

Machine Learning: HW3 Report

November 19, 2016

Professor Hung-Yi Lee

電機三 B03901018 楊程皓

1. Supervised learning

參考 keras example 的 layers 以創建 model, code 如下:

```
1
2 data = []
3 val_data = []
4 answers = np.zeros((4000, 10), dtype=np.int)
5 val_answers = np.zeros((1000, 10), dtype=np.int)
6
7 # parsing data, answers
8 for ...
9
10 # transform to numpy array
11 data = np.array(data)
12 val_data = np.array(val_data)
13 data = data.astype('float32')
14 val_data = val_data.astype('float32')
15 data /= 255
16 val_data /= 255
17
18 model = Sequential()
19 model.add(Convolution2D(32, 3, 3, border_mode='same', input_shape=data.shape[1:]))
20 model.add(Activation('relu'))
21 model.add(Convolution2D(32, 3, 3))
22 model.add(Activation('relu'))
23 model.add(MaxPooling2D(pool_size=(2, 2)))
24 model.add(Dropout(0.25))
25
26 model.add(Convolution2D(64, 3, 3, border_mode='same'))
27 model.add(Activation('relu'))
28 model.add(Convolution2D(64, 3, 3))
29 model.add(Activation('relu'))
30 model.add(MaxPooling2D(pool_size=(2, 2)))
31 model.add(Dropout(0.25))
32
33 model.add(Flatten())
34 model.add(Dense(512))
35 model.add(Activation('relu'))
36 model.add(Dropout(0.5))
37 model.add(Dense(nb_classes))
38 model.add(Activation('softmax'))
39
40 # Use various optimizer, in this case, rmsprop
41 model.compile(loss='categorical_crossentropy', optimizer=optimizer, metrics=['←
```

```
accuracy'])
42
43 # train with 4/5 labeled data
44 model.fit(data, answers, batch_size=batch_size, nb_epoch=nb_epoch, shuffle=True, ←
    validation_data=(val_data, val_answers))
45
46 # save model to model_path
47 model.save(model_path)
```

經過各種 optimizer 測試後，此 model 以 rmsprop 為最佳，nb_epoch = 30，僅以 labeled data train 就可以達到 public data 0.63 (private 0.64) 的準確率。然而以這個 model 做 semi supervised training 卻會得更差的結果。

2. Semi-supervised learning(1)

將 supervised learning 後的 model 做 unlabeled data prediction 後的結果加入原本的 labeled data 作為進一步 training model 的資料，並以 predict 的屬性值線性作為 sample.weight 的參數，並加入 image preprocessing 來優化結果。code 如下：

```
1 data = []
2 training_data = []
3 answers = []
4
5 # parsing labeled data, unlabeled data, answers
6
7 # load supervised learning model from model_path
8 model = load_model(model_path)
9
10 # predict the classes of unlabeled data
11 results = model.predict(training_data, verbose=1)
12
13 # parsing new answers and weights
14
15 # preprocessing and data augmentation
16 datagen = ImageDataGenerator(
17     featurewise_center=False,
18     samplewise_center=False,
19     featurewise_std_normalization=False,
20     samplewise_std_normalization=False,
21     zca_whitening=False,
22     rotation_range=0,
23     width_shift_range=0.1,
24     height_shift_range=0.1,
25     horizontal_flip=True,
26     vertical_flip=False)
27
28 datagen.fit(data)
29
30 # train the model on the batches generated by datagen.flow()
31 model.fit_generator(datagen.flow(data, answers,
32     batch_size=batch_size),
33     samples_per_epoch=data.shape[0],
34     nb_epoch=nb_epoch)
35
36 # save semi-supervised model to ouput_model_path
37 model.save(output_model_path)
```

這裡的 supervised learning model 有用同樣的方式做 preprocessing，但因為單一 supervised learning 結果沒有上述的方式好，故沒有多做敘述。

這裡的 optimizer 選的是 adam, nb_epoch = 30, predict score threshold = 0.8 (不到 threshold 的 data 就不加入下一輪 training 的資料了), public data 準確率為 0.64 (private 0.65), 進步幅度其實並不大...

3. Semi-supervised learning(2)

以 autoencoder 做 semi-supervised learning, 只取前 5000 個最有可能的資料作為下一輪新的資料。code 如下:

```
1 model = Sequential()
2 # dim 3072 -> encoding_dim ( = 512)
3 model.add(Dense(encoding_dim, activation='relu', input_shape=(3072,)))
4 # add num_of_layers ( = 5) layers
5 for i in range(num_of_layers):
6     model.add(Dense(layer_dim, activation='relu'))
7 # dim layer_dim ( = 512) -> 3072
8 model.add(Dense(3072, activation='linear'))
9
10 # use optimizer rmsprop
11 model.compile(loss='mse', optimizer='rmsprop', metrics=['accuracy'])
12 model.fit(data, data,
13         batch_size=batch_size_auto,
14         nb_epoch=nb_epoch_auto,
15         verbose=1)
16     # validation_data=(x_test[0:3000], x_test[0:3000]))
17 score = model.evaluate(x_test[0:3000], x_test[0:3000])
18 print('score', score)
19
20 # define encoder as the middle layer's data
21 encoder = K.function([model.layers[0].input], [model.layers[2].output])
22
23 # encode labeled data
24 data_encoded = encoder([data])[0]
25
26
27 average_data = np.zeros((10, layer_dim), dtype=np.float)
28
29 # calculate average of each cluster of each 512 dim
30 for ...
31
32 # encode unlabeled data
33 x_test_encoded = encoder([x_test])[0]
34
35 # calculate min error of each unlabeled data of each cluster
36 for ...
37
38 # sort unlabeled data by error
39 results = sorted(results, key = getValue)
40
41 # pick first max_size ( = 5000) of results, linear weighted with rank
42 for i in range(max_size):
43     ...
44
45 # load supervised learning model from model_path
46 old_model = load_model(model_path)
47
48 # train the model with added data
```

```
49 old_model.fit(data, answers, batch_size=batch_size, nb_epoch=nb_epoch, shuffle=True, ↵  
    sample_weight=weights)  
50  
51 # save semi-supervised model to ouput_model_path  
52 old_model.save(output_model_path)
```

這裡的 supervised model 我取 adam optimizer 與 nb_epoch = 30, 新一輪的 optimizer 與 nb_epoch 取一樣, 但 private score 僅約為 0.46, 可能參數與 layers 還有待調整。

4. Compare & Analyze

可能因為對 deep learning 了解不夠深, 沒有設計出好的 model 來 train, 故無法作到如 supervised learning 就有 0.8 accuracy 的水準; 而 semi-supervised learning 感覺也沒有抓到核心的關鍵, 加入 weighted 與 threshold 後只比試出的最佳 supervised learning 高一些些 accuracy 而已。autoencoder 的確如老師所言不是很好做, accuracy 在 train 完後反而更糟糕了。

之後如要改進 accuracy 首要就是設計出更好的 model, 並用更好的 clustering 做 semi-supervised。