

Deep Learning in practice

Lecturers: G. CHARPIAT

(January, 2023)

Solution by **DI PIAZZA Théo****Instructions**

- This file is the report of the TP1.

1 Problem 1 - MultiClass

In this exercise, we are interested in the classification problem of the USPS dataset (10 classes, 16x16 pixels for each image, 9298 samples). Below is an example of the samples :



Figure 1: Images from USPS dataset.

The objective is to understand the different architectures of a neural network and the influences of the different hyper-parameters on the performance of the model.

The best model is the one that performed the best on the test sample. It is the model trained with the following configuration : Adam [lr 0.001] on 40 epochs, batch size = 10, MSELoss and 1 CNN (3x3 kernel, 3 channels and 2 for padding) + ReLU as hidden layer.

For our best model, we obtain a precision of 92.23% on the test sample.

Concerning the details of the experiments, **click here to access the Appendix of problem 1**. Remark : *For each experiment, we change the variable one at a time to estimate its "direct impact" on the model. It is important to note that the impact of a hyper-parameter can be different depending on the value of the other hyper-parameters (e.g. the influence of the learning rate and the batch size on the quality of the learning). This has been taken into account for the final model.*

2 Problem 2 - Regression

In this problem, the goal is to predict, from 1460 samples, the house selling prices from 4 features (OverallQual, YearBuilt, TotalBsmtSF, GrLivArea).

	OverallQual	YearBuilt	TotalBsmtSF	GrLivArea	SalePrice
0	7	2003	856	1710	208500
1	6	1976	1262	1262	181500
2	7	2001	920	1786	223500

Figure 2: Header of the dataset used for problem 2.

In the first part, the criterion that is used is the MSE Loss function. Due to a gradient explosion, the data are normalized. First, all the data were normalized: observations and targets. In a second step, only the targets have been normalized. This also solves the problem, and the comments on the results

obtained are the same. As in the previous section, different configurations are tested. For each configuration, only one variable is modified in order to see its "direct" impact on the performance of the trained model. In a second step, the Gaussian Loss is used, in order to predict a probability of distribution rather than a quantity, notably in order to manage uncertainties with the variance.

The selected model is the one with the lowest loss on the validation sample (Gaussian Loss : 0.22) with the following configuration : 3 linear layers of 20 neurons followed by ReLU, on 20 epochs with a batch size of 10, and RMSprop optimizer (lr : 0.001). For our best model, we obtain a Gaussian Loss of 0.33 on the test sample, with normalized data.

Finally, as proposed by the subject, we compare the standard deviation obtained on the test sample with the "Overall Quality" variable: the uncertainty will increase when the feature "OverallQual" increase.

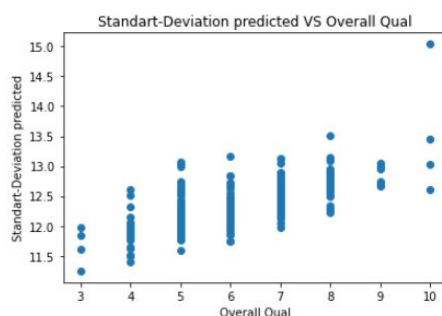


Figure 3: Predicted standard-deviation on test set VS Overall Quality.

Concerning the details of the experiments, [click here to access the Appendix of problem 2](#). For the sake of simplicity, we only present the results obtained with the Gaussian Loss because the observations are similar.

3 Observations

Impact of the number of linear layers From the neural network with 1 single linear hidden layer, we propose to see the impact of adding hidden layers on the model performance. Cela augmente la complexité du calcul et donc les temps de calculs. Concernant les performances, les performances obtenues sont équivalentes.

Impact of the number of neurons From the model with a single linear hidden layer, we modify the number of neurons of this layer to see the impact on the model. As the number of neurons increases, the accuracy of the model increases. Beyond a certain number of neurons, the accuracy does not increase anymore. This increases the complexity of the model and thus the computation time. Quand le nombre de neurones est trop élevé, un risque de sur-apprentissage sur l'échantillon d'entraînement peut être observé, il faut alors trouver un équilibre entre performance et complexité.

Impact of the activation function From the model with a single linear layer followed by an activation function, which precedes the last linear output layer, we compare the results obtained using: ReLU, LeakyReLU, Tanh and Sigmoid. The results seem to be equivalent for ReLU, LeakyReLU and Tanh but underperform for Sigmoid. Indeed, the Sigmoid activation function is generally used after the final layer of the network in order to obtain values interpretable as probabilities, not inside a hidden layer. The execution times are almost equivalent.

Impact of the number of layers in CNN For CNN architectures, we propose to try different numbers of layers. We notice that the results are almost equivalent. Generally, we say "the deeper the better" but it is not always the case, maybe because the classification task is relatively simple. As the number of

layers increases, the complexity of the model increases and the computation time increases.

Impact of batch size For both problems, different batch sizes are tested (from 10 to 1000). The more the batch size increases, the more the precision decreases. Indeed, when the batch size is modified, it is also necessary to adjust the learning rate. The larger the batch size, the shorter the execution time. This is consistent with the fact that more samples are processed before each update of the model parameters.

Impact of the learning rate For each optimizer tested, different learning rates are tried (from 0.001 to 10). The performance of the models varies depending on the optimizer. When the learning rate is too low or too high, the model does not seem to learn correctly. For each optimizer, it is necessary to use an adapted learning rate which allows to decrease the loss in the first epochs. If needed, the learning rate can be adjusted after several epochs to refine the model.

Impact of the number of epochs We compare for 5, 10, 20 and 40 epochs. When the number of epochs increases, the accuracy of the model increases and the execution time also increases. If the number of epochs is too high, it can happen that the train loss continues to decrease but the test loss stagnates or even increases again. In this case, an early stopping is necessary.

Impact of the optimizer We compare 3 optimizers with the adjusted learning rate: SGD (0.1), Adam (0.001) and RMSprop (0.001) The results are very good and almost equivalent, but the execution times are a bit longer for Adam and RMSprop.

Impact of the loss function For the first problem propose to use the `nn.CrossEntropyLoss` instead of the `MSE`. With an Adam optimizer (lr 0.001), batch size of 10 out of 10 epochs, we obtain an accuracy of 94.81% on the validation set. This loss seems well adapted to our problem. Even if in theory, the CrossEntropy loss is more adapted to classification problems, here it is the MSE which allows to obtain the best results.

General comment on the execution time To summarize what has been said before, the execution times here are very short because the data set is of reasonable size and the models used have few parameters. On average, the training of a model on a dozen epochs did not last more than a few tens of seconds, which made it possible to carry out all these executions in less than 1 hour in cumulative execution time. The strategy was to voluntarily choose the configurations likely to interest us, because the possible combinations of the various configurations to be tested are « almost without limit ».

Appendix

Problem 1 - Exercice 1

Impact of the number of linear layers			
Number of hidden layers	Train accuracy	Test accuracy	Execution time
1	93.30%	92.40%	7.65s
2	90.10%	90.09%	8.35s
3	57.25%	58.71%	9.13s

Impact of the number of neurons			
Number of neurons	Train accuracy	Test accuracy	Execution time
20	91.37%	90.47%	7.45s
50	92.98%	91.87%	7.45s
100	93.30%	92.40%	7.65s
150	93.15%	92.25%	7.95s

Impact of the activation function			
Activation function	Train accuracy	Test accuracy	Execution time
ReLU	91.37%	90.47%	7.45s
Tanh	93.25%	92.33%	7.67s
Sigmoid	47.22%	48.31%	8.30s

Impact of the number of CNN layers			
Number of hidden layers	Train accuracy	Test accuracy	Execution time
1	96.02%	94.58%	8.92s
2	95.60%	93.49%	12.16s
3	93.15%	93.15%	15.45s

Impact of the activation function after CNN			
Activation function	Train accuracy	Test accuracy	Execution time
ReLU	96.02%	94.58%	8.92s
LeakyReLU	95.73%	94.58%	9.48s
Tanh	95.73%	94.04%	9.37s

Problem 1 - Exercice 2

Impact of the batch size			
Batch size	Train accuracy	Test accuracy	Execution time
10	96.02%	94.58%	8.92s
50	92.82%	92.02%	6.37s
100	89.12%	89.31%	5.63s
250	32.48%	35.17%	5.39s
500	27.25%	28.81%	5.20s
1000	8.10%	7.28%	5.41s

Impact of the learning rate			
Learning rate	Train accuracy	Test accuracy	Execution time
0.001	26.40%	28.51%	9.51s
0.01	71.45%	73.43%	9.16s
0.1	96.02%	94.58%	8.92s
0.5	96.62%	94.73%	9.36s
1	97.22%	94.97%	8.92s
10	16.13%	17.51%	9.32s

Impact of the number of epochs			
Epochs	Train accuracy	Test accuracy	Execution time
5	96.77%	94.89%	5.61s
10	796.62%	94.73%	9.36s
20	98.90%	96.59%	21.11s
40	99.22%	96.90%	42.04ss

Impact of the optimizer			
Optimizer (lr)	Train accuracy	Test accuracy	Execution time
SGD (0.1)	96.62%	94.73%	9.36s
Adam (0.001)	97.27%	95.12%	11.03s
RMSprop (0.001)	96.95%	94.89%	11.04s

Problem 2 - Gaussian Loss

Impact of the number of layers			
Number of layers	Train Loss	Test Loss	Execution time
1	0.20	0.23	2.87s
2	0.21	0.25	2.86s
3	0.25	0.27	3.78s

Impact of the number of neurons			
Number of neurons	Train Loss	Test Loss	Execution time
5	0.91	0.89	4.08s
10	0.21	0.25	4.39s
20	0.16	0.22	4.40s
40	0.25	0.27	4.42s
50	0.24	0.35	4.71s

Impact of the activation function			
Activation function	Train Loss	Test Loss	Execution time
ReLU	0.18	0.26	4.42s
leakyReLU	0.17	0.23	4.43s
Sigmoid	0.40	0.44	4.12s
Tanh	0.26	0.33	5.13s

Impact of the batch size			
Batch size	Train Loss	Test Loss	Execution time
10	0.17	0.23	4.43s
50	0.24	0.27	1.99s
100	0.83	0.80	1.09s
200	1.02	0.97	0.57s
400	1.19	1.17	0.56s

Impact of the optimizer			
Optimizer (lr)	Train Loss	Test Loss	Execution time
RMSprop (0.001)	0.17	0.23	4.43s
Adam (0.001)	0.17	0.23	5.90s
SGD (0.01)	0.48	0.50	3.92s

Impact of the number of epochs			
Number of epochs	Train Loss	Test Loss	Execution time
10	0.26	0.29	3.14s
20	0.17	0.23	5.90s
50	0.18	0.32	12.85s
100	0.07	0.27	28.40s