

Uppsala University
Department of Information Technology
Division of Systems and Control
NW 2019-02
Last rev. October 15, 2021 by DW

Statistical Machine Learning

Instruction to the laboratory work

Image classification with Deep Learning

Language: Python

Preparation:

Read Section 2 carefully

Solve all preparatory exercises in Section 3

Submit your solutions on Studium (check Studium for the deadline)

Reading instructions:

- SML book: Chapter 6 and Section 5.4.
- This lab-pm: Chapter 1-3.

Name	Assistant's comments	
Program	Year of reg.	
Date		
Passed prep. ex.	Sign	
Passed lab.	Sign	

Contents

1	Introduction	1
2	PyTorch	2
2.1	Installation	2
2.2	Introduction	2
3	Preparation exercises	3
3.1	Softmax and cross-entropy	3
3.2	Dense neural network	4
3.3	Convolutional neural network	5
4	Laboratory exercises	7
4.1	Classification of hand-written digits	7
4.1.1	Download and run the code sample	9
4.1.2	Understand the model	10
4.1.3	Understand the training	11
4.1.4	Train a two-layer neural network	12
4.1.5	Go even deeper!	13
4.1.6	Train longer, be aware of numerical issues!	14
4.1.7	Use convolutional neural networks	15
4.1.8	Improve learning rate* (extras, not mandatory)	18
4.1.9	Regularize with dropout* (extras, not mandatory)	19
4.2	Real world image classification	20

1 Introduction

Deep Learning is a sub-field within machine learning that models high level abstractions of *high-dimensional* and *complex* data. One example of such complex and high-dimensional data is images. Consider a classification model that takes an image as input and the class that this image belongs to (for example “cat”, “dog” or “bike”) as output. Assume further that the image has $100 \times 100 = 10\,000$ pixels where each pixel in the image is considered as one input variable. This makes the problem very high-dimensional (even for such a small image). Further, the relation between these pixel values and the actual class is far from obvious, which makes the data very complex.

Deep neural networks, or just *deep learning*, have in the last 5 to 10 years proven to be very successful to model such complex and high-dimensional data sets. In this lab we will look at a smaller image classification task. We will learn a neural network to read hand-written digits with up to 98% prediction accuracy. In the very end of the lab we will load a pre-trained network that has been trained on a much larger data set with 1 000 different classes.

This laboratory work is based on Lecture 8 and 9 together with Chapter 6 in the SML book. Therefore, it is advisable to have the material from those lectures fresh in mind before starting this laboratory work.

The goal of this laboratory work is to:

- Learn how to build and train a neural network
- Learn how to improve the neural network model and its training.
- Application 1: Learn how to classify hand-written digits using neural network.
- Application 2: See how a state-of-the-art deep neural network performs at classifying real world images
- Get a glimpse of a state-of-the-art software library (PyTorch) for deep learning.

Throughout the lab we will use a software library called *PyTorch*. This library is introduced in Section 2. Section 3 contains the preparatory exercises, and Section 4 contains the exercises that you do during the 4h lab session.

Important: Read Chapter 2, solve the preparatory exercises, and submit your solutions on Studium *before* the lab session. Check Studium for the deadline.

2 PyTorch

PyTorch is an open source software library for machine learning that is based on the machine learning library Torch which is no longer actively developed. PyTorch is developed primarily by Facebook's artificial intelligence research group and used by companies as well as academic research groups. It can be used for general computations with multidimensional arrays on CPUs and GPUs, but it is tailored especially to deep learning and neural networks.

PyTorch is natively written in Python and C++, and well documented APIs exist for these languages. It is not the only deep learning framework, some other state-of-the-art alternatives are TensorFlow and MXNet.

2.1 Installation

PyTorch is already installed on the Linux systems in the computer rooms where the laboratory session is scheduled. You can either use these computers during the lab or bring your own computer.

If you choose to use your own computer, you need to have PyTorch properly installed *before* the lab. The lab assistants will not be able to assist you with the installation process during the lab. Please consult the PyTorch documentation for more information about the installation procedure.

Another option is to use Google Colab to work with PyTorch online. This cloud platform also optionally provides access to GPUs which might speed up some computations.

2.2 Introduction

A Jupyter notebook with an introduction to PyTorch can be downloaded from [here](#). Alternatively, you can open it on Google Colab. Reading and running the notebook is highly recommended, since it introduces important concepts and commands that are required in the lab session.

The official PyTorch tutorials and the PyTorch documentation might be helpful additional resources, in particular if you are looking for a more general introduction to PyTorch. However, the introduction to PyTorch in the provided Jupyter notebook covers everything you should know for the lab session.

3 Preparation exercises

3.1 Softmax and cross-entropy

In Lecture 3 the logistic regression model was introduced for problems with two classes. The class-1 probability $p(y = 1|\mathbf{x})$ was modeled as

$$p(y = 1|\mathbf{x}) = h(z), \quad \text{where} \quad z = \theta_0 + \sum_{j=1}^p x_j \theta_j, \quad \text{and} \quad h(z) = \frac{e^z}{1 + e^z}. \quad (1)$$

Question 3.1: Assume that we have estimated the parameters in (1) $\hat{\theta}_0 \dots \hat{\theta}_p$ using logistic regression from a set of training data. After training, we compute $z = \hat{\theta}_0 + \sum_{j=1}^p x_j \hat{\theta}_j = 1.0$ for a certain test input $\mathbf{x} = [x_1, \dots, x_p]^\top$. What are the probabilities according to the logistic regression model that the corresponding output y belongs to either of the two classes -1 or 1 , i.e., what is $p(y = -1|\mathbf{x})$ and $p(y = 1|\mathbf{x})$?

Answer:

In this lab we will consider classification problems with $M > 2$ classes. For each input \mathbf{x} we define M different *class probabilities* $p(y = 1|\mathbf{x}), \dots, p(y = M|\mathbf{x})$, which are the probabilities that \mathbf{x} belongs to either of the M classes.

We also extend the logistic function in (1) to a function that has these M class probabilities as outputs. We use the *softmax function* mapping $\mathbf{z} = [z_1, \dots, z_M]^\top$ onto M probabilities, i.e., $\mathbb{R}^M \mapsto [0, 1]^M$. The softmax function is defined as

$$\text{softmax}(\mathbf{z}) = \frac{1}{\sum_{l=1}^M e^{z_l}} [e^{z_1} \quad \dots \quad e^{z_M}]^\top. \quad (2)$$

Now we are ready to extend logistic regression model (1) to multiple classes $M > 2$. This extension is defined as

$$p(y = m|\mathbf{x}) = [\text{softmax}(\mathbf{z})]_m, \quad \text{where} \quad z_m = \theta_{m0} + \sum_{j=1}^p \theta_{mj} x_j \quad (3)$$

and where $[\text{softmax}(\mathbf{z})]_m = \frac{e^{z_m}}{\sum_{l=1}^M e^{z_l}}$ is the m th output from the softmax function. In contrast to logistic regression (1), we now have a set of parameters $\theta_{m0}, \dots, \theta_{mp}$ for each class m .

Question 3.2: Consider a problem with three classes $\{1, 2, 3\}$ where we have estimated all parameters $\hat{\theta}_{mj}$ from a training data set using the softmax model

in (3). For a certain test input \mathbf{x} we compute $\mathbf{z} = [z_1, z_2, z_3]^\top$ where $z_1 = 0$, $z_2 = -1$, and $z_3 = 1$. What are the class probabilities for the three different classes, i.e., what are $p(y = 1|\mathbf{x})$, $p(y = 2|\mathbf{x})$, and $p(y = 3|\mathbf{x})$? Which class has the highest probability?

Answer:

Suppose we want to learn the parameters θ based on a training data set $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$. Instead of letting the output \mathbf{y}_i of a data point i be an integer in $\{1, \dots, M\}$, we represent the m th class with a vector $\mathbf{y}_i = [y_{i1} \dots y_{iM}]^\top$, where $y_{ij} = 1$ if $j = m$, and $y_{ij} = 0$ otherwise. This is also known as the *one-hot encoding*. For example, if we have the three classes $\{1, 2, 3\}$ we encode class $m = 1$ as $\mathbf{y}_i = [1 \ 0 \ 0]^\top$, $m = 2$ as $\mathbf{y}_i = [0 \ 1 \ 0]^\top$, and $m = 3$ as $\mathbf{y}_i = [0 \ 0 \ 1]^\top$.

As a measure of fit between true output $\mathbf{y}_i = [y_{i1}, \dots, y_{iM}]^\top$ and the class probabilities $p(y_i = 1|\mathbf{x}_i), \dots, p(y_i = M|\mathbf{x}_i)$ we use the *cross-entropy loss function*

$$L(\mathbf{x}_i, \mathbf{y}_i, \theta) = - \sum_{m=1}^M y_{im} \ln g_{im}, \quad (4)$$

where we use the notation $g_{im} = p(y_i = m|\mathbf{x}_i)$.

Question 3.3: Consider the class probabilities $p(y = 1|\mathbf{x})$, $p(y = 2|\mathbf{x})$, and $p(y = 3|\mathbf{x})$ you got from the previous exercise. Compute the cross-entropy loss $L(\mathbf{x}, \mathbf{y}, \hat{\theta})$ between these probabilities and $\mathbf{y}^{(1)} = [1 \ 0 \ 0]^\top$, $\mathbf{y}^{(2)} = [0 \ 1 \ 0]^\top$, and $\mathbf{y}^{(3)} = [0 \ 0 \ 1]^\top$, respectively. Which one has the lowest cross-entropy?

Answer:

3.2 Dense neural network

Consider a classification problem where the input consists of $p = 144$ input variables $\mathbf{x} = [x_1, \dots, x_p]^\top$ and the output belongs to four classes $y \in \{1, \dots, 4\}$. We want to model the class probabilities $p(y = m|\mathbf{x})$ with a neural network with two dense layers.

$$\mathbf{q} = h \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right), \quad (5a)$$

$$\mathbf{z} = \mathbf{W}^{(2)} \mathbf{q} + \mathbf{b}^{(2)}, \quad (5b)$$

$$\begin{aligned} p(y = 1 | \mathbf{x}) &= [\text{softmax}(\mathbf{z})]_1, \\ &\vdots \\ p(y = M | \mathbf{x}) &= [\text{softmax}(\mathbf{z})]_M. \end{aligned} \quad (5c)$$

The hidden layer $\mathbf{q} = [q_1, \dots, q_U]^T$ has $U = 30$ hidden units. The weight matrices $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ in each of the two dense layers have the size (output units \times input units) and each offset vector $\mathbf{b}^{(1)}$ and $\mathbf{b}^{(2)}$ has the size (output units).

Question 3.4: *What are the sizes of the two weight matrices $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$ and the two offset vectors $\mathbf{b}^{(1)}$, $\mathbf{b}^{(2)}$ in the network above? How many parameters does the network have in total? (Each element in the weight matrices and the offset vectors contains one parameter.)*

Answer:

3.3 Convolutional neural network

Consider the same classification problem as in Section 3.2 but where the 144 input units represent 12×12 grayscale pixels in an image. We want to model this with a convolutional neural network (CNN). In a CNN, the hidden units in each layer are organized in *tensors*¹ of order 3 with the size (rows \times columns \times channels). Each grayscale image in our setting has the size $(12 \times 12 \times 1)$. Hence, each image has only one channel since each grayscale pixel can be represented with one scalar value corresponding to the brightness of that pixel².

The design of a CNN is that in each layer, the units from the previous hidden layer are convolved with a patch of weights, a so-called *filter*. We can have multiple filters (with different parameters) operate on the same input in parallel. Each filter then produces a new output channel for the next hidden layer. Each filter has the size (filter rows \times filter columns \times input channels) and if we stack all parameters in all filters in one weight tensor \mathbf{W} , that tensor has size (filter rows \times filter columns \times input channels \times output channels). Read more about CNNs in Section 6.3 in the SML book. Especially Figure 6.10 in the book might help you to answer the

¹A tensor is a generalization of a matrix to more than two different dimensions. This is also what has given Tensorflow its name since the library is used to calculate and operate on tensors.

²For a color image we would have three channels representing the three RGB colors red, green and blue. Color images are not considered further here.

questions below.

Consider a CNN with two convolutional layers parameterized with $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}$ and two dense layers parameterized with $\mathbf{W}^{(3)}, \mathbf{b}^{(3)}, \mathbf{W}^{(4)}, \mathbf{b}^{(4)}$. The first convolutional layer consists of 4 filters of size (filter rows \times filters columns) $= (5 \times 5)$ and we use the stride $[1,1]$ with zero-padding, i.e. the filter is moving by one step (both row- and column-wise) during the convolution such that the first hidden layer has the same number of rows and columns as the image. As previously stated, each grayscale image in has the size $(12 \times 12 \times 1)$.

Question 3.5: What are the sizes of the weight tensor $\mathbf{W}^{(1)}$, offset vector $\mathbf{b}^{(1)}$, and the first hidden layer $\mathbf{Q}^{(1)}$?

Answer:

In the following convolutional layer we use 8 filters of size 3×3 and the stride $[2,2]$, i.e., the filter is moving by two steps (both row- and column-wise) during the convolution such that the second hidden layer has half as many rows and columns as the previous first hidden layer.

Question 3.6: What is the size of the weight tensor $\mathbf{W}^{(2)}$, offset vector $\mathbf{b}^{(2)}$, and the second hidden layer $\mathbf{Q}^{(2)}$?

Answer:

After the two convolutional layers we implement two dense layers, the first with 60 hidden units and the second with 4 output units before we end with a softmax function to produce the predicted class probabilities.

Question 3.7: What are the sizes of the weight tensors and offset vectors $\mathbf{W}^{(3)}, \mathbf{b}^{(3)}$ and $\mathbf{W}^{(4)}, \mathbf{b}^{(4)}$ belonging to the two dense layers?

Answer:

4 Laboratory exercises

This section contains the laboratory exercises to be executed during the laboratory session. The main lab exercise is to build and train a neural network to classify hand-written digits as presented in Section 4.1. The lab ends in Section 4.2 with an evaluation of a much bigger state-of-the-art network that has been trained on over a million images.

4.1 Classification of hand-written digits

In this part of the lab we will learn how to use a neural network to classify images³. We will consider the so called MNIST data set⁴, which is one of the most well studied data sets within machine learning and image processing.

The data set consist of 60 000 training data points and 10 000 test data points. Each data point consist of a grayscale image with 28×28 pixels of a handwritten digit. The digit has been size-normalized and centered within a fixed-sized image. Each image is also labeled with the digit (0,1,...,8, or 9) it is depicting. In Figure 1 a batch of 100 data points from this data set is displayed.

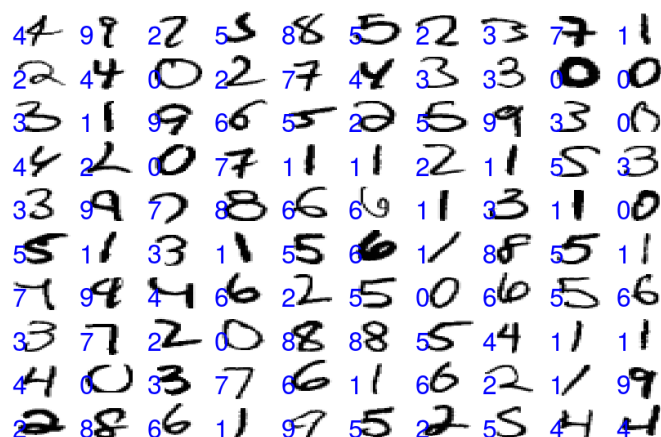


Figure 1: Some samples from the data we will use in the entire Section 4.1. The input is the pixels values of an image (black and white), and the output is the label of the digit it represents (blue).

In this classification task we consider the image as our input $\mathbf{x} = [x_1, \dots, x_p]^T$. Each input variable x_j corresponds to a pixel in the image. In total we have $p = 28 \times 28 = 784$ input variables.

The value of each x_j represents the color of that pixel. The color-value is within the interval $[0,1]$, where $x_j = 0$ corresponds to a black pixel and $x_j = 1$ to a white

³This MNIST lab is inspired by a similar one in the crash-course *Learn TensorFlow and deep learning, without a Ph.D.*

⁴You can find more information about this data set on Wikipedia.

pixel. Anything between 0 and 1 is a gray pixel with corresponding intensity.

We have in total 10 classes representing the 10 digits. We will use the so called one-hot encoding for the output described in Section 3. This means that the $\mathbf{y} = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$ represents the digit “0”, $\mathbf{y} = [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$ represents the digit “1”, and so forth. Consequently, the output \mathbf{y} is of dimension 10.

Based on a set of training data $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ with images and labels, the problem is to find a good model for the class probabilities

$$p(y = m|\mathbf{x}), \quad m = 1, \dots, 10, \quad (6)$$

i.e. the probabilities that an unseen image \mathbf{x} belongs to each of the 10 classes.

4.1.1 Download and run the code sample

Download the Jupyter notebook `mnist_onelayer.ipynb` from [here](#), launch Jupyter, and open the notebook. Alternatively, you can open it on Google Colab.

Task 4.1 Run the code as it is. When the training is done after a while, three figures appear, see Figure 2. ○

Tip 4.1 If you use Google Colab, you can run your code on GPUs by going to Runtime → Change runtime type and selecting GPU. It will give you a significant speedup, in particular when training deep neural networks and CNNs! ○

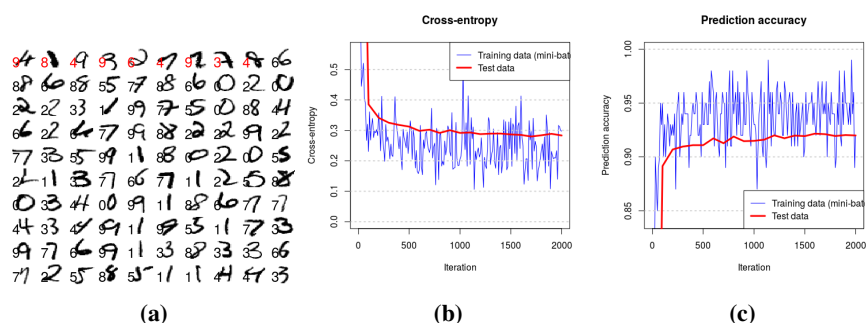


Figure 2: Three figures generated by the code. Figure 2a: Prediction performance on 100 randomly selected test points. Figure 2b: Cost function (cross-entropy) on test/training data. Figure 2c: Prediction accuracy on test/training data.

In Figure 2a, 100 randomly selected test images (out of the total 10 000 test images) are displayed together with their predicted label. For those images that have been incorrectly classified, the label is colored red.

To train the network we minimize a cost function that tells how bad we are at predicting the training data correctly. The cost function for this problem is the cross-entropy (more about that in Section 4.1.3). Figure 2b displays the cost function on test and training data with the iteration number on the x -axis.

In Figure 2c the prediction accuracy on test and training data⁵ is displayed. The prediction accuracy on test data is the performance measure that we are mostly interested in. The cross-entropy and prediction accuracy on both training data and test data as displayed on the figures are also printed in the terminal during training.

Question 4.1: What classification accuracy do you get on the test data?

Answer:

⁵The prediction accuracy on training data is evaluated on only 100 randomly selected training samples used during training (called mini-batch). That is why the training accuracy is so “noisy”.

4.1.2 Understand the model

The sample code is an implementation of a one-layer neural network with softmax transformation of the output. For a certain class m and data point i with $\mathbf{x}_i = [x_{i1}, \dots, x_{ip}]^T$ the model of the class probabilities is

$$p(y_i = m | \mathbf{x}_i; \boldsymbol{\theta}) = \frac{e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} \quad \text{where} \quad z_{im} = b_m + \sum_{j=1}^p x_{ij} W_{jm}, \quad (7)$$

where $\boldsymbol{\theta}$ is a vector with all the parameters W_{jm} and b_m in the model. With n training data points $\{\mathbf{x}_i, y_i\}_{i=1}^n$ and $m = 1, \dots, M$ classes we can write the model in matrix notation

$$\mathbf{G} = \text{softmax}(\mathbf{X}\mathbf{W} + \mathbf{b}), \quad (8)$$

where

$$\mathbf{G} = \begin{bmatrix} \mathbf{g}_1^T \\ \vdots \\ \mathbf{g}_n^T \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} W_{11} & \dots & W_{1M} \\ \vdots & & \vdots \\ W_{p1} & \dots & W_{pM} \end{bmatrix}, \quad \mathbf{b} = [b_1 \quad \dots \quad b_M],$$

with $\mathbf{g}_i^T = [p(y_i = 1 | \mathbf{x}_i; \boldsymbol{\theta}), \dots, p(y_i = M | \mathbf{x}_i; \boldsymbol{\theta})]$, and where \mathbf{W} and \mathbf{b} are the *weight matrix* and the *offset vector*, respectively. The offset vector \mathbf{b} is added to all n rows in (8). Also the softmax function is applied row by row as

$$\text{softmax}(\mathbf{z}_i^T) = \frac{1}{\sum_{l=1}^M e^{z_{il}}} [e^{z_{i1}} \quad \dots \quad e^{z_{iM}}], \quad \text{where} \quad \mathbf{z}_i^T = \mathbf{x}_i^T \mathbf{W} + \mathbf{b}. \quad (9)$$

Remember, the i th row in \mathbf{X} (i.e., \mathbf{x}_i^T) contains the 784 pixel values of image i and the i th row in \mathbf{G} (i.e., \mathbf{g}_i^T) contains the 10 probabilities that the i th image belongs to each of the 10 classes. The weight matrix \mathbf{W} and the offset vector \mathbf{b} contain all the parameters of the model. Note that \mathbf{W} and \mathbf{b} is the transposed version of the equivalent weight matrix and offset vector in (5b). We choose to do so in order to avoid transposing \mathbf{W} and \mathbf{b} in the code.

Task 4.2 Make sure you understand the model above. Read the code in the section **The model** and make sure that you can map the model presented above to what is implemented. If something is unclear, ask the lab supervisors! \circ

Question 4.2: How many parameters does this model have?

Answer:

4.1.3 Understand the training

To find good parameters, we need to train the model. This requires a cost function. The cost function describes the “distance” between the outputs $\mathbf{y}_i = [y_{i1} \dots y_{iM}]^T$ and the predicted class probabilities $p(y_i = 1 | \mathbf{x}_i; \boldsymbol{\theta}), \dots, p(y_i = M | \mathbf{x}_i; \boldsymbol{\theta})$. Using the notation g_{im} for each such predicted class probability, we can write the cross-entropy cost function

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \sum_{m=1}^M y_{im} \ln g_{im}, \quad (10)$$

where $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}\}$ are the parameters and where $g_{im} = p(y_i = m | \mathbf{x}_i; \boldsymbol{\theta})$. Remember, if \mathbf{y}_i encodes the hand-written digit “3”, all ten elements in \mathbf{y}_i will be zero except the fourth element, which will be one. Each g_{im} is a value between 0 and 1 which corresponds to the probability that the image i belongs to class m .

The training of the network is done by minimizing the cost function (10) in a loop. In each iteration, we compute the gradient $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ of the cost function with respect to the parameters and update the parameters by going a small step in the opposite direction of the gradient

$$\boldsymbol{\theta}^{(t+1)} := \boldsymbol{\theta}^{(t)} - \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}. \quad (11)$$

This is called gradient descent and γ is the learning rate. In each gradient step we do not use all $n = 60\,000$ training data to compute the cost function (10) and its gradient. Instead, we use 100 randomly selected training data points in each iteration. We call such a group of training data a *mini-batch*.

Next iteration we randomly selected 100 new training data points (of the ones which have not yet been selected) and compute the gradient based on this group of data. We continue until all training data points have been used. One such sweep through the training data is called an *epoch*. After one epoch is completed, we start the process over again. This training procedure is called *stochastic gradient*.

Task 4.3 Make sure you understand the training procedure described above. Read through the code in the section **The training**. Map the training procedure to what is written in the code. Ask the lab supervisors if anything is unclear! \circ

Question 4.3: How many iterations does it take until we have seen all training data points, i.e., how many iterations are included in each epoch? How many epochs do you train for in the code?

Answer:

4.1.4 Train a two-layer neural network

We will now add a second layer to the network. We get the model

$$\mathbf{Q} = h(\mathbf{XW}^{(1)} + \mathbf{b}^{(1)}), \quad (12a)$$

$$\mathbf{G} = \text{softmax}(\mathbf{QW}^{(2)} + \mathbf{b}^{(2)}), \quad (12b)$$

where \mathbf{Q} is an intermediate layer of hidden units. For this we need two weight matrices $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$ and two offset vectors $\mathbf{b}^{(1)}$, $\mathbf{b}^{(2)}$ ⁶.

```
U = 200 # number of hidden units
self.W1 = nn.Parameter(0.1 * torch.randn(784, U))
self.b1 = nn.Parameter(torch.zeros(U))
self.W2 = nn.Parameter(0.1 * torch.randn(U, 10))
self.b2 = nn.Parameter(torch.zeros(10))
```

We keep the softmax as activation function on the last layer, but use the logistic function, also called the *sigmoid function*⁷, after the first layer. The sigmoid function is defined as $h(x) = 1/(1 + e^{-x})$. The model will then be

```
Q = torch.sigmoid(X.mm(self.W1) + self.b1)
G = F.softmax(Q.mm(self.W2) + self.b2, dim=1)
```

When we grow the network deeper, it is important to initialize the weights randomly. In the lines above, each weight is initialized by a value drawn from a normal distribution with standard deviation 0.1, which works fine.

Task 4.4 Save `mnist_onelayer.ipynb` as a new file `mnist_twolayers.ipynb`. Add one more layer with 200 hidden units. ◦

Question 4.4: *What classification accuracy on test data do you get? Try some other numbers of hidden units (ranging from 10 to 750). How does it affect the performance? What happens if you initialize the weights with zeros as we did for the single layer neural network?*

Answer:

⁶For convenience, the weight matrices and offset vectors can be summarized in a `nn.Linear` object in PyTorch. In this computer lab, however, we implement all weight matrices and offset vectors explicitly.

⁷This is the same function as the logistic function (1). However, here it appears in a different context, hence the different name.

4.1.5 Go even deeper!

Now we can continue to go even deeper!

Task 4.5 Save `mnist_twolayers.ipynb` as a new file `mnist_fivelayers.ipynb`. Continue to add even more layers. Add in total five layers with 200, 100, 60, and 30 hidden units between each layer. ○

Question 4.5: *What classification accuracy on test data do you get?*

Answer:

When you go deeper there are some things to be aware of

- The sigmoid activation function causes issues in deep networks. It squashes all inputs into $[0,1]$ and when doing so repeatedly in multiple layers the gradients with respect to the parameters in the deepest layers might get close to zero. Instead, you might consider using the *Rectified Linear Unit* (ReLU) activation function $h(x) = \max(0, x)$. To use ReLU, simply replace all `torch.sigmoid` in the code with `F.relu`.
- Initialize all weights randomly! If you have not yet done so, do that now! For the offset vectors, when using ReLU:s, the best practice is to initialize them to a small positive values such that it operates in the non-negative range of ReLU. We can initialize all offset parameters with 0.1 as

```
b5 = nn.Parameter(torch.ones(10)/10)
```

- In problems where we have a lot of parameters (how many do we have now?) we also have many so-called “saddle points”. The gradient descent algorithm tends to get stuck at such saddle points. A better optimizer for dealing with this problem is the Adam-optimizer. To train with that routine, simply replace `optim.SGD` with `optim.Adam`. When swapping optimization routine also change the learning rate from 0.5 to 0.003.

Task 4.6 Improve the training of the network by using the tips above. Does this improve the prediction accuracy? ○

Task 4.7 Is the model trained after 2 000 iterations or do you think you can get an even better classification accuracy if you train longer? Try to train for 10 000 iterations. ○

By training this deeper model longer, you might suddenly get very poor performance due to numerical issues. Look on the next page why this happens and how to resolve the problem.

4.1.6 Train longer, be aware of numerical issues!

If you train the network longer, some elements in \mathbf{g}_i start getting very close to zero (meaning that the network is very certain that image i does not belong to that class). This is a problem in (10) since $\ln(0)$ will then output NaN (not-a-number), so also the whole cross-entropy, and the network will not be able to train anything.

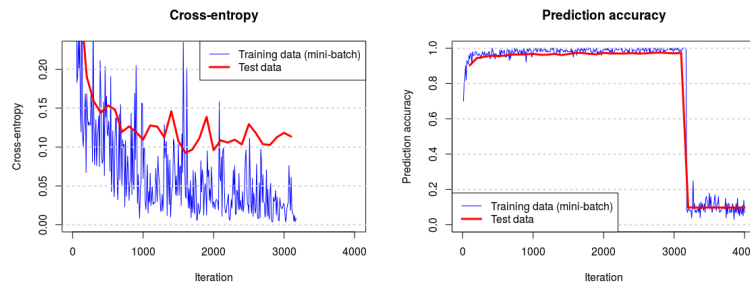


Figure 3: If some elements in \mathbf{g}_i get very close to zero, the logarithm function outputs not-a-number and the performance crashes.

To circumvent this numerical issue, the softmax (3) and the cross-entropy (10) can be computed in one step. Since the softmax (3) is an exponential (with normalization) no logarithm of a small number is needed since $\ln(e^t) = t$ even if, say, $t = -100$ and $e^{-100} \approx 0$.

In PyTorch, softmax and cross-entropy can be computed in one step with `F.cross_entropy`. First we have to replace the lines

```
G = F.softmax(Q4.mm(self.W5) + self.b5, dim=1)
return G
```

with

```
Z = Q4.mm(self.W5) + self.b5
return Z
```

to obtain the so-called logits Z , and then we use `F.cross_entropy` instead of `crossentropy` to compute the cross-entropy based on these logits.

Task 4.8 Implement the suggested change above to make the code more robust. Try again to train for 10 000 iterations. ○

Question 4.6: *What classification accuracy on test data do you get?*

Answer:

4.1.7 Use convolutional neural networks

All models that we considered so far started with putting all the 28×28 pixels in each image into a long vector with 784 elements. This partially destroys the spatial information present in the images. In contrast, a convolutional neural network (CNN) exploits this information, which enables us to achieve an even better classification performance. With this adjustment you should be able to reach approximately 98.5% prediction accuracy!

We will use a CNN with three convolutional layers and two final dense layers. The settings for the three convolutional layers are given in Table 1 and a graphical illustration of the whole network in Figure 4.

Table 1: Architecture of the three convolutional layers.

	Layer 1	Layer 2	Layer 3
Number of filters/output channels	4	8	12
Filter rows and columns	(5×5)	(5×5)	(4×4)
Stride	[1,1]	[2,2]	[2,2]
Padding	[2,2]	[2,2]	[1,1]

In this course, we use the convention that a weight tensor W in a convolutional layer has the size (filter rows \times filter columns \times input channels \times output channels). PyTorch, however, demands weight tensors of the size (output channels \times input channels \times filter rows \times filter columns). Hence the weight tensor and offset vector for the first convolutional layer are implemented as

```
U1 = 4
self.W1 = nn.Parameter(0.1 * torch.randn(U1, 1, 5, 5))
self.b1 = nn.Parameter(torch.ones(U1)/10)
```

The corresponding model implementation for that convolutional layer is

```
Q1 = F.relu(F.conv2d(X, self.W1, bias=self.b1, stride=1,
padding=2))
```

Note that we use the *non-vectorized* version of the images and do not reshape the input dimensions into a vector!

We use zero-padding such that we get equally many columns and rows for $Q1$ as we have for X , and hence specify `padding=2`. The correct padding for the other convolutional layers can be found in Table 1.

Continue reading on the next page!

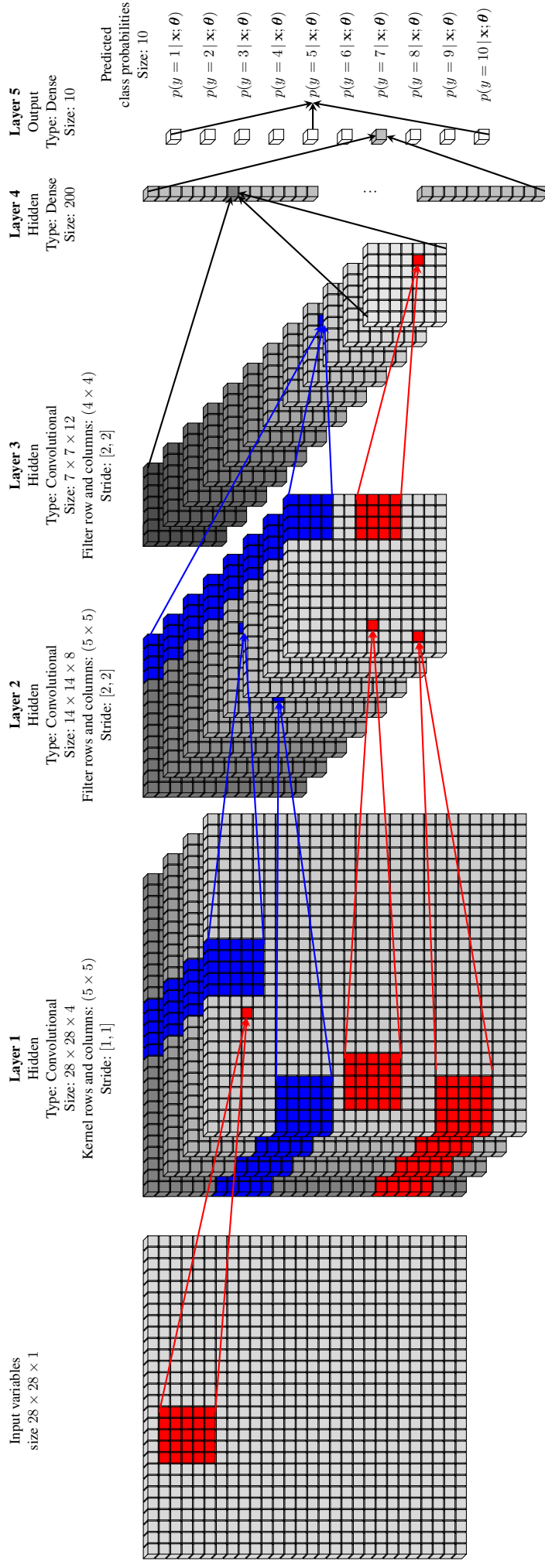


Figure 4: The architecture of the convolutional neural network in the lab. It consists of three convolutional layers and two dense layers. The leftmost square is the input image where each cube represents a pixel. The following three hidden convolutional layers organized in rows, columns and channels, where each cube represents a pixel. A certain hidden unit in one of these hidden layers depends on a patch of rows and columns and all channels from the previous layer, here illustrated as colored regions. Each hidden unit is multiplied with parameters in a filter. The same filter is used for all hidden units in a certain channel. A different filter is used to produce a different channel, here illustrated with the different colors red and blue. The filter size is 5×5 , 5×5 and 4×4 for the three convolutional layers. The network ends with two dense layers and a final softmax function to produce class probabilities as output.

After the three convolutional layers, we use two dense layers. Before we can apply the first dense layer, all hidden units in the third convolutional layer needs to be flattened into a long vector. This can be done with the command

```
Q3flat = Q3.view(-1, U3flat)
```

where the tensor `Q3` holds the hidden units in the third convolutional layer and has the size (batch size \times channels \times rows \times columns). The command reshapes `Q3` into a tensor `Q3flat` of size (batch size \times `U3flat`). The vectorized `Q3flat` is then the input to the first dense layer.

Question 4.7: *How many hidden units are there in total in the third convolutional layer, i.e. what is `U3flat` supposed to be in the code above? Hint 1: We have applied a stride of 2 twice were each of them reduced the number of rows and columns in the original image by a factor of 2. Hint 2: Have a look at Figure 4.*

Answer:

Finally, the first dense layer (layer 4) should have 200 hidden units.

Task 4.9 Save `mnist_fivelayers.ipynb` as a new file `mnist_CNN.ipynb`. Replace the first three dense layers in the previous code with three convolutional layers using the settings in Table 1 for each of these three layers. Add the reshaping command according to what is stated above. Update the forth and fifth layer according to the instructions above. Train for at least 4 000 iterations. Note that the training of the CNN takes a bit longer, be patient! ○

Tip 4.2 If you use Google Colab, you can get a significant speedup by running your code on GPUs (go to Runtime \rightarrow Change runtime type and select GPU). ○

Question 4.8: *What prediction accuracy do you achieve?*

Answer:

The next two pages contain tasks that are not mandatory (but, of course, interesting!). The mandatory part of the lab continues in Section 4.2 on page 20.

4.1.8 Improve learning rate* (extras, not mandatory)

In the end of the training, the learning rate $\gamma = 0.005$ is really too fast. Neither the prediction accuracy nor the cross-entropy on test data really converges and we will not get down to the best minimum of the cost function. You would therefore prefer having a lot smaller γ , but with a very small γ the training takes too long. One solution is to start learning fast (to get approximately close to the minimum) and then slow down. Consider adjusting the learning rate γ as

$$\gamma^{(t)} = \gamma_{\min} + (\gamma_{\max} - \gamma_{\min})e^{-\frac{t}{2000}}, \quad (13)$$

where $\gamma_{\max} = 0.003$, $\gamma_{\min} = 0.0001$, and t being the iteration number. This means that we start with a learning rate of $\gamma_{\max} = 0.003$ and approach $\gamma_{\min} = 0.0001$ as $t \rightarrow \infty$.

You can update the learning rate of a PyTorch `optimizer` by running

```
for p in optimizer.param_groups:
    p['lr'] = gamma_new
```

where `gamma_new` is the new learning rate⁸.

Task (optional) 4.10 Improve the training by adjusting the learning rate according to the formula above or any other formula that you come up with! Train for at least 6 000 iterations. You should now be able to push the prediction accuracy on the test data above 99%! ◦

Question 4.9: *Do you manage to get 99% prediction accuracy on test data?*

Answer:

Question 4.10: *At this point, the network can start overfitting on training data. How can you see that? What could you do to avoid this overfitting from happening?*

Answer:

⁸PyTorch itself also provides different so-called learning rate schedulers that can be used to adjust the learning rate.

4.1.9 Regularize with dropout* (extras, not mandatory)

When we extend our network and start seeing signs of over-fitting we know that there is room for improvement! One way to avoid the over-fitting is to reduce the size of the network. However, in practice a better strategy is to extend it a bit more and add regularization to the network. In Lecture 9 we talk about a popular regularization method for neural networks called *dropout*.

In dropout we remove at each training iteration a random selection of hidden units together with its incoming and outgoing links. We then train the remaining part of the network as if the removed units were not present. At each training iteration a new random selection of hidden units are removed. During test time all hidden units are used but their outputs are scaled with the probability that they were present during training. Read more about it in Section 6.4 in the SML book.

Here we choose to add dropout to the last hidden layer, since this will affect the weights in the first dense layer where most of the weights in our network reside. This is implemented by adding a dropout layer to our network

```
self.dropout = nn.Dropout(p=pzero)
```

and applying it to the last hidden unit

```
Q4dropout = self.dropout(Q4)
```

where `pzero` is the probability that a hidden unit in that hidden layer is zeroed. During training we keep 75% of the hidden units on average but during testing we want to keep all hidden units. Therefore we activate the dropout layer during training by calling `net.train()` and deactivate it during testing by calling `net.eval()`.

Task (optional) 4.11 Regularize the network by implementing dropout according to the description above. Train for 6 000 iterations. ○

Question 4.11: *What classification performance do you achieve? Does it seem like you are doing less overfitting?*

Answer:

Task (optional) 4.12 Play around with number of layers, channels, stride, filter size, dropout, learning rate etc. to achieve an even better classification performance. For example, extending the network to 6, 12, and 24 filters of sizes (6×6) , (5×5) , and (4×4) in the three convolutional layers might give an even better performance. If you overfit again, fight back with some more dropout. ○

4.2 Real world image classification

We have implemented and trained a five layer CNN. We also learned how to deal with several practical issues that appear when training a deeper network. Working with the MNIST data set, we went from $\approx 92\%$ prediction accuracy on test data up to more than $\approx 99.2\%$. This is quite close to the world record of 99.77%! See the full “leaderboard” on <http://yann.lecun.com/exdb/mnist/>.

The MNIST data set is a fairly small data set in a deep learning context. In this final lab exercise, we will use a network that has been trained on 1.2 million images provided by ImageNet⁹ used in the *Large Scale Visual Recognition Challenge* 2012-2014, see Figure 5 for a few training data examples. Each image is labeled (by hand!) with presence or absence of one of 1000 object categories¹⁰. Simonyan and Zisserman (2014) provided the winning contribution of the 2014 competition called VGG16. See Table 2 for a comparison between the model used previously in this lab trained on MNIST, and the VGG16 model trained on the ImageNet data.



Figure 5: 100 images from ImageNet

Table 2: Comparison between the MNIST classification problem in the lab with the VGG16 network trained on ImageNet.

	This lab	VGG16
Data set	MNIST	ImageNet (only a subset)
Training data size	60 000	1 200 000
Test data size	10 000	150 000
Nr of pixels	$28 \times 28 = 784$	$224 \times 224 = 50\,176$
Nr of image channels (colors)	1 (only gray-scale)	3 (RGB images)
Nr of class labels	10	1 000
Nr of layers	5	16
Nr of parameters	122 260	138 000 000
Training time	8 min ¹¹	2-3 weeks ¹²

To get good weights for such a network, a lot of training time and computational

⁹Website: <http://www.image-net.org/>.

¹⁰Class categories: <http://image-net.org/challenges/LSVRC/2014/browse-synsets>.

¹¹With 10 000 iterations on a laptop

¹²On a system equipped with four NVIDIA Titan Black GPUs

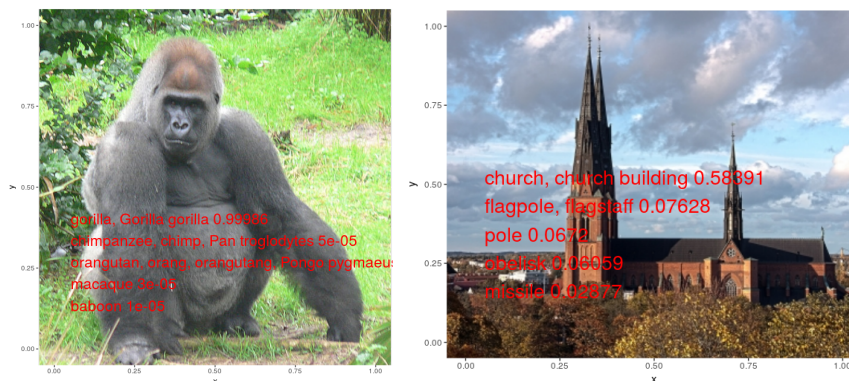


Figure 6: A few results from the VGG16 classifier trained on the ImageNet data set.

resources are required. However, when the network has been trained already, the class of a new unseen image can be predicted fast. Here, we will predict the class for images of our choice using the pre-trained network VGG16 presented above.

Task 4.13 Download the Jupyter notebook `VGG16_classification.ipynb` from here, launch Jupyter, and open the notebook. Alternatively, you can open it on Google Colab.

Look at the code (you don't have to understand the details). Note that the code will load pre-trained weights instead of training them on a training data set. Run the notebook and study the output. ○

Task 4.14 Analyze some images by changing the line where the image is loaded. How well does the model perform on your choice of images? Figure 6 displays some results that we got. ○

References

Simonyan, K., & Zisserman, A. (2014). *Very Deep Convolutional Networks for Large-Scale Image Recognition* (tech. rep.). arXiv:1409.1556.