

Module 44 : Sequence models in Pytorch

- outline :
- i) introduce dataset
 - ii) data processing : train-test split, encoding, visualization
 - iii) basic RNN - testing inference
 - iv) evaluation and training
 - v) LSTM
 - vi) GRU
 - vii) Exercises
- ↑
preparing the data.
↓
to check if everything
is working

*) importing libraries : from io import open
import os, string, random, time, math
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

from sklearn.model_selection import train_test_split

import torch
import torch.nn as nn
import torch.optim as optim

```
device_gpu = torch.device("cuda:0" if torch.cuda_is_available()  
from IPython.display import clear_output else "cpu")  
                                                                  ↓  
                                                                  used to clear o/p of cell  
                                                                  ↓  
                                                                  for ex. needing to show multiple plots,  
                                                                  but retaining only the last instance,  
                                                                  then this construct can be used to  
                                                                  clear the intermediate o/p's programmatically
```

- *) Task : Given a name o/p the language of the name typically .
 - ↓
 - name : varying length string
 - multi classification task. → softmax based classifier
 - ↓
 - sequence to class task. (sequence labelling task)

*) Dataset :

```
languages = []
```

```
data = []
```

```
X = []
```

```
Y = []
```

with open ('name2lang.txt', 'r') as f :

for line in f :

line = line.split(',')

name = line[0].strip()

lang = line[1].strip()

if not lang in languages :

languages.append(lang)

X.append(name)

Y.append(lang)

data.append((name, lang))

} storing it in diff
formats.

n-languages = len(languages)

print(languages)

print(data[0:10])

→ train-test split

X-train, X-test, Y-train, Y-test =

train-test-split(X, Y, test_size = 0.2
, random_state = 0
, stratify = Y)

* encoding names and languages: languages basically represent classes, so we encode them using one hot encoding, whereas the input strings will require an encoding of each character, and then stacking of these vector encodings as one entity. Characters are encoded since they are the one sent to the model one after the other.

all_letters = string.ascii_letters + ".,;"

n_letters = len(all_letters)

↓

used to dimension our sequence models

to accept inputs of upto n letters

↓

one hot encoding for each character

→ encoding fn for string

def name_rep(name):

rep = torch.zeros(len(name), 1, n_letters)

for index, letter in enumerate(name):

pos = all_letters.find(letter)

rep[index][pos] = 1

return rep

```
def lang_rep(lang):  
    return torch.tensor([languages.index(lang)],  
                       dtype=torch.long)  
↓  
index label corresponding to the  
language.
```

→ Basic visualization

```
count = {}  
for l in languages:  
    count[l] = 0  
for d in data:  
    count[d[1]] += 1
```

```
plt_ = sns.barplot(list(count.keys()),  
                   list(count.values))
```

```
plt_.set_xticklabels(plt_.get_xticklabels(),  
                     rotation=90) → rotate the labels  
plt.show() to avoid visual  
overlap.
```

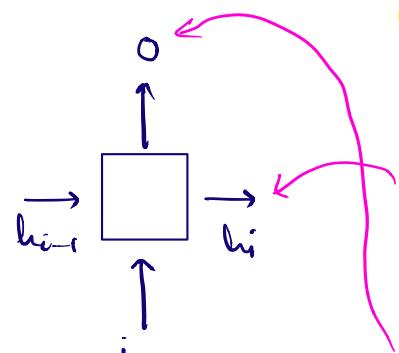
figuring out the baseline accuracy (accuracy on random class selection) for a dataset.

7) Basic net and testing inference

```
class RNN_net (nn.Module):
```

corresponding to size of 1-hot encoded vector

```
def __init__(self, input_size, hidden_size,  
           output_size):  
    self.W1 = np.random.rand(input_size, hidden_size)  
    self.b1 = np.zeros((1, hidden_size))  
    self.W2 = np.random.rand(hidden_size, output_size)  
    self.b2 = np.zeros((1, output_size))
```



super(RNN_net, self). __init__() → parent module's init
self.hidden_size = hidden_size
self.i2h = nn.Linear(input_size + hidden_size,
hidden_size)
self.i2o = nn.Linear(input_size + hidden_size,
output_size)

$$\text{ref. softmax} = \text{nn.LogSoftmax} \quad (\text{dim}=1)$$

i/p vector
↑

→ initial hidden state

```
def forward (self, input_, hidden):
```

```
combined = torch.cat((input_, hidden), 1)
```

hidden = self.size (combined)

output = self.x2o(combined)

output = self.softmax(output)

return output, hidden

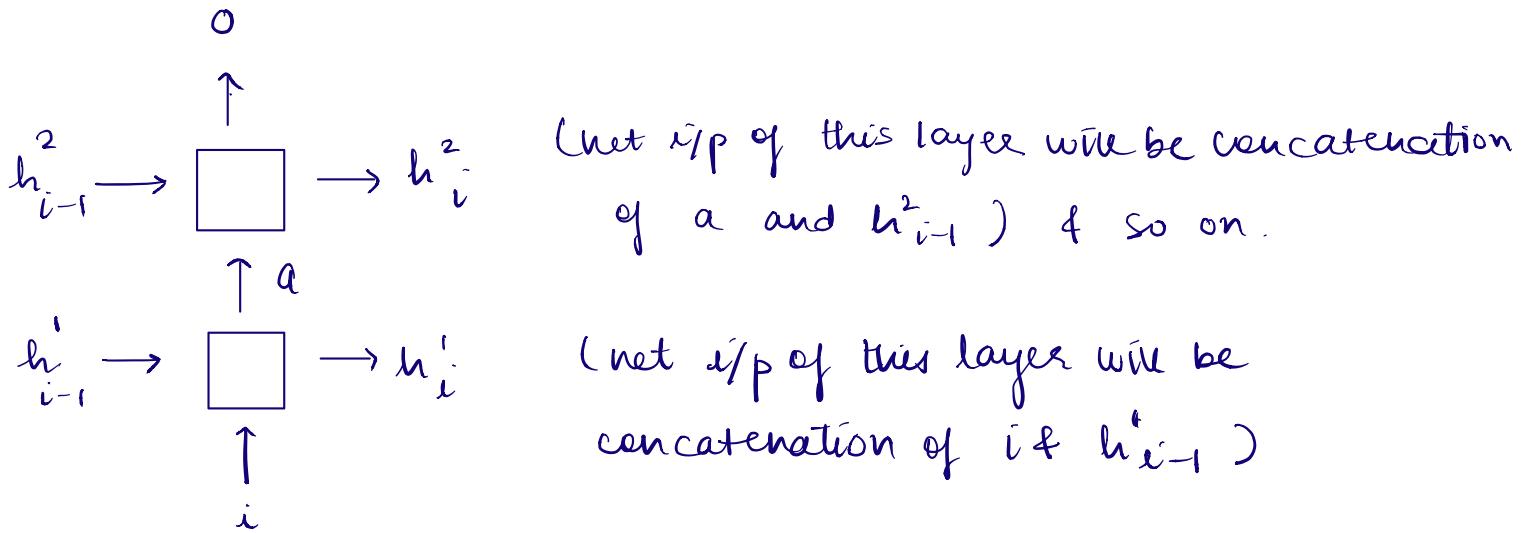
→ can be avoided
for inference,
softmax is exp.

can be done by defining
own boolean.

```
def init_hidden(self):
```

```
    return (torch.zeros(1, self.hidden_size))
```

→ Extending the idea to 2 layers.



→ hidden_size is a hyperparameter, and it is to be tuned based on the complexity of the task and the amount of data at hand.

→ greater the value, more is the capacity for complexity with the chance of overfitting
 $n\text{-hidden} = 128$

```
net = RNN(net(n_letters, n_hidden, n_languages))
```

def infer(net, name):
 → a little nontrivial as compared to DNN,
 or CNN.

```
    net.eval()
```

```
    name_onehot = name_reps(name) → one hot encoding repr.  
    of the i/p -
```

hidden = net . init - hidden () → initializing hidden unit
→ one character at a time

for i in range (name - ohe . size () [0]) :
 output , hidden = net (name - ohe [i] , hidden)

return output → all opps but the last one ignored,
where the last one returns the
probability distribution corresponding to
the labels of the task.

*) name - ohe [i] represents a vector corresponding to (1 , n - letters) ,
The reason we had the middle dimension , is because we
always want to give an input to the RNN which is (1 , x) .

output = infer (net , " Adam ") → inference run
index = torch . argmax (output) → fetching class with max
probability
print (output , index)

*) evaluate model

→ to remove bias due to correlation between samples being
in a sequence .

def dataloader (n_points , X , Y) :

 to - net = [] → no . of samples to sample from dataset
 for i in range (n - points) :

```
index_ = np.random.randint(len(x_))
name, lang = x_[index_], y_[index_]
to_net.append((name, lang, name_rep(name),
                lang_rep(lang)))
return to_ret
```

one hot encodings

```
dataloader(2, x-train, y-train)
```

```
def eval(net, n-points, k, x_, y_):
    data_ = dataloader(n-points, x_, y_)
    correct = 0
```

for name, language, name-one, lang-rep in data_:
 output = infer(net, name)
 val, indices = output.top(k)

if lang-repr in indices:
 correct += 1

accuracy = correct / n-points
return accuracy

```
eval(net, 1000, 3, x-test, y-test)
```

*) Training

computation for one batch.
↑

```
def train ( net , opt , criterion , n-points ) :  
    opt . zero-grad ()  
    total-loss = 0  
  
    data = dataloader ( n-points , X-train , Y-train )  
  
    for name , lang , name-ohe , lang-rep in data :  
        hidden = net . init-hidden ()  
        for i in range ( name-ohe . size [0] ) :  
            output , hidden = net ( name-ohe [i] ,  
            hidden )  
  
            loss = criterion ( output , lang-rep )  
            loss . backward ( retain-graph = True ) → backpropagation through the  
            dynamically created  
            compute graph  
            total-loss += loss  
  
    opt . step ()  
  
    return total-loss / n-points
```

```
criterion = nn.NLLLoss()
```

```
opt = optim.SGD(net.parameters(), lr=0.01,  
momentum=0.9)
```

1. / time

```
train(net, opt, criterion, 200)
```

```
eval(net, 1000, l, λ-test, γ-test)
```

good practice to have
separate infer, train &
↑ training setup

*) full training setup : multiple batches , observing how loss is
changing.

```
def train_setup(net, lr=0.01, n_batches=100,  
batch_size=10, momentum=0.9,  
display_freq=5):
```

```
criterion = nn.NLLLoss()
```

→ can be cross entropy if softmax
not applied on o/p. → beneficial
for inference
separation
for
speedup.

```
loss_arr = np.zeros(n_batches + 1)
```

```
for i in range(n_batches):
```

→ this workflow can
also contain checkpointing
and saving of
model.

```
loss_arr[i+1] = (loss_arr[i]+i
```

```
+ train(net, opt, criterion, batch_size)) / (i+1)
```

```

if i > display_freq == display_freq - 1:
    clear_output(wait = True)
    print('Iteration', i, 'Top-1:', eval(net, len(x_test), 1, x_test, y_test),
          'Top-2:', eval(net, len(x_test), 2, x_test, y_test))

plt.figure()
plt.plot(loss_arr[1:i], '-+')
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.show()
print('\n\n')

```

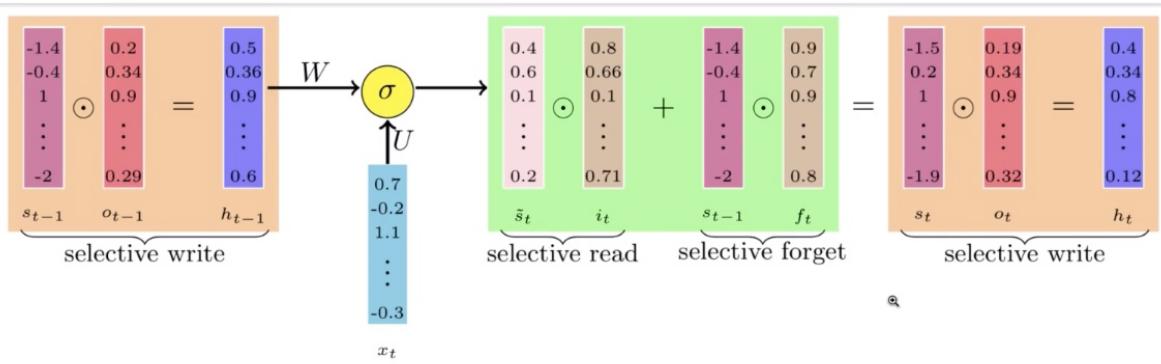
→ n-hidden = 128

```

net = RNN_net(n_letters, n_hidden, n_languages)
train_setup(net, lr=0.0005, n_batches=100,
           batch_size=256)

```

* LSTM



Applies a multi-layer long short-term memory (LSTM) RNN to an input sequence.

For each element in the input sequence, each layer computes the following function:

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{(t-1)} + b_{hi}) \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{(t-1)} + b_{hf}) \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{(t-1)} + b_{hg}) \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{(t-1)} + b_{ho}) \\ c_t &= f_t * c_{(t-1)} + i_t * g_t \\ h_t &= o_t * \tanh(c_t) \end{aligned}$$

→ LSTM block
fn. handles major computations.

where h_t is the hidden state at time t , c_t is the cell state at time t , x_t is the input at time t , $h_{(t-1)}$ is the hidden state of the layer at time $t-1$ or the initial hidden state at time 0, and i_t, f_t, g_t, o_t are the input, forget, cell, and output gates, respectively. σ is the sigmoid function, and $*$ is the Hadamard product.

In a multilayer LSTM, the input $x_t^{(l)}$ of the l -th layer ($l \geq 2$) is the hidden state $h_t^{(l-1)}$ of the previous layer multiplied by dropout $\delta_t^{(l-1)}$ where each $\delta_t^{(l-1)}$ is a Bernoulli random variable which is 0 with probability `dropout`.

in terms of func signature, tried to keep things same to as they were for RNN.

```
def __init__(self, input_size, hidden_size, output_size):
    super(LSTM_net, self).__init__()
    self.hidden_size = hidden_size
    self.lstm_cell = nn.LSTM(input_size,
                           hidden_size)
    self.h2o = nn.Linear(hidden_size,
                        output_size)
```

self.softmax = nn.LogSoftmax(dim=2)

↑ take benefit of vectorization

multiple i/p's at same time

def forward(self, input, hidden): to get performance benefit
 out, hidden = self.lstm_cell(input.view(100,-1), hidden) batches
 ← (input.view(100,-1), hidden)
 no. of levels |

LSTM cell allows for multiple layers of transformations

```
output = self.h2o(hidden[0])
```

↑
comprises of the
hidden cell
state

```
output = self.softmax(output)
```

```
return output.view(1, -1), hidden
```

→ to be sent
for the next
character

```
def init_hidden(self):
```

```
return (torch.zeros([1, 1, self.hidden_size]),  
       torch.zeros(1, 1, self.hidden_size))
```

↓

different from RNN, now we have a tuple of two
tensors. one for hidden state and one for cell
state.

n-hidden = 128

```
net = LSTM_net(n_letters, n_hidden, n_languages)
```

```
train_setup(net, lr=0.0005, n_batches=100, batch_size=256)
```

GRU

*> GRU cell

CLASS `torch.nn.GRU(*args, **kwargs)`

[SOURCE]

Applies a multi-layer gated recurrent unit (GRU) RNN to an input sequence.

For each element in the input sequence, each layer computes the following function:

$$r_t = \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr})$$

$$z_t = \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz})$$

$$n_t = \tanh(W_{in}x_t + b_{in} + r_t * (W_{hn}h_{(t-1)} + b_{hn}))$$

$$h_t = (1 - z_t) * n_t + z_t * h_{(t-1)}$$

→ only state to save

where h_t is the hidden state at time t , x_t is the input at time t , $h_{(t-1)}$ is the hidden state of the layer at time $t-1$ or the initial hidden state at time 0, and r_t, z_t, n_t are the reset, update, and new gates, respectively. σ is the sigmoid function, and $*$ is the Hadamard product.

In a multilayer GRU, the input $x_t^{(l)}$ of the l -th layer ($l \geq 2$) is the hidden state $h_t^{(l-1)}$ of the previous layer multiplied by dropout $\delta_t^{(l-1)}$ where each $\delta_t^{(l-1)}$ is a Bernoulli random variable which is 0 with probability dropout.

```

class GRU_net(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super(GRU_net, self).__init__()
        self.hidden_size = hidden_size
        self.gru_cell = nn.GRU(input_size, hidden_size)
        self.h2o = nn.Linear(hidden_size, output_size)
        self.softmax = nn.LogSoftmax(dim=2)

    def forward(self, input_, hidden):
        out, hidden = self.gru_cell(input_.view(1, 1, -1),
                                    hidden)
        output = self.h2o(hidden)
        output = self.softmax(output)
        return output.view(1, -1), hidden

    def init_hidden(self):
        return torch.zeros(1, 1, self.hidden_size)

n_hidden = 128
net = GRU_net(n_letters, n_hidden, n_languages)
train_setup(net, lr=0.0005, n_batches=100, batch_size=256)

```

topk accuracy may be misleading
↑

-) exercises:
 - 0. compute language wise accuracy, confusion matrix
 - i, j : if ground truth was i, how many times I am classifying it as j
 - 1. play with hyper-parameters (including hidden layer size) & trying to get better results.
 - 2. observe gradient explosion at higher learning rates
 - 3. increase performance by moving to GPUs
 - 4. Think on how to increase performance further (benefit of batching)

- *) SIMD : Single instruction multiple data (benefit of batching on GPU)