

ML4Science: Week 4 Meeting

Calafà, M., Mescolini, G., Motta, P.

École Polytechnique Fédérale de Lausanne (EPFL)

Machine Learning Project 2

14 Dec 2021

Tutor: Dr. Michele Bianco



Outlines

1 Lightening the computational effort (part 2)

- Data cleaning

2 Process reloading

- Process reloading

3 Other additions

- Plots
- Other additions

4 Issues and Questions

- Questions

Lightening the computational effort (part 2)

Data cleaning

With the garbage collector (`gc.collect()`) we cleaned run-time the memory previously allocated for:

- Macro-tensor X (with all data from the dataset)
- Sub-tensors of X (train, test ...)
- Temporary datasets

In order to keep in memory the data-loader only. Many GB saved, lower risk of sudden kernel stops.

Process reloading

Process reloading

The training process is very slow if repeated for many epochs.

Idea: save locally the main parameters at the end of every epoch.

Two benefits:

- 1 Possibility to stop the computation and restart it later
- 2 Backup of the information for undesirable breaks

Process reloading

It is important to distinguish between the very first run from the following reloadings:

```

1
2  if (first_run == True):
3      net = CNN.CNN()
4      optimizer = optim.Adam(net.parameters(), lr=1e-3)
5      current_epoch = 0
6      final_epoch = epochs
7      prev_loss = 10**2 # high initial value
8      all_test_losses = []
9      all_train_losses = []
10     all_R2_train = []
11     all_R2_test = []
12  else:
13     #Resume the training
14     PATH = './model/last_model.pt'
15     net = CNN.CNN()
16     optimizer = optim.Adam(net.parameters(), lr=1e-3)
17     checkpoint = torch.load(PATH)
18     net.load_state_dict(checkpoint['model_state'])
19     optimizer.load_state_dict(checkpoint['optimizer_state'])
20     scheduler.load_state_dict(checkpoint['scheduler_state'])
21     current_epoch = checkpoint['epoch'] + 1
22     prev_loss = checkpoint['loss']
23     final_epoch = current_epoch + epochs
24     ...

```

Other additions

Plots

We added the correlation plot ($x_{i_{predicted}}$ vs $x_{i_{true}}$) and the loss trend with respect to the epoch, both for training and for testing.

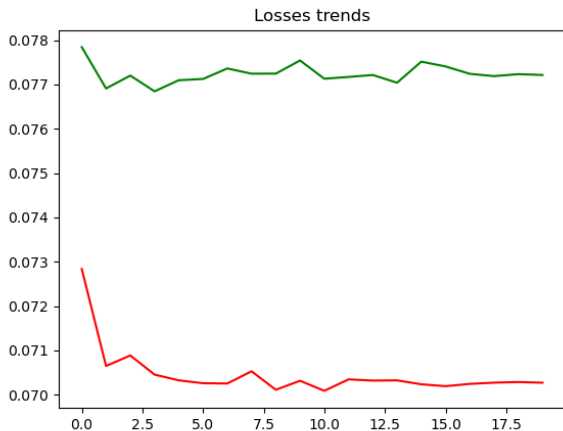
```

1
2 def correlation_plot(x_pred, x_true):
3     plt.plot(x_true, x_true, 'r') # y = x
4     plt.plot(x_true, x_pred, 'b') # our actual prediction
5     sigma = np.std(x_pred)
6     plt.plot(x_true, x_pred + sigma, 'r-')
7     plt.plot(x_true, x_pred - sigma, 'r-')
8
9 def plot_losses(epochs, loss_tr, loss_te):
10    plt.plot(epochs, loss_tr, 'r') # training losses
11    plt.plot(epochs, loss_te, 'g') # test losses
12    plt.title('Losses trends')
13    plt.show()

```

Plots

Here we report the above mentioned loss plot obtained with our results (green is test, red is training):



Other additions

■ Scheduler

```

1 from torch.optim.lr_scheduler import ReduceLROnPlateau
2 # ...
3 scheduler = ReduceLROnPlateau(optimizer = optimizer, mode = 'min', factor
    = 0.1, patience = 7, min_lr = 1e-7)
4 # ...
5 scheduler.step(loss_test)

```

■ R2 score

```

1 from sklearn.metrics import r2_score
2 # ...
3 R2 = r2_score(y_train.detach(), output.detach())

```

■ Extraction of Validation set

```

1 # split dataset into training (2500), validation set (500) and prediction
  set (300)
2 X_train, X_pred, y_train, y_pred = train_test_split(X, y, test_size=0.1,
    random_state=2021)
3 X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train,
    test_size=0.2, random_state=2021)

```

Issues and Questions

Issues

- 1 Issues with Paolo's Laptop: we could only reach 20 epochs. Is there a possibility of usage of EPFL Cluster?
- 2 About February conference