

The movie review dataset is a built-in feature of keras, and so we can take advantage of a large set of data which has already been labeled. This will allow us to create models and see how effective they are.

2 A reference for keras

You can find an introduction to keras in the book:

Francois Chollet,
Deep Learning with Python,
Manning, 2018,
ISBN: 9781617294433
<https://www.manning.com/books/deep-learning-with-python>

The IMDB movie review exercise is covered in section 3.4, *Classifying movie reviews: a binary classification example*.

3 Where can I run keras?

We have three options for running the movie review example with keras. You can run it

1. on your laptop;
2. at the PSC interactively, using the `interact -gpu` option;
3. at the PSC noninteractively (“batch mode”), using a SLURM script;

4 What is the IMDB data and how is it used?

The IMDB dataset consists of 50,000 movie reviews and 50,000 labels. Each review is labeled “0” if it was judged to be negative, or “1” if positive. These labels were supplied by humans who read the reviews. Our goal is to come up with a procedure that can automatically produce a label for a movie review, and which will closely match the behavior of human readers.

A dictionary was created from all the words in all the reviews, and each word in that dictionary was given an index. Then each review was used to generate a corresponding file of numbers, where each word was replaced by its index.

Thus the text of a movie review is now numeric. For technical reasons, we wish to consider only the 10,000 most common words, so each numeric file is modified to eliminate unusual words. When we actually process a movie review, we do one last step: we replace the file of numbers by a vector of length 10,000, where entry i of the vector is set to 1 if `word[i]` appeared at least once in the review. The reason for doing this is that the neural network needs to process vectors of a uniform size. We keep the neural network happy by making every movie review a 10,000 entry vector of 0’s and 1’s.

We divide the data into three sets: training, validation, and testing data. We will build a model with the training data, and then use it to predict the labels on the validation data. The prediction failures will be used to adjust the model. We will do this adjustment a fixed number of times (perhaps 10 or 20 “epochs”) and then declare the model ready for testing.

We now hold the model fixed, and try it out on the testing data. If the training procedure was done well, then the model should have good accuracy in predicting the labels for the testing data. If the model does poorly, then we must go back and adjust our model and repeat the entire process.

5 Peeking at the reviews

The movie reviews in the IMDB dataset are no longer human readable; they are just lists of numbers. However, it is possible to use the dataset to decode any review. The file *imdb_decode.py* can be used this way:

```
python3
from imdb_decode import imdb_decode
imdb_decode(7)

? lavish production values and solid performances in this straightforward adaption of jane ?
  satirical classic about the marriage game within and between the classes in ? 18th
  century england northam and paltrow are a ? mixture as friends who must pass through ?
  and lies to discover that they love each other good humor is a ? virtue which goes a
  long way towards explaining the ? of the aged source material which has been toned down
  a bit in its harsh ? i liked the look of the film and how shots were set up and i
  thought it didn't rely too much on ? of head shots like most other films of the 80s and
  90s do very good results
```

The question marks in the listing indicate unusual words that were not in the top 10,000 most common. You should be able to guess that this review is labeled 1 **positive**. We hope that our movie classifier will also be able to correctly label it.

6 Getting the IMDB data

If you are planning to work on your laptop, then the IMDB example program will automatically use the Internet to download a copy of the IMDB dataset that it needs, storing it in a location on your machine.

If you are working on the PSC, however, programs running on the computational nodes are not allowed to use the Internet, so we have to get a copy of the dataset beforehand. I have a copy in my account, and you can get a copy for yourself by logging into your PSC account and issuing the following commands:

```
cp ~/jburkard/imdb_copy.sh imdb_copy.sh
bash imdb_copy.sh
```

The datafiles will be copied to a hidden subdirectory in your account. If you want to check that they are there, try this command:

```
ls .keras/datasets
```

You should see two files: *imdb.npz* and *imdb_word_index.json*. The *imdb.npz* is binary (not text!) but the *imdb_word_index.json* actually contains the list of words and their assigned indices.

On your PSC account, you could peek at this file by the command:

```
more .keras/datasets/imdb_word_index.json
```

which would list the beginning of the file:

```
{"fawn": 34701, "tsukino": 52006, "nunnery": 52007, "sonja": 16816,
"vani": 63951, "woods": 1408, "spiders": 16115, "hanging": 2345,
"woody": 2289, ...
```

7 keras on your laptop

You may find it tricky to install keras on your laptop. However, if you can get it set up, you may prefer to do your work there. The movie review is large, but not enormous, so it should run fairly quickly for you.

Installation information can be found on the keras website `keras.io`.

Mac and Linux users can try these install commands:

```
sudo pip install tensorflow
sudo pip install keras
```

while Windows users may try to install with:

```
pip install tensorflow
pip install keras
```

If your installs are successful, then you should download the file `imdb.py` from the lab website. Then you can execute it on your laptop with a command like:

```
python3 imdb.py
```

On Mac and Linux, you can save the “interesting” output to a separate text file:

```
python3 imdb.py > imdb.txt
```

8 keras with an interactive PSC GPU

Copy the file `imdb_psc.py` into your home directory at PSC. There are two ways to do this:

- Download the file from the lab webpage to your own system. Then use `sftp` to transfer it to PSC;
- OR...;
- Use `ssh` to log into your PSC account, and copy my version, by typing

```
cp ~/jburkard/imdb_psc.py imdb_psc.py
```

Assuming you have already copied the IMDB datasets to your PSC account, (see above) you are ready to run the example.

I assume you are logged in to one of the PSC login nodes. Now request interactive access to a GPU. You may have to wait some time before you are granted access. Then type something like this:

```
interact -gpu
module load anaconda3/2019.03
source activate keras-gpu      <— This logs you onto a gpu.
python imdb_psc.py             <— You will see lots of messages.
python imdb_psc.py > imdb.txt  <— Repeat, but save interesting output.
exit                           <— log off from the gpu
```

You don’t have to run twice; it’s just good to run once and see if the program works at all. Then a second run allows you to save the output you are interested in.

Also, you don’t actually have to log out right away. You can stay logged in to the interactive GPU, run your program, modify it with `nano`, and rerun it, so that you can carry out various experiments. But be sure to log out when you are done so that the next user can access that GPU.

9 keras with a noninteractive PSC GPU

If you don’t want to wait for access to an interactive GPU, you can run the IMDB example by submitting it to the PSC job scheduler. Along with the python/keras script, you also need a special job file, which we have prepared for you, called `imdb_psc.sh`.

Copy the files `imdb_psc.py` and `imdb_psc.sh` into your home directory at PSC. There are two ways to do this:

- Download the files from the lab webpage to your own system. Then use `sftp` to transfer them to PSC;
- OR...;
- Use `ssh` to log into your PSC account, and copy my versions, by typing

```
cp ~/jburkard/imdb_psc.py imdb_psc.py
cp ~/jburkard/imdb_psc.sh imdb_psc.sh
```

Assuming you have already copied the IMDB datasets to your PSC account (see above), you are ready to submit your job.

As you are logged in on one of the PSC login nodes, simply type the command

```
sbatch imdb_psc.sh
```

This will send your job to be executed, and give you a six digit number that identifies the job. You can check right away that your job has gotten onto the list of things to be executed by typing

```
squeue -u joeuser
```

where you want to replace `joeuser` by your PSC username.

If you want, you can log out now and check on your job later, by logging back in and issuing the `squeue` command.

If the job doesn't show up in the queue anymore, it has run. At the very least, you will have a file `imdb_psc.log`, which will mostly contain uninteresting system messages, and which we will only check if your job did not run properly. If your job did run, you will also have a file `imdb_psc.txt`, which will contain the printed output from your keras example. You can examine this using the command

```
more imdb_psc.txt
```

or you can use the `nano` editor to view it. Another option is to start an `sftp` session on your laptop and then request that a copy be brought back to your home system:

```
get imdb_psc.txt
```

10 Understanding the code

You will be working with the file `imdb.py`, or, if you are running on the PSC, the slightly different file `imdb_psc.py`. In both cases, the program can be thought of as having these parts:

1. Load the data, and prepare it for use;
2. Describe the model;
3. Create the model;
4. Use the model on training and validation data;
5. Evaluate the model on new test data;

We will not worry about how the data is loaded and prepared, except to note that keras will need to download the data from the Internet (if you are running on your laptop), or will expect to find the data in your home directory (if you are running at the PSC). It will then rearrange and split the data so that it has the right shape for the neural network, and is divided into training, validation, and test sets.

The model is described as a sequential model, with two hidden layers, each with 16 hidden units, and a `relu` activation. The first layer expects an input vector of length 10,000; in other words, one of our movie reviews. Our output layer uses the `sigmoid` function, which returns a value between 0 and 1, the probability that the movie review is negative or positive.

```

model = models.Sequential()

model.add ( layers.Dense ( 16, activation = 'relu', input_shape = (word_num,) ) )
model.add ( layers.Dense ( 16, activation = 'relu' ) )
model.add ( layers.Dense ( 1, activation = 'sigmoid' ) )

```

The model is created by choosing an optimizer, loss function, and a metric. We have discussed the `rmsprop` optimizer in class. The `binary_crossentropy` loss function is a way of measuring the difference between two probability distributions.

```

model.compile (
    optimizer = 'rmsprop',
    loss = 'binary_crossentropy',
    metrics = ['accuracy'] )

```

Once the model is created, we apply it to the training and validation data, and save a report in a dictionary called `history`:

```

history = model.fit (
    partial_x_train ,
    partial_y_train ,
    epochs = 20,
    batch_size = 512,
    validation_data = (x_val, y_val) )

```

We can use `history` to print out the final values of the validation loss and accuracy. The code prints these values like this:

```

Model loss and accuracy on validation data:
Final validation loss 0.7035827110290528
Final validation accuracy 0.8651999831199646

```

Now that the model has been trained, we want to see how well it can handle new data. We carry out this experiment on the test data, for which we know the correct result, and we report how well our model does:

```

results = model.evaluate ( x_test , y_test )

for i in range ( len ( model.metrics_names ) ):
    print ( model.metrics_names[i], results[i] )

```

The test loss and test accuracy are printed out something like this:

```

Model loss and accuracy on test data:
loss 0.3238627934074402
accuracy 0.87308

```

11 Choose an experiment

The validation accuracy and test accuracy measure how well our classifier performs on data for which it had not been trained. In the example code, after the 20th epoch, the validation accuracy was about 0.865 and the test accuracy was about 0.873.

The accuracy of our model is affected by the parameter choices that were made in the program. In each of the following experiments you are to vary one of the parameters in the model, and note the resulting validation and test accuracy for each parameter choice.

The list of experiments for you to choose from includes:

1. The example used the command `model.add(layers.Dense())` twice, specifying the `relu` activation function. Compare using `relu` versus using the `tanh` activation function.
2. The example used the command `model.add(layers.Dense())` twice, to set up two hidden layers. Compare using one, two, or three hidden layers.
3. The example used the command `model.add(layers.Dense())` with 16 units in the two hidden layers. Compare using 16, 32, or 64 hidden units in each layer.
4. The example used a `model.compile()` command in which the optimizer was `rmsprop`. Compare using `rmsprop` versus just one of the other optimizers on this list: `sgd`, `adagrad`, `adadelta`, `adamax`.
5. The example used a `model.compile()` command in which the loss function was `cross_entropy`. Compare the results for `cross_entropy` versus using `mse`.
6. The example used a `model.fit()` command in which 20 epochs of training were carried out. Compare your results using 5, 10, and 20 epochs.

12 Computing Assignment #11

Do **two** of the experiments from the above list. Write a short report in which you explain:

1. which two experiments you carried out;
2. where you did your experiments: laptop, interactive gpu, noninteractive gpu;
3. the validation and test accuracies for your several cases;

Example:

MATH 1900 Report for Joe User

I ran experiments #1 and #6.

I ran the experiments on my laptop.

Here are my tables:

Experiment #1:

activation	validation	test
<code>relu</code>	0.86	0.87
<code>tanh</code>	0.73	0.68

Experiment #6:

epochs	validation	test
5	0.80	0.71
10	0.82	0.75
20	0.86	0.87

Email your report to Dr Schneier mhs64@pitt.edu before Wednesday, 4 December.