

Training Artificial Neural Networks

Ozgur Gulsuna

^aMiddle East Technical University, Electrical and Electronics Engineering, Ankara, Turkey

Introduction

Insert here your abstract text.

Keywords: Type your keywords here, separated by semicolons ;

1. Basic Concepts

Here introduce the paper, and put a nomenclature if necessary, in a box with the same font size as the rest of the paper. The paragraphs continue from here and are only separated by headings, subheadings, images and formulae. The section headings are arranged by numbers, bold and 10 pt. Here follows further instructions for authors.

1.1. Which Function ?

An ANNs classifier that is trained with cross-entropy loss approximates the conditional probability distribution function. More specifically, for an input data, the output of the classifier is a probability distribution for the classes. The cross-entropy loss function is a measure between the predicted probability distribution and the true distribution. The form of the loss function is decreasing, smooth and differentiable, which makes it easier to optimize using gradient-based methods. This form is also known as the negative log-like function.

1.2. Gradient Computation

1.3. Some Training Parameters and Basic Parameter Calculations

1. The batch refers to a subset of the training data that is used to compute the weights for one iteration. More specifically, the batch size is the number of training samples in a batch. The epoch on the other hand refers to the number of times the entire training data is used to update the weights. In training, there are generally multiple epoch iterations where the weights are updated with different batches/subsets of the training data.
2. For the N number of training samples, the number of batches per epoch is N/B , where B is the batch size. A little side note that the solution is rounded up to the higher integer if N/B is not an integer.
3. For the optimization iterations, such as SGD, for E number of epochs, the total number of iterations is $E \times N/B$. Again, a practical side note states that the N/B is rounded up to the higher integer.

E-mail address: ozgur.gulsuna@metu.edu.tr

1.4. Computing Number of Parameters of ANN Classifiers

1. Starting from the initial layer of the MLP, we have D_{in} number of input neurons and H_1 number of neurons in first hidden layer. Also there are biases associated with each neuron. Therefore, the number of parameters of the each layer is,

$$\begin{aligned} \text{Input Layer} &= D_{in} \times H_1 + H_1 \\ \text{Hidden Layers} &= H_1 \times H_2 + H_2 \\ &\dots \\ \text{More Hidden Layers} &= H_{k-1} \times H_k + H_k \\ \text{Output Layer} &= H_k \times D_{out} + D_{out} \end{aligned}$$

The total sum can be written as, where K is the number of hidden layers.

$$\text{Total Number of Parameters} = D_{in} \times H_1 + \sum_{k=2}^K (H_{k-1} \times H_k + H_k) + H_k \times D_{out} + D_{out}$$

2. CNN structure is more complicated. The number of parameters of a CNN layer is calculated as follows:
For the input layer, the number of parameters is,

$$\text{Input Layer} = (H_{in} \times W_{in} \times C_{in} \times C_1) + C_1$$

where H_{in} and W_{in} are the height and width of the input image, and C_{in} is the number of channels of the input image. Each input of layer is the output of the previous layer. For the convolutional layers, the number of parameters is calculated as,

$$\text{Convolutional Layer} = H_k \times W_k \times C_{k-1} \times C_k + C_k$$

Combination of all layers is,

$$\text{Convolutional Layers} = \sum_{k=2}^K H_k \times W_k \times C_{k-1} \times C_k + C_k$$

Here all the parameters are summed up. The output is assumed to be the last index of the array. The final equation for the total number of parameters is,

$$\text{Total Number of Parameters} = (H_{in} \times W_{in} \times C_{in} \times C_1) + C_1 + \sum_{k=1}^K H_k \times W_k \times C_{k-1} \times C_k + C_k$$

2. Implementing a Convolutional Layer with NumPy

The section involves implementing conv2d function using NumPy for forward propagation and testing it on a small batch of MNIST dataset. We downloaded and loaded input and kernel files, and created an output image using the part2Plots function. The implementation code can be found in the appendix named my_conv2d.py. We confirmed the correctness of our implementation by the output image.

2.1. Experimental Work

The generated output for the convolution over the MNIST dataset is shown in Figure 1.



Fig. 1. Convolution over the number 8 of the MNIST dataset

2.2. Results and Discussion

1. The Convolutional Neural Networks are important for couple of reasons. First of all, when the input shape of the CNN is selected as 2D, it is well suited image processing. Second, the CNN is able to learn the spatial features of the input image. This also means that the CNN is able to learn and extract the features of the input image without any manual work. CNN's are also able to recognize the features of the input image even if the input image is rotated or scaled. Since the features are extracted from image, partial occlusion of the image affect the performance of the CNN less.
2. The kernel of a Convolutional Layer is essentially a matrix of weights that is convolved (inversely correlated) over with the input data to extract features. The size of the kernel refers to the number of rows and columns in the matrix. It corresponds to the reception of the filter, meaning that higher sizes can extract more complex features.
3. The output image shows that the convolution of pre-presented kernels for the number 8 of the MNIST dataset. Basically each filter is convoled over the images to grasp different meanings. Each column is another kernel with each row is different input image.
4. The numbers in the same column look like each other since they both have a representation of the same number and the same kernel is able to extract the features related to the number 8 other than the specific image.
5. The numbers in the same row look different although the input image the same. This is because the kernels are different and they are able to extract different features from the same image.
6. For more specific examples, the third column kernel represents that an 8 has two "islands" of white patches in the middle but the size, shape and location of these pathes differ for each 8 although all of them represent the same thing in a different manner. Another column such as 6, implies the white track like feature of the number 8. In this sense some features are more dinstinctive than other however when different of these combined make the action work even though they do not seem to represent a clear feature. This is similar to human behaviour as we associate the similar patterns to the general inputs and this is the importance of the convolutaional layers, the features can be learned in a sense.

3. Experimenting ANN Architectures

3.1. Experimental Work

This experimental work focuses on testing various Artificial Neural Network (ANN) architectures for a classification task. The models will use adaptive moment estimation (Adam) with default parameters for the optimizer. The datasets will be preprocessed and split into three sets: training, validation, and testing. The ANN architectures to be tested are 'mlp 1', 'mlp 2', 'cnn 3', 'cnn 4', and 'cnn 5', each with their specific layers. For each architecture, the ANN will be trained for 15 epochs, and training loss, accuracy, validation accuracy, and test accuracy will be recorded. The best test accuracy and weights of the first layer will be recorded, and a dictionary object will be created and saved for each architecture. Performance comparison plots will be created, and the weights of the first layer of all architectures will be visualized.

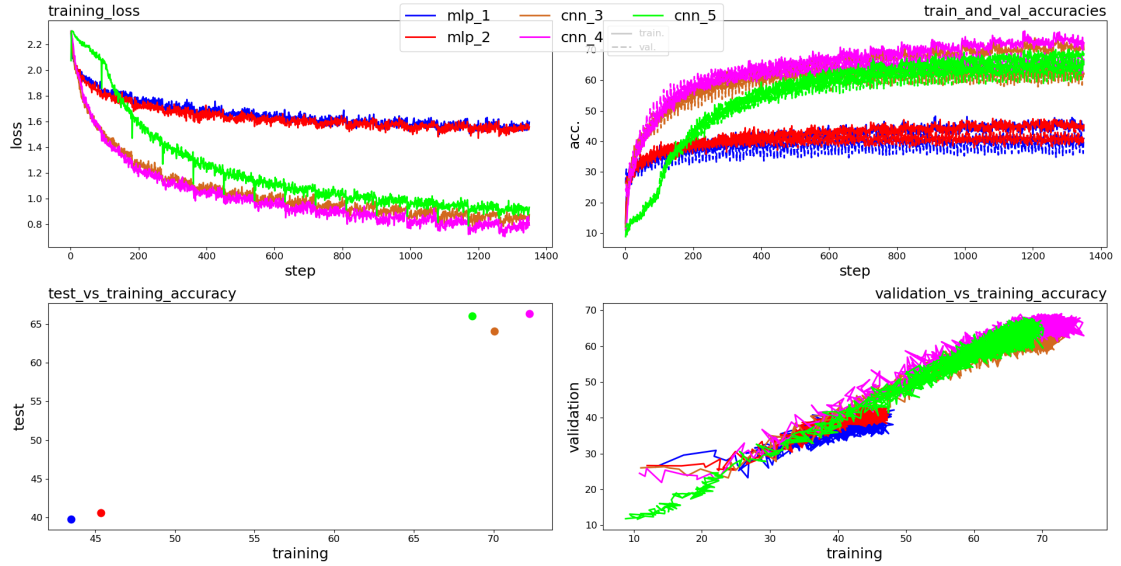
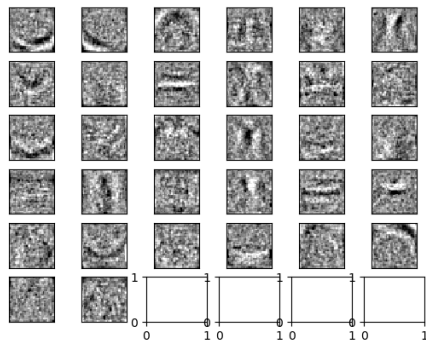


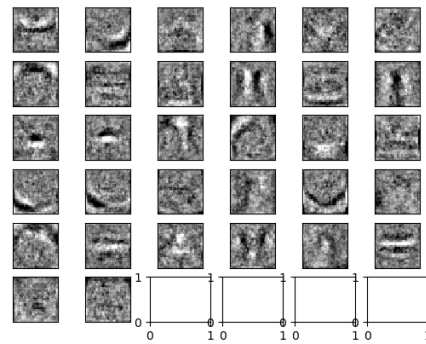
Fig. 2. Performance Comparison Plots for the ANN Architectures

3.2. Results and Discussion

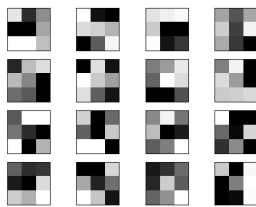
1. Generalization performance refers to the ability of the model to to classify the unseen data. It is used as the ability to recognize patterns and apply this knowledge to the new data.
2. The plots showing the results with the "test" data show the generalization performance since the test data is not used in the training process and the model is not familiar with the test data. Validation data is also seemed to be a good indicator of the generalization performance. However, although the validation data is not directly used in training process, it is used to tune the parameters of the model. Therefore, the model is familiar with the validation data in a sense. The first two curves and the last x-y plot give hint about the comparative generalization performance since these show the results with the validation data. The third scatter plot on the other hand is used with the test data, hence have the correct generalization performance. However it only shows the best run hence the variety of the results are a topic of discussion and the plot does not show that information.
3. Copmarative results show that the convolutional architectures perform better than the multi-layer perceptron architectures. This is because the convolutional layers are able to extract spatial features from the data with more grasping ability. The "mlp_1" and "mlp_2" are very similar in terms of performance although they have different size of parameters and number of layers. The "cnn_3" and "cnn_4" are also very similar again the latter has more layers. The "cnn_5" is more of a slow learner and could not get to the same level of performance as the "cnn_4" but with more epochs it is seem to be able to surpass the "cnn_4" since the gradient of the accuracy is increasing.
4. Higher the number of parameters, it is generally easier for model to learn complex features. However it also means that the model is more prone to overfitting. This is because the model is able to learn the training data in more depth, like its noise characteristics not the required features. Hence, it is not able to generalize well. This is called overfitting, the models with higher parameters have more "memory" that they can memorize the unwanted characteristics that is specific to the training data. Another aspect is the distribution of the parameters, the convolutaional layers use the parameters more efficiently and make more of the increased number of parameters without easily falling into overfitting trap.
5. The depth of the architecture is also relevalant with the distribution of the parameters, how they organized in an architecture. The models with more depth are able to learn more complex features as well, however they are harder to train in terms of computation. The extremum of depth parameter results in not overfitting but underfitting. This is because the model is not able to learn the features of the data in depth and generalize well.



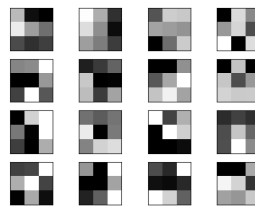
a) mlp_1



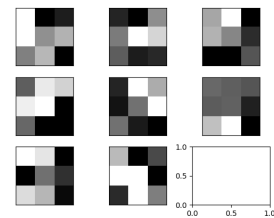
b) mlp_2



c) cnn_3



d) cnn_4



e) cnn_5

Avoid hyphenation at the end of a line. Symbols denoting vectors and matrices should be indicated in bold type. Scalar variable names should normally be expressed using italics. Weights and measures should be expressed in SI units. All non-standard abbreviations or symbols must be defined when first mentioned, or a glossary provided.

3.3. File naming and delivery

Please title your files in this order ‘procedia acronym_conference acronym_authorslastname’. Submit both the source file and the PDF to the Guest Editor.

Artwork filenames should comply with the syntax “aabbbbbb.ccc”, where:

- a = artwork component type
- b = manuscript reference code
- c = standard file extension

Component types:

- gr = figure
- pl = plate
- sc = scheme
- fx = fixed graphic

3.4. Footnotes

Footnotes should be avoided if possible. Necessary footnotes should be denoted in the text by consecutive superscript letters¹. The footnotes should be typed single spaced, and in smaller type size (8 pt), at the foot of the page in

¹ Footnote text.

Fig. 3. (a) first picture; (b) second picture.

which they are mentioned, and separated from the main text by a one line space extending at the foot of the column. The ‘Els-footnote’ style is available in the “TeX Template” for the text of the footnote.

Please do not change the margins of the template as this can result in the footnote falling outside printing range.

4. Illustrations

All figures should be numbered with Arabic numerals (1,2,3,...). Every figure should have a caption. All photographs, schemas, graphs and diagrams are to be referred to as figures. Line drawings should be good quality scans or true electronic output. Low-quality scans are not acceptable. Figures must be embedded into the text and not supplied separately. In MS word input the figures must be properly coded. Preferred format of figures are PNG, JPEG, GIF etc. Lettering and symbols should be clearly defined either in the caption or in a legend provided as part of the figure. Figures should be placed at the top or bottom of a page wherever possible, as close as possible to the first reference to them in the paper. Please ensure that all the figures are of 300 DPI resolutions as this will facilitate good output.

The figure number and caption should be typed below the illustration in 8 pt and left justified [**Note:** one-line captions of length less than column width (or full typesetting width or oblong) centered]. For more guidelines and information to help you submit high quality artwork please visit: <http://www.elsevier.com/artworkinstructions> Artwork has no text along the side of it in the main body of the text. However, if two images fit next to each other, these may be placed next to each other to save space. For example, see Fig. 1.

5. Equations

Equations and formulae should be typed in MathType, and numbered consecutively with Arabic numerals in parentheses on the right hand side of the page (if referred to explicitly in the text). They should also be separated from the surrounding text by one space

$$X_r = \dot{Q}_{rad}'' / (\dot{Q}_{rad}'' + \dot{Q}_{conv}'')$$

$$\rho = \frac{\vec{E}}{J_c(T = \text{const.}) \cdot \left(P \cdot \left(\frac{\vec{E}}{E_c} \right)^m + (1 - P) \right)} \quad (1)$$

6. Online license transfer

All authors are required to complete the Procedia exclusive license transfer agreement before the article can be published, which they can do online. This transfer agreement enables Elsevier to protect the copyrighted material for the authors, but does not relinquish the authors’ proprietary rights. The copyright transfer covers the exclusive rights to reproduce and distribute the article, including reprints, photographic reproductions, microfilm or any other reproductions of similar nature and translations. Authors are responsible for obtaining from the copyright holder, the permission to reproduce any figures for which copyright exists.

Acknowledgements

Acknowledgements and Reference heading should be left justified, bold, with the first letter capitalized but have no numbers. Text below continues as normal.

Appendix A. An example appendix

Authors including an appendix section should do so before References section. Multiple appendices should all have headings in the style used above. They will automatically be ordered A, B, C etc.

A.1. Example of a sub-heading within an appendix

There is also the option to include a subheading within the Appendix if you wish.

References

- [1] Filippini, Massimo, and Lester C. Hunt. (2011) “Energy demand and energy efficiency in the OECD countries: a stochastic demand frontier approach.” *Energy Journal* **32** (2): 59–80.
- [2] Filippini, Massimo, and Lester C. Hunt. (2012) “US residential energy demand and energy efficiency: A stochastic demand frontier approach.” *Energy Economics* **34** (5): 1484–1491.
- [3] Weyman-Jones, Thomas, Júlia Mendonça Boucinha, and Catarina Feteira Inácio. (2015) “Measuring electric energy efficiency in Portuguese households: a tool for energy policy.” *Management of Environmental Quality: An International Journal* **26** (3): 407–422.
- [4] Saunders, Harry (2009) “Theoretical Foundations of the Rebound Effect”, in Joanne Evans and Lester Hunt (eds) *International Handbook on the Economics of Energy*, Cheltenham, Edward Elgar
- [5] Sorrell, Steve (2009) “The Rebound Effect: definition and estimation”, in Joanne Evans and Lester Hunt (eds) *International Handbook on the Economics of Energy*, Cheltenham, Edward Elgar

Instructions to Authors for LaTeX template:

1. ZIP mode for LaTeX template:

The zip package is created as per the guide lines present on the URL <http://www.elsevier.com/author-schemas/preparing-crc-journal-articles-with-latex> for creating the LaTeX zip file of Procedia LaTeX template. The zip generally contains the following files:

- ecrc.sty
- elsarticle.cls
- elsdoc.pdf
- .bst file
- Manuscript templates for use with these bibliographic styles
- Generic and journal specific logos, etc.

The LaTeX package is the main LaTeX template. All LaTeX support files are required for LaTeX pdf generation from the LaTeX template package.

Reference style .bst file used for collaboration support: In the LaTeX template packages of all Procedia titles a new “.bst” file is used which supports collaborations downloaded from the path <http://www.elsevier.com/author-schemas/the-elsarticle-latex-document-class>

2. Reference style used in Computer Science:

| Title | Reference style |
|-------|----------------------|
| PROCS | 3 Vancouver Numbered |