

Assignment: Deep Learning

2ID90, edition 2017-2018 Q3

1 Overview

In this assignment you are going to gain practical experience with basic aspects of DEEP LEARNING, in particular with deep artificial neural networks. You will be given a library containing skeleton implementations for handling neural networks. You will gradually learn to use the library and extend it with additional functionality. Finally, you will design and tune a *convolutional neural network* for image classification. The tasks marked ♣ have to be reported on. Make sure that you describe HOW you realized a task, report on your RESULTS, and supply EVIDENCE (algorithmic details, formulas, tables, figures, etcetera).

2 Installation

A MAVEN project¹, called DeepLearning, is provided. It contains the source code for our deep learning library. To install the project

1. Download the assignment, see section 10;
2. Unzip it in a suitable folder of your own choice;
3. Open the project DeepLearning in your IDE;
4. Compile the project: this step will download several libraries from online maven repositories.

If the project compiles, your installation is completed successfully.

¹ This project has been developed and tested in NetBeans 8.2. It should also work in other IDE's that support maven. However, your mileage may vary.

3 Regression: Three valued function

Copy the following empty experiment to folder src/main/java/-experiments.

```
1 package experiments;
2 import nl.tue.s2id90.dl.experiment.Experiment;
3 import java.io.IOException;
4
5 public class FunctionExperiment extends Experiment {
6     // (hyper)parameters
7     // ...
8     public void go() throws IOException {
9         // you are going to add code here.
10    }
11    public static void main(String[] args) throws IOException {
12        new FunctionExperiment().go();
13    }
14 }
```

In a first step we are going to read some data. The data is going to be processed in so-called batches². Each batch consists of batchSize

² Sometimes also called mini-batches.

labeled input records. Add `batchSize`³ to the hyperparameters⁴ and add the following lines to the `go()` method.

```
// read input and print some information on the data
2 InputReader reader = GenerateFunctionData.THREE_VALUED_FUNCTION
  (batchSize);
System.out.println("Reader info:\n" + reader.toString());
```

³ `int batchSize = 32;`

⁴ Hyperparameters are parameters of the learning process that are not learned.

Use your IDE to fill in missing imports⁵. Running the resulting program should give you an output similar to the following⁶.

```
reader class      : GenerateFunctionData
2 batch size      : 16
#batches          : 62
4 #training pairs : 1000
#validation pairs : 100
6 input shape     : (1,1)
output shape      : (1,3)
8 headers         : [x, x, x^2, exp(x)]
```

⁵ In Netbeans, Right-Mouse-Click menu/Fix imports.

⁶ To print 10 sample records: `reader.getValidationData(10).forEach(System.out::println);`

From the printed headers you can see that each record has 4 fields: a feature vector with 1 field ("x") and three result fields ("x", "x^2", and "exp(x)"). The feature or input vector has `SHAPE (1,1)`: that is, it is one vector with 1 element; the output has `SHAPE (1,3)`: one vector with 3 elements⁷. In our program we can pick up these values as follows.

```
int inputs = reader.getInputShape().getNeuronCount();
2 int outputs = reader.getOutputShape().getNeuronCount();
```

⁷ See section 7 for more information about shapes.

Next we construct a neural network with inputs input neurons and outputs output neurons. We do so in a separate method:

```
Model createModel( int inputs, int outputs ) {
2   Model model = new Model(new InputLayer("In", new TensorShape(inputs), true));
   model.addLayer(new SimpleOutput("Out", new TensorShape(inputs), outputs, new MSE(), true));
4   return model;
}
```

This model⁸ has an input layer with *inputs* neurons and is FULLY CONNECTED to an output layer that has *inputs* connections coming in to each of its *outputs* neurons. The output layer has LOSS FUNCTION MSE, mean squared error. The last step in creation of a model consists of initializing its weights⁹. After that we have a fully initialized model, but we still need to train it. For training we need two additional hyperparameters¹⁰:

⁸ `System.out.println(model);` prints a summary of the model.

⁹ `model.initialize(new Gaussian());`

¹⁰ `int epochs=10; float learningRate = 0.01f;`

- **epochs:** The parameter *epochs* is the number of epochs that a training takes. In an epoch all the training samples are presented once to the neural network.

- `learningRate`: Parameter for the gradient descent optimization method.

Add the following code to your experiment's `go()` method to first create the SGD¹¹ optimizer and then call your experiment's training method.

```
// Training: create and configure SGD && train model
2 Optimizer sgd = SGD.builder()
    .model(model)
4   .validator(new Regression())
    .learningRate(learningRate)
6   .build();

8 trainModel(model, reader, sgd, epochs, o);
```

¹¹ SGD stands for Stochastic Gradient Descent, a method for optimization of the network weights. Its validator computes a validation of the network, which is used to indicate the performance of the network. The validator is used after training with a batch using the data in that batch, and after an epoch with the full validation set. Here, the validator class `Regression` just computes the mean squared error. So, lower validation here means higher performance.

After compiling and running¹² your experiment, its output should be similar to the following, ignoring the individual batches.

```
Validation after epoch 1: (batch: 31; validation: 0.36179352)
2 Validation after epoch 2: (batch: 62; validation: 0.19996504)
Validation after epoch 3: (batch: 93; validation: 0.13828962)
4 Validation after epoch 4: (batch: 124; validation: 0.13278994)
Validation after epoch 5: (batch: 155; validation: 0.15740058)
6 Validation after epoch 6: (batch: 186; validation: 0.19607270)
Validation after epoch 7: (batch: 217; validation: 0.23963349)
8 Validation after epoch 8: (batch: 248; validation: 0.28260258)
Validation after epoch 9: (batch: 279; validation: 0.32205325)
10 Validation after epoch 10: (batch: 310; validation: 0.35727578)
```

¹² To appreciate what the `trainingModel` method does, read its source code. In netbeans you can jump to the source code by `Ctrl-Left` click on the method name.

Clearly, this optimization process is not converging.

-  **Examine** the hyperparameters to see what their effect is.

Changing the topology of the network, may also make the network more powerful. To do so, add an additional layer to the network, as shown below, where the integers m and n are the number of incoming connections and the number of neurons of the layer, respectively. Adapt the rest of the network accordingly.

```
model.addLayer(new FullyConnected("fc1", new TensorShape(m), n, new RELU()));
```

-  **Examine** the new hyperparameter(s) to see what their effect is.

4 Classification: *fashion-MNIST*

We are now going to look at the fashion-MNIST dataset; see Figure 1. This dataset contains labeled grayscale images of 28x28 pixels. There are 60,000 training images and 10,000 validation images. The 10 labels are: 0: T-shirt/top, 1: Trouser, 2: Pullover, 3: Dress, 4: Coat, 5: Sandal, 6: Shirt, 7: Sneaker, 8: Bag, and 9: Ankle boot.

- Create a new experiment `ZalandoExperiment` with learning rate 0.01, batch size 32, and epochs 5.

- Read the Zalando data.

```
1 // read input and print some information on the data
  InputReader reader = MNISTReader.fashion(batchSize);
3 System.out.println("Reader info:\n" + reader.toString());
```

- Print the values of a record to get some idea of the data:

```
1 // print a record
  reader.getValidationData(1).forEach(System.out::println);
```

- Show a few images to get more acquainted with the dataset.

```
// add a field to your Experiment class
2 String[] labels= {
    "T-shirt/top", "Trouser", "Pullover", "Dress", "Coat",
4    "Sandal", "Shirt", "Sneaker", "Bag", "Ankle boot"
  };
6 ...
// add a constructor to your experiment that enables the gui
8 ZalandoExperiment() { super(true); }
10 ...
// add in the method go():
12 ShowCase showCase = new ShowCase(i -> labels[i]);
  FXGUI.getSingleton().addTab("show case", showCase.getNode());
14 showCase.setItems(reader.getValidationData(100));
```

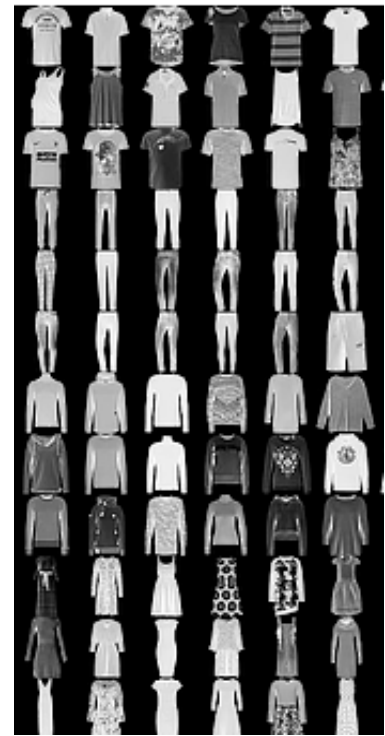


Figure 1: some images from the fashion MNIST dataset, made available by Zalando research.

- Build a fully connected network, having an input layer, a flatten layer, and an output layer. The flatten layer flattens the shape¹³ of an input image into a linear shape $(1, m)$ ¹⁴. For the output layer, use a softmax activation function that produces a probability distribution over the output classes. Use CrossEntropy as the loss function. See the implementation below.

```
// add flatten layer after input layer
2 model.addLayer(new Flatten("Flatten", new TensorShape(...));
// add output layer
4 model.addLayer(new OutputSoftmax("Out",
    new TensorShape(m), n, new CrossEntropy()));
6 );
```

- Finally, create a suitable optimizer; note, that we need a different validator¹⁵ than in the regression case.

```
Optimizer sgd = SGD.builder()
2   .model(model)
   .learningRate(learningRate)
4   .validator(new Classification())
   .build();
6 trainModel(model, reader, sgd, epochs, o);
```

¹³ See Section 7.

¹⁴ For the output layer, the integers m and n are the number of incoming connections and the number of neurons of the layer, respectively.

¹⁵ The Classification validator computes the fraction of correct predictions. So, for this validator a higher value is better.

- ♣ **Optimize** the hyperparameters learningRate and batchSize.
- ♣ To improve your results **implement** MEAN SUBTRACTION. See course notes CS231n part 2 for background on data preprocessing. See Section 8 for details and use the following code skeleton to realize this.

```

1 public class MeanSubtraction implements DataTransform {
2     Float mean;
3     @Override public void fit(List<TensorPair> data) {
4         if (data.isEmpty()) {
5             throw new IllegalArgumentException("Empty dataset");
6         }
7         for(TensorPair pair: data) {
8             ...
9         }
10        ...
11    }
12    @Override public void transform(List<TensorPair> data) {
13        // To do
14    }
15 }
    
```

- ♣ Finally, several variations of gradient descent exist. GRADIENT DESCENT WITH MOMENTUM is one of them; See course notes CS231n part 3 on SGD. **Implement** gradient descent with momentum. See Section 9 for details on adding your own gradient descent variant. For gradient descent with momentum you need to memorize the previous update during the gradient descent. This update is of the same dimension as your gradient. A safe way to construct one is as follows.

```

1 INDArrary update;
2 void update(...) {
3     // on the first call of this method, create update vector.
4     if (update==null) update=gradient.dup('f').assign(o);
5     ...
6 }
    
```

5 Classification: Square, Circle, or Triangle?

- ♣ Read the SCT dataset¹⁶ with 28x28 grayscale images containing either a square, a circle or a triangle. **Apply** your Zalando network to this dataset and **tune** your hyperparameters.

¹⁶ The SCT dataset is generated on the spot, if you choose big data sizes, its generation is slow.

```

1 int seed = 11081961, trainingDataSize=1200, testDataSize=200;
2 InputReader reader = new PrimitivesDataGenerator(batchSize,
3     seed, trainingDataSize, testDataSize);
    
```

- ♣ To improve the performance of the network, **design** a network using convolutional layers, possibly combined with pooling layers. Use the following layer constructors to create these layers.

```

/* Convolution2D layer constructor */
2 public Convolution2D(String layerName, TensorShape inputShape,
   int kernelSize, int noFilters, Activation activation
4 );

/* Pooling layer constructor */
6 public PoolMax2D(String layerName, TensorShape inputShape,
   int stride
8 );

```

Our convolutional layers have stride 1, and zero-padding such that width and height of output images equal those of the input images. Note that this limits the choice of the kernel size. For the pooling layer the stride has to be a divisor of the image size.

- **♣ Implement L2 weight decay.** See CS231n course notes. Make a custom UpdateFunction that combines a gradient descent update with L2 weight decay. Do not apply weight decay for the bias parameters. Use the following as a code skeleton for a weight decay class.

```

1 public class L2Decay implements UpdateFunction {
   float decay;
3   UpdateFunction f;
   public Decay(Supplier<UpdateFunction> supplier, float decay) {
5       this.decay = decay;
       this.f = supplier.get();
7   }

   @Override
9   public void update(...) {
11      ...
   }
13 }
...
15 // while building an optimizer
Optimizer sgd = SGD.builder() ...
17 .updateFunction(
   () -> new Decay(GradientDescentWithMomentum::new, 0.0001f)
19 ).build();

```

- **♣ Implement ADADELTA**, a gradient descent alternative. See resources, for a paper describing this method. The main advantage of AdaDelta is that it automatically chooses the learning rate.

6 Convolution: CIFAR10

- Read the Cifar10 dataset, a dataset containing 28x28 RGB images with 10 different labels: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck; see Figure 2. This dataset contains 50,000 training images and 10,000 test images.

```

1 Cifar10Reader reader = new Cifar10Reader( batchSize, 10);

```

The second parameter in the above constructor is the number of labels; so, in the above case the whole dataset is read, but if this parameter is 5, only images classified with the first 5 labels are read.

- Inspect the data.
- Normalize the data.
- ♣ **Design** a network for this dataset.
- ♣ **Tune**: parameters, #layers, layer sizes, etcetera; targeting 70+ % accuracy.
- ♣ **Optionally, implement**
 - new layers: Dropout; see CS231n notes;
 - Batch Normalization (hard);
 - ...

7 Appendix: Shape

A `SHAPE` is a tuple that describes how many numbers an array has in each dimension. For instance, shape $(3, 5)$ stands for a 2-dimensional array with 3 rows and 5 columns; see Figure 3.

When describing the input (or output) of a neural network there often is an additional 1 at the beginning of the shape. This is due to the fact that the network can actually handle, for efficiency reasons, multiple vectors simultaneously. So, for example if the input of a network has shape $(1, 3, 5)$ it actually accepts, for instance 2 vectors of shape $(3, 5)$ as input, but only when they are combined in an array of shape $(2, 3, 5)$; see Figure 4.

When defining the layers of a network we typically need to specify the shape of its input. For that we use the class `TensorShape`. We then write, for instance, `new TensorShape(10)` to specify that a layer has 10 inputs, or `new TensorShape(20, 25, 3)` to specify that the input is a 20×25 image with 3 layers, r, g, and b. There are two peculiarities with this class:

- For the constructor of the class `TensorShape`, the leading 1, as described in the previous paragraph, is left out. So, `new TensorShape(10)`, results in shape $(1, 10)$.
- The order of the arguments of the `TensorShape` class differs from the above introduced shapes. For example, `new TensorShape(20, 25, 3)`, gives shape $(1, 3, 20, 25)$.

8 Appendix: Data preprocessing

Data transformation in our deep learning library is done by implementing the `DataTransform` interface.

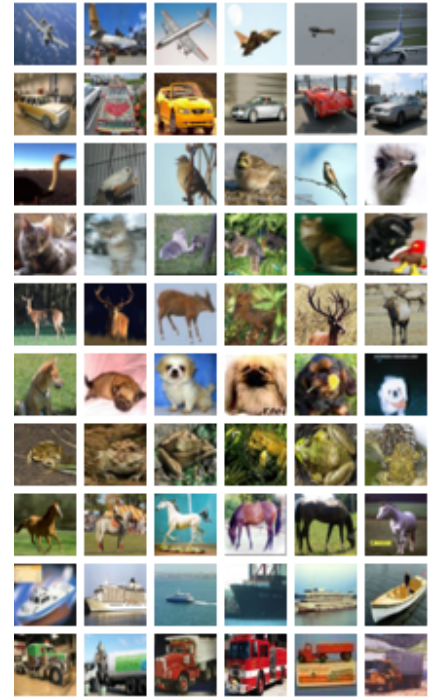


Figure 2: Some Cifar10 images.

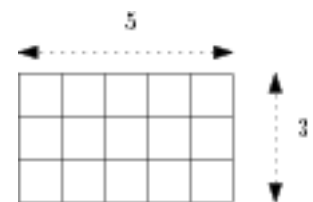


Figure 3: Array with shape $(3, 5)$.

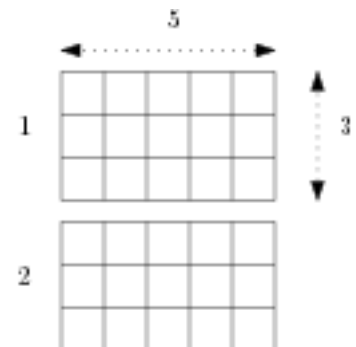


Figure 4: Array with shape $(2, 3, 5)$. Conceptually, the same as two rows, each containing a 3×5 matrix.


```

1 public interface DataTransform {
2     /** computes statistics for the dataset consisting of input-output pairs, these statistics
3      * are used in the transform method. Note that the stats are only based on the input.
4      * @param pairs dataset */
5     void fit(List<TensorPair> pairs);
6
7     /** transforms the dataset, using the statistics calculated by the fit method.
8      * @param pairs dataset */
9     void transform(List<TensorPair> pairs);
10 }
    
```

Data is then, for instance, transformed as follows:

```

1 DataTransform dt = new MyDataTransform();
2 dt.fit(myTrainingData);
3 dt.transform(myTrainingData);
4 dt.transform(myValidationData);
    
```

Each data element is a tensor pair, consisting of an input tensor and an output tensor. The data transformation only operates on the input tensor. The data¹⁷ inside a tensor is stored in an n-dimensional array `INDArray`. Such arrays and operations on it are implemented in a highly optimized library, called `N-DIMENSIONAL ARRAYS FOR JAVA`, or in short `ND4J`¹⁸.

There are typically two types of operations on an `INDArray` `nda`; those that return their result in a newly created array, and those that store their results 'in place' in `nda` itself. The latter operations typically have a suffix 'i'. For instance, `nda.addi(1)` adds 1 to each element of `nda`. The statement `nda.add(1)` on the other hand does not change `nda` at all, but returns a newly created array.

¹⁷ To get the `INDArray` in a tensor:
`Tensor t; ...; INDArray a = t.getValues();`

¹⁸ For a quick introduction see the user guide and syntax at <https://nd4j.org/>. For the full api of `ND4j` see its [javadoc](#).

9 Appendix: Gradient Descent variants

For your own gradient descent variant you need to create a class that implements the `UpdateFunction` interface and inform¹⁹ the `SGD` optimizer to use it. As a result the optimizer will automatically create two separate `UpdateFunction` objects for each layer in the network: One for updating the weights and one for updating the bias.

```

1 class MyGradientDescentVariant implements UpdateFunction { ... }
2 ....
3 Optimizer sgd = SGD.builder() ...
4     .updateFunction(MyGradientDescentVariant::new).build();
    
```

¹⁹ The optimizer actually needs a parameter-less method that can create an `UpdateFunction` object. Two ways to realize that are: 1) a reference to a constructor of a class, e.g. `MyGradientDescent::new`; or, a lambda function, e.g. `() -> new MyGradientDescentVariant()`.

Below the interface `UpdateFunction` is given.

```

1 public interface UpdateFunction {
2     /** A typical implementation of this interface does a gradient descent step, like
3      * array ← array - (learningRate/batchSize) * gradient.
4      * Other implementations may decide to ignore e.g. the learningRate.
5      * As a side effect the method update makes all components of the gradient vector zero.
6      * Typically, this is done by the call: gradient.assign(0).
7      * @param array array that is to be updated
8      * @param isBias true is array represents bias values, as opposed to weights.
9      * @param learningRate learning rate for gradient descent
10     * @param batchSize number of samples whose resulting gradients are accumulated in gradient
11     * @param gradient accumulated gradient */
12     void update(INDArray array, boolean isBias, float learningRate, int batchSize, INDArray gradient)
13 }
    
```


The default implementation of this interface is given below. The Nd4j method called there is a generalized vector addition: $\vec{v} \leftarrow \vec{v} + \mu \vec{w}$, implemented in a native BLAS²⁰ library.

²⁰ Basic Linear Algebra Subprograms

```

1 public class GradientDescent implements UpdateFunction {
2     /* Does a gradient descent step with factor 'minus learningRate' and corrected for batchSize. */
3     @Override public
4     void update(INDArray array, boolean isBias, float learningRate, int batchSize, INDArray gradient) {
5         float factor = -(learningRate/batchSize);
6         Nd4j.getBlasWrapper().level1().axpy( array.length(), factor, gradient, array );
7         // array ← array + factor * gradient
8         gradient.assign(0);
9     }
10 }

```

10 Appendix: Resources

- Assignment text and programming resources are available from the [2017-2018 2ID90 canvas website](#).
- Course notes: CS231n Convolutional Neural Networks for Visual Recognition [part 1](#), [part 2](#), [part 3](#); Andrej Karpathy, Stanford Computer Science.
- [ADADELTA: An Adaptive Learning Rate Method](#); Matthew D. Zeiler.