

▼ Data downloading

```
1 !wget --no-check-certificate --content-disposition https://raw.githubusercontent.com/Ju
2
3 !unzip language_data.zip
4 !rm language_data.zip
```

▼ import packages

```
1 !pip install unidecode

1 import os
2 import time
3 import math
4 import glob
5 import string
6 import random
7
8 import torch
9 import torch.nn as nn
10
11 from rnn.helpers import time_since
12
13 %matplotlib inline

1 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
```

▼ Language recognition with an RNN

If you've ever used an online translator you've probably seen a feature that automatically detects the input language. While this might be easy to do if you input unicode characters that are unique to one or a small group of languages (like "你好" or "γεια σας"), this problem is more challenging if the input only uses the available ASCII characters. In this case, something like "těší mě" would become "tesi me" in the ascii form. This is a more challenging problem in which the language must be recognized purely by the pattern of characters rather than unique unicode characters.

We will train an RNN to solve this problem for a small set of languages that can be converted to romanized ASCII form. For training data it would be ideal to have a large and varied dataset in different language styles. However, it is easy to find copies of the Bible which is a large text translated to different languages but in the same easily parsable format, so we will use 20 different copies of the Bible as training data. Using the same book for all of the different languages will hopefully prevent minor overfitting that might arise if we used different books for each language (fitting to common

✓ 47s completed at 10:51 AM



```
1 from unicode import unicode as unicodeToAscii
2
3 all_characters = string.printable
4 n_letters = len(all_characters)
5
6 print(unicodeToAscii('těší mě'))
```

tesí me

```
1 # Read a file and split into lines
2 def readFile(filename):
3     data = open(filename, encoding='utf-8').read().strip()
4     return unicodeToAscii(data)
5
6 def get_category_data(data_path):
7     # Build the category_data dictionary, a list of names per language
8     category_data = {}
9     all_categories = []
10    for filename in glob.glob(data_path):
11        category = os.path.splitext(os.path.basename(filename))[0].split('_')[0]
12        all_categories.append(category)
13        data = readFile(filename)
14        category_data[category] = data
15
16    return category_data, all_categories
```

The original text is split into two parts, train and test, so that we can make sure that the model is not simply memorizing the train data.

```
1 train_data_path = 'language_data/train/*_train.txt'
2 test_data_path = 'language_data/test/*_test.txt'
3
4 train_category_data, all_categories = get_category_data(train_data_path)
5 test_category_data, test_all_categories = get_category_data(test_data_path)
6
7 n_languages = len(all_categories)
8
9 print(len(all_categories))
10 print(all_categories)
11
12
13
14
15
16
17
18
19
20
21 ['lithuanian', 'xhosa', 'english', 'romanian', 'portuguese', 'norwegian', 'esperanto']
```

Data processing

```
1 def categoryFromOutput(output):
2     top_n, top_i = output.topk(1, dim=1)
```

```

6 # turn string into long tensor
7 def stringToTensor(string):
8     tensor = torch.zeros(len(string), requires_grad=True).long()
9     for c in range(len(string)):
10         tensor[c] = all_characters.index(string[c])
11     return tensor
12
13 def load_random_batch(text, chunk_len, batch_size):
14     input_data = torch.zeros(batch_size, chunk_len).long().to(device)
15     target = torch.zeros(batch_size, 1).long().to(device)
16     input_text = []
17     for i in range(batch_size):
18         category = all_categories[random.randint(0, len(all_categories) - 1)]
19         line_start = random.randint(0, len(text[category]) - chunk_len)
20         category_tensor = torch.tensor([all_categories.index(category)], dtype=torch.long)
21         line = text[category][line_start:line_start+chunk_len]
22         input_text.append(line)
23         input_data[i] = stringToTensor(line)
24         target[i] = category_tensor
25     return input_data, target, input_text

```

Implement Model

For this classification task, we can use the same model we implement for the generation task which is located in `rnn/model.py`. See the `MP4_generation.ipynb` notebook for more instructions. In this case each output vector of our RNN will have the dimension of the number of possible languages (i.e. `n_languages`). We will use this vector to predict a distribution over the languages.

In the generation task, we used the output of the RNN at every time step to predict the next letter and our loss included the output from each of these predictions. However, in this task we use the output of the RNN at the end of the sequence to predict the language, so our loss function will use only the predicted output from the last time step.

Train RNN

```
1 from rnn.model import RNN
```

```

1 chunk_len = 100
2
3 BATCH_SIZE = 100
4 n_epochs = 2000
5 hidden_size = 120
6 n_layers = 4
7 learning_rate = 0.003

```

TODO: Fill in the train function. You should initialize a hidden layer representation using your RNN's `init_hidden` function, set the model gradients to zero, and loop over each time step (character) in the input tensor. For each time step compute the output of the of the RNN and the next hidden layer representation. The cross entropy loss should be computed over the last RNN output scores from the end of the sequence and the target classification tensor. Lastly, call backward on the loss and take an optimizer step.

```

1 def train(rnn, target_tensor, data_tensor, optimizer, criterion, batch_size=BATCH_SIZE)
2     """
3     Inputs:
4     - rnn: model
5     - target_target: target character data tensor of shape (batch_size, 1)
6     - data_tensor: input character data tensor of shape (batch_size, chunk_len)
7     - optimizer: rnn model optimizer
8     - criterion: loss function
9     - batch_size: data batch size
10    Returns:
11    - output: output from RNN from end of sequence
12    - loss: computed loss value as python float
13    """
14    # output, loss = None, None
15    #####
16    #         YOUR CODE HERE         #
17    #####
18    # with torch.no_grad():
19    data_tensor = data_tensor.to(device)
20    hidden = rnn.init_hidden(batch_size, device=device)
21    rnn.zero_grad()
22    for c in range(chunk_len):
23        # output, hidden = rnn.forward(data_tensor, hidden)
24        output, hidden = rnn(data_tensor[:,c], hidden)
25
26    loss = criterion(output, target_tensor.squeeze())
27    loss.backward()
28    optimizer.step()
29    #####          END          #####
30
31    return output, loss.item()
32

```

```

1 def evaluate(rnn, data_tensor, seq_len=chunk_len, batch_size=BATCH_SIZE):
2     with torch.no_grad():
3         data_tensor = data_tensor.to(device)
4         hidden = rnn.init_hidden(batch_size, device=device)
5         for i in range(seq_len):
6             output, hidden = rnn(data_tensor[:,i], hidden)

```

```
14         return output, loss.item()

1  n_iters = 10000 #2000 #100000
2  print_every = 50
3  plot_every = 50
4
5
6  # Keep track of losses for plotting
7  current_loss = 0
8  current_test_loss = 0
9  all_losses = []
10 all_test_losses = []
11
12 start = time.time()
13
14 optimizer = torch.optim.Adam(rnn.parameters(), lr=learning_rate)
15
16
17 number_correct = 0
18 for iter in range(1, n_iters + 1):
19     input_data, target_category, text_data = load_random_batch(train_category_data, c
20     output, loss = train(rnn, target_category, input_data, optimizer, criterion)
21     current_loss += loss
22
23     _, test_loss = eval_test(rnn, target_category, input_data)
24     current_test_loss += test_loss
25
26     guess_i = categoryFromOutput(output)
27     number_correct += (target_category.squeeze()==guess_i.squeeze()).long().sum()
28
29     # Print iter number, loss, name and guess
30     if iter % print_every == 0:
31         sample_idx = 0
32         guess = all_categories[guess_i[sample_idx]]
33
34         category = all_categories[int(target_category[sample_idx])]
35
36         correct = '✓' if guess == category else 'x (%s)' % category
37         print('%d %d%% (%s) %.4f %.4f %s / %s %s' % (iter, iter / n_iters * 100, time
38         print('Train accuracy: {}'.format(float(number_correct)/float(print_every*BAT
39         number_correct = 0
40
41     # Add current loss avg to list of losses
42     if iter % plot_every == 0:
43         all_losses.append(current_loss / plot_every)
44         current_loss = 0
45         all_test_losses.append(current_test_loss / plot_every)
```

250 2% (0m 30s) 0.0707 0.0710 i cieli. Quando la udire la sua voce, v e un rumor u al
Train accuracy: 0.9134

300 3% (0m 43s) 0.1365 0.1173 lkeskare ham. Og se, to blinde sade ved Vejen, og da de
Train accuracy: 0.9244

350 3% (0m 50s) 0.1764 0.1601 a no coracao para fazer por Jerusalem. Nao havia comigo
Train accuracy: 0.938

400 4% (0m 57s) 0.0550 0.0424 t op rad som blir til skam for ditt hus, lagt op rad or
Train accuracy: 0.938

450 4% (1m 5s) 0.1527 0.1085 ju edhe per pak kohe; ecni gjersa keni drite, qe te mos
Train accuracy: 0.9342

500 5% (1m 12s) 0.1983 0.2706 luyordu. Umma, Afek ve Rehov; koyleriyle birlikte yirmi
Train accuracy: 0.9432

550 5% (1m 19s) 0.0934 0.0697 udesta. Missa elattekkin, ette saa syoda lintujen etteka
Train accuracy: 0.9304

600 6% (1m 26s) 0.1016 0.1012 igliuoli d'Aaronne, sara unto per succedergli, fara anc
Train accuracy: 0.9436

650 6% (1m 33s) 0.1681 0.1083 forte dig ud af AEgypten, af Traellehuset. Du ma ikke
Train accuracy: 0.9362

700 7% (1m 41s) 0.1030 0.0707 c Gie-ho-va cam dong Gia-ha-xi-en, con trai Xa-cha-ri,
Train accuracy: 0.9498

750 7% (1m 48s) 0.0889 0.0838 murte igjen for mig, sa jeg ikke kan komme ut; han gjo
Train accuracy: 0.958

800 8% (1m 55s) 0.0741 0.0665 '. Dhe ai vepro i keshtu dhe dora e tij iu shendosh si
Train accuracy: 0.9624

850 8% (2m 2s) 0.0855 0.0613 s vel klause: "Ka jai padaryti?" Gehazis atsiliepe: "Ji
Train accuracy: 0.9602

900 9% (2m 9s) 0.0778 0.0592 adu tautai ar karalystei ja statyti ir itvirtinti, bet j
Train accuracy: 0.9718

950 9% (2m 17s) 0.2178 0.1646 dan, viidenkymmenen ja kymmenen paallikoiksi. He ratkoi
Train accuracy: 0.9658

1000 10% (2m 24s) 0.0971 0.0919 pore m, aquele, o Espirito da verdade, ele vos guiara
Train accuracy: 0.9528

1050 10% (2m 31s) 0.0731 0.0565 ezzel borita be. Es a rudakat betola az oltar oldalai
Train accuracy: 0.9564

1100 11% (2m 38s) 0.2225 0.1830 es brought Aaron's sons, and clothed them with coats,
Train accuracy: 0.9606

1150 11% (2m 45s) 0.1095 0.0788 vim, ze me uz neuvidi nikdo z vas, k nimz jsem na svy
Train accuracy: 0.9678

1200 12% (2m 52s) 0.1361 0.1039 apenbaret dem det. For hans usynlige vesen, bade hans
Train accuracy: 0.959

1250 12% (3m 0s) 0.0771 0.0758 a unjengaye. Akabuyanga noko uYehova ekuvutheni komsin
Train accuracy: 0.963

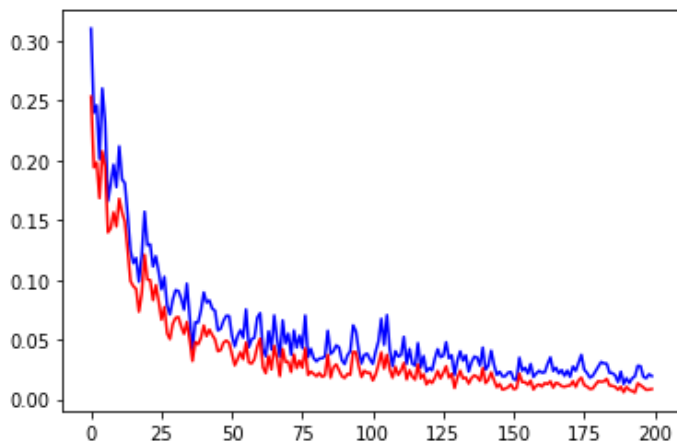
1300 13% (3m 7s) 0.0911 0.0694 o, beza kwaYuda naseYerusalem; ngokuba uYarobheham noo
Train accuracy: 0.9714

1350 13% (3m 14s) 0.0945 0.0791 sidererete come impuri. Pero, una fonte o una cistern
Train accuracy: 0.9668

1400 14% (3m 21s) 0.1444 0.1039 apatat astfel fiinta`, -zice Domnul. -,Iata spre cine
Train accuracy: 0.9722

1450 14% (3m 28s) 0.0105 0.0097 uizo. Porque para todo proposito ha tempo e juizo; po
Train accuracy: 0.9784

[<matplotlib.lines.Line2D at 0x7f8fb35cde50>]



Evaluate results

We now visualize the performance of our model by creating a confusion matrix. The ground truth languages of samples are represented by rows in the matrix while the predicted languages are represented by columns.

In this evaluation we consider sequences of variable sizes rather than the fixed length sequences we used for training.

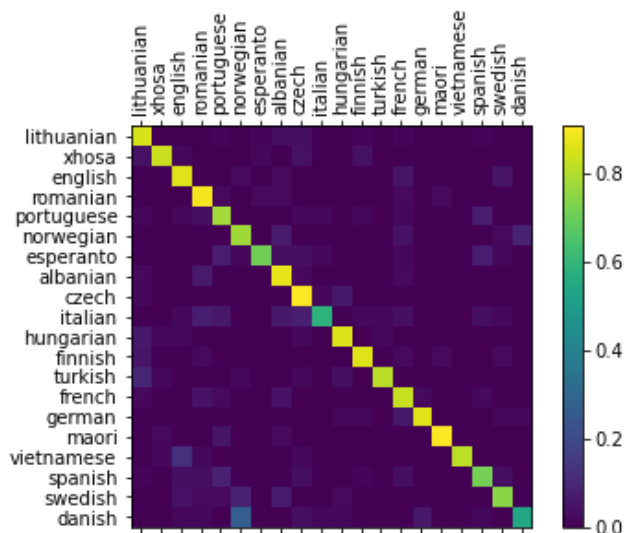
```
1 eval_batch_size = 1 # needs to be set to 1 for evaluating different sequence lengths
2
3 # Keep track of correct guesses in a confusion matrix
4 confusion = torch.zeros(n_languages, n_languages)
5 n_confusion = 1000
6 num_correct = 0
7 total = 0
8
9 for i in range(n_confusion):
10     eval_chunk_len = random.randint(10, 50) # in evaluation we will look at sequences o
11     input_data, target_category, text_data = load_random_batch(test_category_data, chun
12     output = evaluate(rnn, input_data, seq_len=eval_chunk_len, batch_size=eval_batch_si
13
14     guess_i = categoryFromOutput(output)
15     category_i = [int(target_category[idx]) for idx in range(len(target_category))]
16     for j in range(eval_batch_size):
```

```

30 ax = fig.add_subplot(111)
31 cax = ax.matshow(confusion.numpy())
32 fig.colorbar(cax)
33
34 # Set up axes
35 ax.set_xticklabels([''] + all_categories, rotation=90)
36 ax.set_yticklabels([''] + all_categories)
37
38 # Force label at every tick
39 ax.xaxis.set_major_locator(ticker.MultipleLocator(1))
40 ax.yaxis.set_major_locator(ticker.MultipleLocator(1))
41
42 plt.show()

```

Test accuracy: 0.797



You can pick out bright spots off the main axis that show which languages it guesses incorrectly.

Run on User Input

Now you can test your model on your own input.

```

1 def predict(input_line, n_predictions=5):

```



```
17     predictions.append([value, all_categories[category_index]])
18
19 predict('This is a phrase to test the model on user input')
20
> This is a phrase to test the model on user input
(8.24) albanian
(3.16) english
(2.05) french
(-0.90) romanian
(-0.93) hungarian
```

Output Kaggle submission file

Once you have found a good set of hyperparameters submit the output of your model on the Kaggle test file.

```
1  ### DO NOT CHANGE KAGGLE SUBMISSION CODE ###
2  import csv
3
4  kaggle_test_file_path = 'language_data/kaggle_rnn_language_classification_test.txt'
5  with open(kaggle_test_file_path, 'r') as f:
6      lines = f.readlines()
7
8  output_rows = []
9  for i, line in enumerate(lines):
10     sample = line.rstrip()
11     sample_chunk_len = len(sample)
12     input_data = stringToTensor(sample).unsqueeze(0)
13     output = evaluate(rnn, input_data, seq_len=sample_chunk_len, batch_size=1)
14     guess_i = categoryFromOutput(output)
15     output_rows.append((str(i+1), all_categories[guess_i]))
16
17  submission_file_path = 'kaggle_rnn_submission.txt'
18  with open(submission_file_path, 'w') as f:
19     output_rows = [('id', 'category')] + output_rows
```

