
MLP Coursework 1: Learning Algorithms and Regularization

sXXXXXXXX

Abstract

The abstract should be 100–200 words long, providing a concise summary of the contents of your report.

1. Introduction

This document provides a template for the MLP coursework 1 report. This template structures the report into sections, which you are recommended to use, but can change if you wish. If you want to use subsections within a section that is fine, but please do not use any deeper structuring. In this template the text in each section will include an outline of what you should include in each section, along with some practical LaTeX examples (for example figures, tables, algorithms). Your document should be no longer than **five pages**, with an additional page (or more!) allowed for references.

The introduction should place your work in context, giving the overall motivation for the work, and clearly outlining the research questions you have explored – in this case questions arising from an investigation of learning rate algorithms, learning rate schedulers, and regularization/weight decay. This section should also include a concise description of the Balanced EMNIST task and data – be precise: for example state the size of the training, validation, and test sets.

2. Baseline systems

In this section, in order to get the best performance of stochastic gradient descent (SGD) in 2-5 hidden layers as the baseline, I debugged the learning rate hyperparameter. Firstly I will compare these systems using the valid data set and find the best learning rate.

From the data in the table above, we compare the respective acc(valid) of 2-5 hidden layers, we can get the best learning rate as shown in the following table.

Then I use the test data to compare the accuracy of these best systems.

The accuracy rate is the highest and 3 hidden layers and 4 hidden layers, the accuracy rate is about 0.83, and the accuracy of 2 hidden layers and 5 hidden layers is only about 0.82. Compared with the valid data, these systems have a 1

In the following experiments, I will use this part of the best 3 hidden SGD as the baseline to compare the performance

Table 1. Add caption

2 Hidden Layers				
learning_rate	acc(train)	acc(valid)	error(train)	error(valid)
0.1	9.32E-01	8.15E-01	1.72E-01	8.84E-01
0.05	9.24E-01	8.28E-01	2.01E-01	6.45E-01
0.01	8.71E-01	8.35E-01	3.84E-01	5.12E-01
0.0075	8.59E-01	8.30E-01	4.27E-01	5.19E-01
0.005	8.44E-01	8.22E-01	4.88E-01	5.59E-01

Table 2. Add caption

3 Hidden Layers				
learning_rate	acc(train)	acc(valid)	error(train)	error(valid)
0.1	9.32E-01	8.15E-01	1.72E-01	8.84E-01
0.05	9.24E-01	8.28E-01	2.01E-01	6.45E-01
0.01	8.71E-01	8.35E-01	3.84E-01	5.12E-01
0.0075	8.81E-01	8.42E-01	3.50E-01	4.92E-01
0.005	8.44E-01	8.22E-01	4.88E-01	5.59E-01

Table 3. Add caption

4 Hidden Layers				
learning_rate	acc(train)	acc(valid)	error(train)	error(valid)
0.1	9.29E-01	8.17E-01	1.73E-01	9.84E-01
0.05	9.36E-01	8.27E-01	1.58E-01	8.94E-01
0.01	9.00E-01	8.39E-01	2.81E-01	5.18E-01
0.0075	8.88E-01	8.42E-01	3.16E-01	4.89E-01
0.005	8.65E-01	8.28E-01	3.92E-01	5.27E-01

Table 6. Add caption

SGD for test set		
Layers(learning_rate)	acc(train)	acc(test)
2 Hidden Layers(0.01)	8.71E-01	8.23E-01
3 Hidden Layers(7.5e-3)	8.81E-01	8.34E-01
4 Hidden Layers(7.5e-3)	8.88E-01	8.31E-01
5 Hidden Layers(5e-3)	8.75E-01	8.22E-01

of subsequent algorithms.

Table 4. Add caption

5 Hidden Layers				
learning_rate	acc(train)	acc(valid)	error(train)	error(valid)
0.1	9.27E-01	8.17E-01	1.81E-01	9.47E-01
0.05	9.36E-01	8.29E-01	1.57E-01	8.82E-01
0.01	8.99E-01	8.32E-01	2.70E-01	5.47E-01
0.0075	8.88E-01	8.32E-01	3.09E-01	5.41E-01
0.005	8.75E-01	8.33E-01	3.58E-01	5.16E-01
0.004	8.63E-01	8.27E-01	3.98E-01	5.35E-01

In this section you should report your baseline experiments for EMNIST. No need for theoretical explanations of things covered in the course, but if you go beyond what was covered please explain what you did with references to relevant paper(s) if appropriate. In this section you should aim to cover the both the “what” and the “why”: *what* you did, giving sufficient information (hyperparameter settings, etc.) so that someone else (e.g. another student on the course) could reproduce your results; and *why* you performed the experiments you are reporting - what you are aiming to discover, what is the motivation for the particular experiments you undertook. You should also provide some discussion and interpretation of your results.

Your experimental sections should include graphs (for instance, figure 1) and/or tables (for instance, table 7)¹, using the figure and table environments, in which you use `\includegraphics` to include an image (pdf, png, or jpg formats). Please export graphs as **vector graphics** rather than **raster files** as this will make sure all detail in the plot is visible. Matplotlib supports saving high quality figures in a wide range of common image formats using the `savefig` function. **You should use `savefig` rather than copying the screen-resolution raster images outputted in the notebook.** An example of using `savefig` to save a figure as a PDF file (which can be included as graphics in a LaTeX document is given in the coursework document.

Table 5. Add caption

the best learning rate for 2-5 Hidden Layers	
number of Hidden layers	the best learning rate
2 Hidden layers	1.00E-02
3 Hidden layers	7.50E-03
3 Hidden layers	7.50E-03
4 Hidden layers	5.00E-03

If you need a figure or table to stretch across two columns use the `figure*` or `table*` environment instead of the `figure` or `table` environment. Use the `subfigure` environment if you want to include multiple graphics in a single figure.

3. Learning algorithms – RMSProp and Adam

In this section you should compare RMSProp and Adam with gradient descent, introducing these learning rules ei-

¹These examples were taken from the ICML template paper.

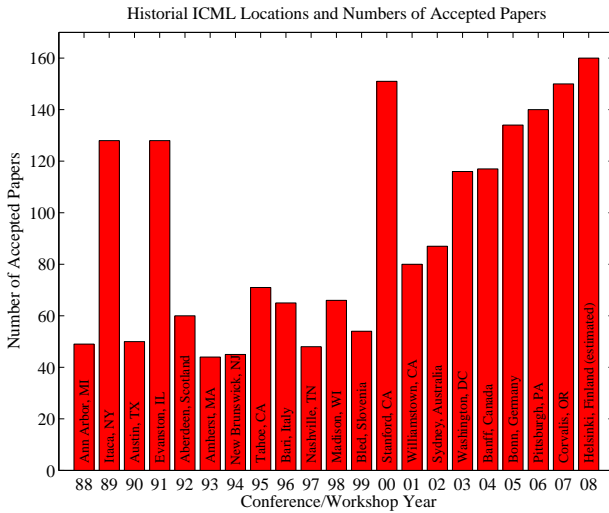


Figure 1. Historical locations and number of accepted papers for International Machine Learning Conferences (ICML 1993 – ICML 2008) and International Workshops on Machine Learning (ML 1988 – ML 1992). At the time this figure was produced, the number of accepted papers for ICML 2008 was unknown and instead estimated.

DATA SET	NAIVE	FLEXIBLE	BETTER?
BREAST	95.9± 0.2	96.7± 0.2	✓
CLEVELAND	83.3± 0.6	80.0± 0.6	×
GLASS2	61.9± 1.4	83.8± 0.7	✓
CREDIT	74.8± 0.5	78.3± 0.6	
HORSE	73.3± 0.9	69.7± 1.0	×
META	67.1± 0.6	76.5± 0.5	✓
PIMA	75.1± 0.6	73.9± 0.5	
VEHICLE	44.9± 0.6	61.5± 0.4	✓

Table 7. Classification accuracies for naive Bayes and flexible Bayes on various data sets.

ther as equations or as algorithmic pseudocode. If you present the different approaches as algorithms, you can use the `algorithm` and `algorithmic` environments to format pseudocode (for instance, Algorithm 1). These require the corresponding style files, `algorithm.sty` and `algorithmic.sty` which are supplied with this package.

SGD algorithm equation:

$$\begin{aligned} g_t &\leftarrow \nabla J_i(\theta_{t-1}) \\ \theta_t &\leftarrow \theta_{t-1} - \eta g_t \end{aligned} \quad (1)$$

RMSProp algorithm equation:

$$\begin{aligned} g_t &\leftarrow \nabla J_i(\theta_{t-1}) \\ G_t &\leftarrow \gamma G_t + (1 - \gamma) g_t \odot g_t \\ \theta_t &\leftarrow \theta_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \end{aligned} \quad (2)$$

Adam algorithm equation:

$$\begin{aligned} g_t &\leftarrow \nabla J_i(\theta_{t-1}) \\ G_t &\leftarrow \gamma G_t + (1 - \gamma) g_t \odot g_t \\ \theta_t &\leftarrow \theta_{t-1} - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \end{aligned} \quad (3)$$

Algorithm 1 Bubble Sort

Input: data x_i , size m

repeat

Initialize $noChange = true$.

for $i = 1$ **to** $m - 1$ **do**

if $x_i > x_{i+1}$ **then**

Swap x_i and x_{i+1}

$noChange = false$

end if

end for

until $noChange$ is *true*

You should, in your own words, explain what the different learning rules do, and how they differ. You should then present your experiments and results, comparing and contrasting stochastic gradient descent, RMSProp, and Adam. As before concentrate on the “what” (remember give enough information so someone can reproduce your experiments), the “why” (why did you choose the experiments that you performed – you may have been motivated by your earlier results, by the literature, or by a specific research question), and the interpretation of your results.

In every section, you should present your results in a way that makes it easy for a reader to understand what they mean. You should facilitate comparisons either using graphs with multiple curves or (if appropriate, e.g. for accuracies) a results table. You need to avoid having too many figures, poorly labelled graphs, and graphs which should be comparable but which use different axis scales. A good presentation will enable the reader to compare trends in the same graph – each graph should summarise the results relating to a particular research (sub)question.

Your discussion should interpret the results, both in terms of summarising the outcomes of a particular experiment, and attempting to relate to the underlying models. A good report would have some analysis, resulting in an understanding of why particular results are observed, perhaps with reference to the literature. Use bibtex to organise your references – in this case the references are in the file `example-refs.bib`. Here is an example reference (Langley, 2000).

4. Cosine annealing learning rate scheduler

In this section you should present the cosine annealing scheduler, supported using equations or algorithmic pseudocode. Explain the approach in your own words. Following this present your experiments, again remembering to include the “what”, the “why”, and the interpretation of results. Again make sure you provide enough information about the hyperparameters to enable someone else to reproduce your results.

5. Regularization and weight decay with Adam

In this section you should present your experiments on regularization and weight decay with Adam. In your own words explain the problem using L2 regularization with Adam (using both your own results and the discussion by Loshchilov & Hutter (2017)), and the solution proposed by Loshchilov & Hutter. Report the experiments you have carried out, being careful to be clear what research question is being addressed by each experiment – again you should include *what* experiments you performed (include details of hyperparameters, etc.), *why* you performed them (what was the motivation for the experiments, what research questions are you exploring), and the interpretation and discussion of your results.

6. Conclusions

Based on your experiments, what is the setup (models, algorithms) that produces your best system on the Balanced EMNIST data?

You should draw conclusions from the experiments, related to the research questions outlined in the introduction (section 1). You should state the conclusions clearly and concisely. It is good if the conclusion from one experiment influenced what you did in later experiments – your aim is to learn from your experiments. Extra credit if you relate your findings to what has been reported in the literature.

A good conclusions section would also include a further work discussion, building on work done so far, and referencing the literature where appropriate.

References

- Langley, P. Crafting papers on machine learning. In Langley, Pat (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Loshchilov, Ilya and Hutter, Frank. Fixing weight decay regularization in Adam. *arXiv preprint arXiv:1711.05101*, 2017. URL <https://arxiv.org/abs/1711.05101>.