# Contents

# List of Figures

# 1 Project Description

## 1.1 Basics

We use the following notation:

- Matrices are written in uppercase bold e.g. **X**.

- Vectors are written in lowercase bold e..g **x**.

- scalars are written in lowercase or uppercase.  Lowercase indicates that it is a counting variable and uppercase that it is one of the limits in an finite set.

- For all functions , e.g. $f(x)$, where $x = $ **X** , we apply the function element wise.

- Dot product is indicated by simple concatenation of two matrices/vectors e.g. **Xy**.

- Element wise multiplication is indicated as $\mathbf{X} \otimes \mathbf{Z}$.

Neural networks are a statistical model for regression or classification.  Even though naturally they do not produce any classification output, we can model the regressed output to be used as a classifier.

Overall we have two different modes: Training and evaluation. In the training phase we use the back propagation algorithm to train the network and adjust it so that it produces output, which is close to our target output.

Furthermore a neural network contains $L$ hidden layers, which are called hidden since their output is not directly observable.

The weights are in our simple case randomly initialized and then updated using the back propagation rule.

A basic neural network has therefore the following parameters:

- The weights from an neuron $k$ in layer $a$ to a neuron $j$ in layer $b$

- The input to the network, usually a vector **x** with $k$ values

- The target which should be estimated, usually a vector **t** with $q$ values.

- The learning rate $\eta$, which is the indicator how fast the network learns.

Moreover we can add some advanced techniques as the momentum to speed up the training the momentum

### 1.1.1 Gradient Descent

Gradient descent is a first-order optimization algorithm. To find a local minimum of a function using gradient descent, it takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point. If instead one takes steps proportional to the positive of the gradient, one approaches a local maximum of that function; the procedure is then known as gradient ascent [Kiw01, Qia99].

Gradient descent is based on the observation that if the multivariable function $F(\mathbf{x})$ is defined and differentiable in a neighborhood of a point **a**, then $F(\mathbf{x})$ decreases fastest if one goes from **a** in the direction of the negative gradient of F at $\mathbf{a}, -\nabla F(\mathbf{a})$ [Yua99]. It follows that, if

$$\mathbf{b} = \mathbf{a} - \gamma \nabla F(\mathbf{a}) \tag{1.1}$$

for $\gamma$ small enough, then $F(\mathbf{a}) \geq F(\mathbf{b})$. With this observation in mind, one starts with a guess $\mathbf{x}_0$ for a local minimum of F, and considers the sequence $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \ldots$ such that

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma_n \nabla F(\mathbf{x}_n), \ n \geq 0. \tag{1.2}$$

Finally we have [Cau47]:

$$F(\mathbf{x}_0) \geq F(\mathbf{x}_1) \geq F(\mathbf{x}_2) \geq \cdots, \tag{1.3}$$

### 1.1.2 Minibatch Stochastic Gradient Descent

Both statistical estimation and machine learning consider the problem of minimizing an objective function that has the form of a sum:

$$Q(w) = \sum_{i=1}^{n} Q_i(w) \tag{1.4}$$

where the parameter $w^*$ which minimizes Q(w) is to be estimated. Each summand function $Q_i$ is typically associated with the $i$-th observation in the data set (used for training).

When used to minimize the above function, a standard (or "batch") gradient descent method would perform the following iterations :

$$w := w - \eta \nabla Q(w) = w - \alpha \sum_{i=1}^{n} \nabla Q_i(w) \tag{1.5}$$

where $\eta$ is the learning rate.

But since in many cases, the update of one gradient would take to estimate the gradient over the whole dataset, the computation would take too much time. Often times we only need a rough estimate of the gradient, not a precise one.

What stochastic gradient descent does is to remove the sum and only uses one sample to estimate the gradient.

$$w := w - \eta \nabla Q(w) = w - \alpha \nabla Q(w) \tag{1.6}$$

In practice this turns out to be less effective than batch gradient descent, since we only get a very rough estimate of the gradient. Therefore the minibatch gradient descent technique is used. The difference here is that while batch gradient descent uses $n$ samples and stochastic gradient descent uses only $1$, minibatch gradient uses $b$ training samples, where $b << n$. Typical sizes for $b$ are in the range if $b \in 2, \dots, 100$.

Therefore in minibatch stochastic Gradient descent we randomly subset the training set and only take some samples out of it and estimate the gradient out of the seen sample set.

The important thing to do is to randomize the training set, since otherwise only a small portion of the training set will be seen and leads to wrong gradients.

### 1.1.3  The Forward Backward procedure in detail

To train the network we use the standard back-propagation algorithm, which essentially does a forward and a backward pass of the network in a row. Firstly we initialize the

weights in a matrix.

$$\mathbf{W}_{j \times k} = \begin{bmatrix} w_{11} & w_{12} & \ldots & w_{1k} \\ w_{21} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ w_{j1} & \ldots & \ldots & w_{jk} \end{bmatrix}, \mathbf{x}_{1 \times k} = \begin{bmatrix} a_1 \\ \vdots \\ \vdots \\ a_k \end{bmatrix} \tag{1.7}$$

**Forward propagation** is one of the two passes the network needs to do ( hence it's name ). In forward propagation the network calculates the predicted output of the network.

The output of the hidden layer $j$ of the network can be calculated as the weighted sum of the inputs

$$\mathbf{nnet}_j = \sum_{k=1}^{n} w_{kj} x_k = \mathbf{Wx} \tag{1.8}$$

Even though this method should already work, we add a bias to every node to increase the learning speed and void that the network stops learning.

$$\mathbf{nnet}_j = \sum_{k=1}^{n} w_{kj} x_k + b_k = \mathbf{Wx} + \mathbf{b} \tag{1.9}$$

r

Later we see that back-propagation needs some variables, which can be already pre-computed during the forward step. These two variables are the output $\mathbf{o}_j$ of layer $j$ and the derivative of the output w.r.t to the weighted sum $\mathbf{nnet}_j$. We represent the derivatives in a matrix $\mathbf{D}$ and the output of the current layer

$$\mathbf{D}_j = \frac{\mathbf{o}_j}{\mathbf{nnet}_j} \tag{1.10}$$

In case of sigmoid activation function, we obtain:

$$\mathbf{D}_j = \varphi\left(\mathbf{nnet}_j\right)\left(1 - \varphi\left(\mathbf{nnet}_j\right)\right) \tag{1.11}$$

For the outputs, we simply feed forward the network and store the output $\mathbf{o}_j$ in out

buffered lists.

$$\mathbf{o}_j = \varphi\left(\mathbf{Wx} + \mathbf{b}\right) \tag{1.12}$$

As soon as the feed forward is done, we produced an output of the network, which we can denote as $\mathbf{o}_L$. Now we need to decide how close our output is to the targetted output $\mathbf{t}$. For this we use a cost function in usual cases it is sufficient to use $MSE$ as such, which is defined as:

$$\mathbf{err} = \frac{1}{2}\left(\mathbf{o}_L - \mathbf{t}\right)^2 \tag{1.13}$$

Here (1.13) we need to make sure that the last output layer has the same dimensions as the target output.

**Backpropagation**    starts when all layers are trained. The idea behind back propagation is that the hidden layers do not produce any output. So we cannot modify their weights directly, since we don't know how large the error was. With back propagation we calculate the error at the output layer $L$ and then propagate this error to the hidden layers back. Therefore we update all $L - 1$ weight matrices and biases.
At first we calculate the differences between the target output and the estimated output.

$$\boldsymbol{\delta}_L = \mathbf{err} \otimes \mathbf{D}_L \tag{1.14}$$
$$\nabla\mathbf{W}_L = \mathbf{o}_L \boldsymbol{\delta}_L^T \tag{1.15}$$

We store both, the deltas and the nablas to later update the weights ( with the $\nabla$s) and the biases (using the $\delta$). From here on we then calculate the other layers backwards, we do:

$$\boldsymbol{\delta}_i = \mathbf{D}_i \otimes \left(\mathbf{W}_{i+1}^T \boldsymbol{\delta}_{i+1}\right) \tag{1.16}$$
$$\nabla\mathbf{W}_i = \mathbf{o}_i \boldsymbol{\delta}_i^T \tag{1.17}$$

Whereas we again store the deltas and the nablas to later update the weights.
To update the weights, we use the usual gradient descent update rule:

$$\mathbf{W}^* = \mathbf{W} - \alpha\nabla\mathbf{W}_i^T \tag{1.18}$$

$$\tag{1.19}$$

To update the biases, we use essentially the same update rule, but only consider the given $\delta_i$s.

$$\mathbf{b}^* = \mathbf{b} - \alpha\boldsymbol{\delta}_i^T \tag{1.20}$$

Finally, if we would like to improve the converging speed we can apply momentum. Momentum adds extra "velocity" towards the gradient curve, by using the last estimated value of the gradient $\mathbf{W}_{i-1}^*$ ($i$ denotes the current iteration) and applying on that the momentum ($\alpha$) as:

$$\mathbf{W}_{i+1} = \mathbf{W}_i - \eta\nabla\mathbf{W}_i + \alpha\mathbf{W}_i \tag{1.21}$$

The momentum is initialized with 0 and takes effect after one iteration of gradient descent.

## 1.2 Implementation details

OpenCL offers multiple implementation types. It is natively written in C, but has also a C++ wrapper onboard. We used in our project the C++ wrapper since it is more naturally to use that in a C++ project. First one needs to include the necessary header into any class.

```
#include <CL/cl.hpp>
```

Listing 1.1: OpenCL C++ header

**The interface to OpenCL** was written in C++11 and makes heavily use of the current variadic args feature. For this small scale task, the interface is definitely too complex, but it can be used for any other OpenCL task.

```
__kernel void mult(const int wSrc, __global const float* A,__global const ↩
    float* B,__global float* output)
{
    const int idx = get_global_id(0);
    const int idy = get_global_id(1);

```

```
6     output[idy*wSrc+idx] = A[idy*wSrc+idx] * B[idy*wSrc+idx];
7 }
```

Listing 1.2: Example Kernel function

OpenCL uses externally defined functions ( coined as kernels ) to compile the code and run it, during the runtime. Therefore one needs to write a kernel for its needs. An example kernel can be seen at 1.2. In our case we defined various kernels, for dot products and other matrix operations.

Since OpenCL is a multi device and platform GPU/CPU interface, one needs to first figure out which kind of platform (e.g. NVIDIA,AMD,Intel) the current machine is using and then decide which accelerator will be used 1.2.

```
1  void exampleplatform(const char * programpath){
2  std::vector<cl::Platform> all_platforms;
3  cl::Platform::get(&all_platforms);
4  //Assume having only one platform
5  cl::Platform defaultplatform = all_platforms[0];
6  std::vector<cl::Device> all_gpu_devices;
7  //CL_DEVICE_TYPE_GPU can also be CL_DEVICE_TYPE_CPU for CPU
8  defaultplatform.getDevices(CL_DEVICE_TYPE_GPU, &all_gpu_devices);
9  //Assume having only one device / take the first
10 cl::Device device = all_gpu_devices.front();
11 //Init the current context with the device
12 cl::Context(device);
13 //We wrote a helper function here to get the string data out of the kernel↩
       file
14 const char* content = util::file_contents(programpath);
15 }
```

Listing 1.3: Get Platforms and devices

Moreover since OpenCL is compiled during runtime, one needs to extract the code from the kernel file (e.g. kernel.cl) and give this content to the OpenCL.

As soon as the context is initialized we can run our kernels on the device. The problem here is that the kernel is defined with it's own parameters and types, but we dont know these during the compile time. To have a universal interface we used the variadicargs feature to allow the programmer a very straight forward way to use any kernel function.

```cpp
1  void runKernel(const char *kernelname){
2  cl::Program::Sources sources;
3  //contents is the already read out content of the kernels file (e.g. ←
       kernels.cl)
4  sources.push_back(std::make_pair(contents,strlen(contents)+1));
5  //Init the program with the context and the source
6  cl::Program program(context,sources);
7  //Build the program on the device
8  program.build({device});
9  cl::CommandQueue queue(context,device);
10 //The operator allows us to set args to the kernel
11 cl::Kernel kernel_operator(program,kernelname);
12 //kernel_operator allows us to send arguments to the kernel by calling
13 //kernel_operator.setArgs(ARGNUMBER,ARGUMENT);
14
15 //Init space on the device using cl::Buffers
16 //Send Arguments to the device ....
17 // quene.enqueueWriteBuffer() .....
18 //Wait until the arguments did arrive on the device
19 queue.finish();
20
21 //Execute!
22 cl::Event event;
23 queue.enqueueNDRangeKernel(kernel_operator,cl::NullRange,cl::NDRange(10),←
       cl::NullRange,NULL,&event);
24 //Wait until the execution has finished
25 event.wait();
26 // Read out the results by calling
27 // quene.enqueueReadBuffer()
28 }
```

Listing 1.4: Kernel usage

The usual Kernel execution can be seen in 1.2. Even though it is recommended to use OpenCL in this fashion, we cannot cope with different arguments for the kernel. Therefore we would need to init a device and context every time ( or at least as a singleton ) and then rewrite the argument passing for every parameter independently. This would take a lot of time, if the GPU is used extensively.

Our solution to that problem is as follows:

```
1  class OpenCL{
2
3  // Constructor ..... etc
4  // Hook for the iteration
5  template<std::size_t P=0,typename... Tp>
6     typename std::enable_if<P == sizeof...(Tp), void>::type addkernelargs(←
          std::tuple<Tp ...>&& t,cl::Kernel &k,cl::CommandQueue &,std::vector<←
          cl::Buffer> &outputbuffers) const{
7     // Do nothing
8     }
9
10 //  Start of the iteration
11 template<std::size_t P = 0, typename... Tp>
12    typename std::enable_if< P < sizeof...(Tp), void >::type addkernelargs(←
          std::tuple<Tp...> && t,cl::Kernel &kernel,cl::CommandQueue &,std::←
          vector<cl::Buffer> &outputbuffers) const{
13        // Type
14         typedef typename std::tuple_element<P, std::tuple<Tp...> >::type ←
             type;
15
16         // Add the value of the current item from std::get<P> to the args ←
             in kernel
17         // This function decides which type the kernel arg is
18         addkernelarg(P, std::get<P>(t), kernel,queue,outputbuffers);
19
20         // Recurse to get the remaining args
21         addkernelargs<P + 1, Tp... >(std::forward<std::tuple<Tp...>>(t), ←
             kernel,queue,outputbuffers);
22    }
23
24 // Adding Std::vector as type to the kernel args list
25    template<typename T>
26    void addkernelarg(std::size_t i, std::vector<T> const & arg, cl::Kernel←
          & kernel,cl::CommandQueue &) const{
27       cl::Buffer buffer(this->context,CL_MEM_READ_WRITE,arg.size()*sizeof(←
           T));
28       queue.enqueueWriteBuffer(buffer,CL_FALSE,0,sizeof(T)*arg.size(),&(←
           arg[0]));
29       kernel.setArg(i,buffer);
30
31    }
```

```
32

33

34  // Adding any array into the kernel args
35      template<typename T,std::size_t N>
36      void addkernelarg(std::size_t i, T const (& arg)[N], cl::Kernel & ←
            kernel,cl::CommandQueue &) const{
37          cl::Buffer buffer(this ->context,CL_MEM_READ_WRITE,N*sizeof(T));
38          queue.enqueueWriteBuffer(buffer,CL_FALSE,0,sizeof(T)*N,&arg);
39          kernel.setArg(i,buffer);
40      }

41

42      // Adding any constant to the kernel
43      template<typename T>
44      void addkernelarg(std::size_t i, T const & arg, cl::Kernel & kernel,cl←
            ::CommandQueue &) const{
45          cl::Buffer buffer(this ->context,CL_MEM_READ_WRITE,arg.size()*sizeof(←
                T));
46          queue.enqueueWriteBuffer(buffer,CL_FALSE,0,sizeof(T)*arg.size(),&(←
                arg[0]));
47          kernel.setArg(i,buffer);
48      }

49

50  }
```

Listing 1.5: Add arguments to kernel dynamic way

This recipe can be used to do the same actions for the reading buffers out after the transfer has finished. Therefore we can create a highly dynamic wrapper class for OpenCL.

## 1.2.1 Using the Interface

# Bibliography

[Cau47] Cauchy, Augustin: *Méthode générale pour la résolution des systèmes d'équations simultanées*. 1847. -- 536 -- 538 S. `http://gallica.bnf.fr/ark:/12148/bpt6k2982c.image.f540.pagination.langEN`

[Kiw01] Kiwiel, Krzysztof C.: Convergence and efficiency of subgradient methods for quasiconvex minimization. In: *Mathematical Programming* 90 (2001), März, Nr. 1, 1--25. `http://dx.doi.org/10.1007/PL00011414`. -- DOI 10.1007/PL00011414. -- ISSN 0025--5610

[Qia99] Qian, Ning: On the momentum term in gradient descent learning algorithms. In: *Neural Networks* 12 (1999), Nr. 1, S. 145 -- 151

[Yua99] Yuan, Ya-xiang: Step-sizes for the gradient method. In: *AMS/IP Studies in Advanced Mathematics* 42 (1999), Nr. 2