = np.zeros((n_batch, n_ctx), dtype=np.float32)
29      start = encoder['_start_']
30      for i, statement in enumerate(statements):
31          if not isinstance(statement[0], int):
32              statement = sum(statement, [])
33          full = [start] + statement[:max_len] + [c1f token]

```