

Machine Learning

Lecture 9

An Artificial Neural Network consists of an input layer, a number of hidden layers and finally an output layer.

Assume the data are vectors in \mathbb{R}^d i.e. d -dimensional tuples of real numbers.

We are counting the layers from the top so the output layer is layer 0 and if there are a total of N layers, the input layer is layer N , so the hidden layers are numbered $1, 2, \dots, N - 1$

Each layer consists of 3 pieces of data

- A matrix W
- A bias vector \underline{b}
- An activation function $\mathbb{R} \rightarrow \mathbb{R}$, which can be for instance the logistic function σ or the hyperbolic tangent \tanh . Nowadays the most common activation function is the *ReLU*, $x \mapsto \max(x, 0)$

A data vector \underline{x} travels through the network.

$$\underline{x} \mapsto (\underline{u}_N = (\underline{x}W_N + \underline{b}_N)) \mapsto (\phi_N(\underline{u}_N) = \underline{z}_{N-1}) \mapsto (\underline{u}_{N-1} = (\underline{z}_{N-1}W_{N-1} + \underline{b}_{N-1})) \mapsto \dots$$

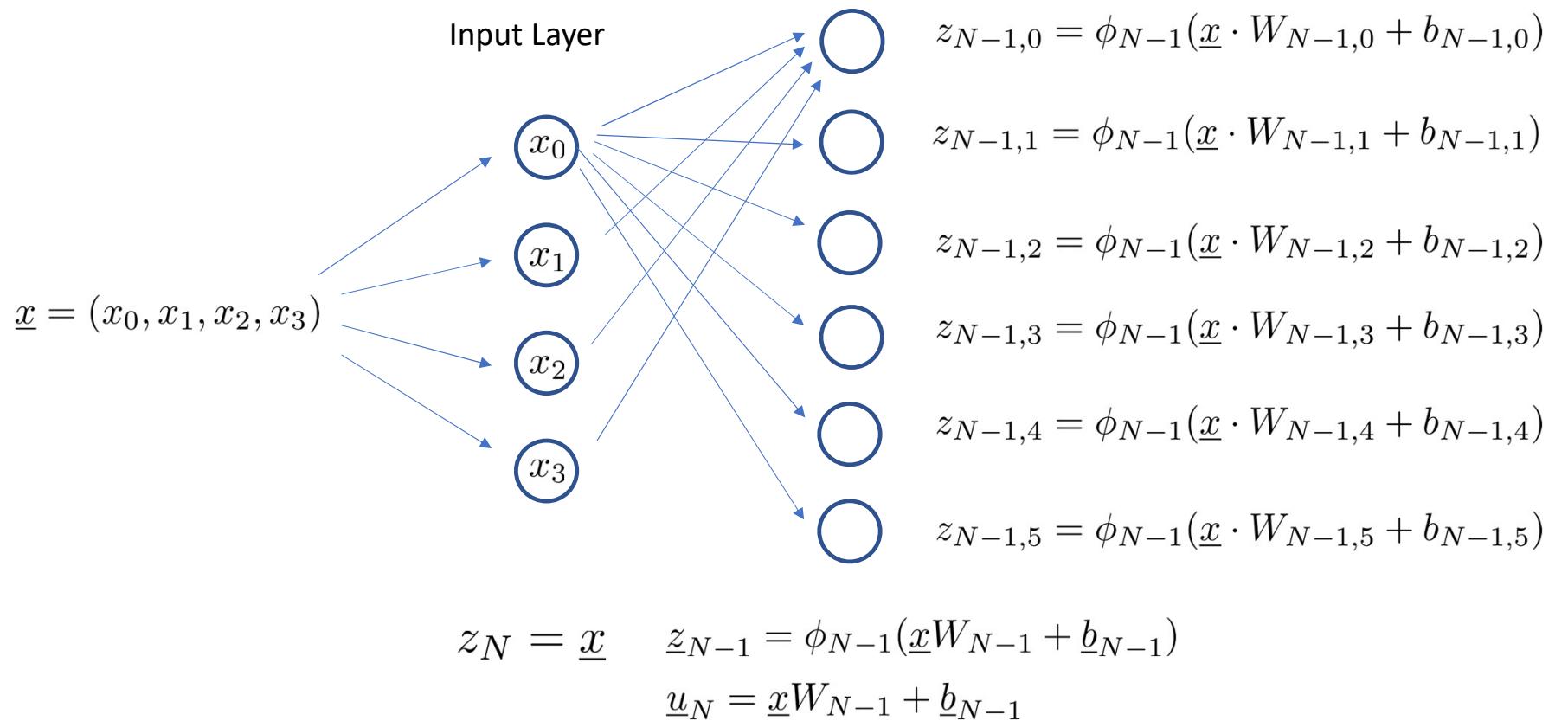
where the matrix W_i , the bias vector \underline{b}_i and the activation function ϕ_i are the data belonging to layer i .

Remark that if W_i is a $d_i \times e_i$ dimensional matrix then \underline{b}_i is an e_i -dimensional vector and W_{i-1} must have e_i rows i.e. W_{i-1} is $d_{i-1} \times e_{i-1}$ with $d_{i-1} = e_i$

Hidden Layer N-1

$W_{N-1} = (W_{N-1,0} \ W_{N-1,1} \ W_{N-1,2} \ W_{N-1,3} \ W_{N-1,4} \ W_{N-1,5})$, of dimension 4×6 so each $W_{N-1,i}$ is a column vector of dimension 4.

Bias vector $\underline{b}_{N-1} = (b_{N-1,0}, b_{N-1,1}, \dots, b_{N-1,5})$ and activation function ϕ_{N-1}

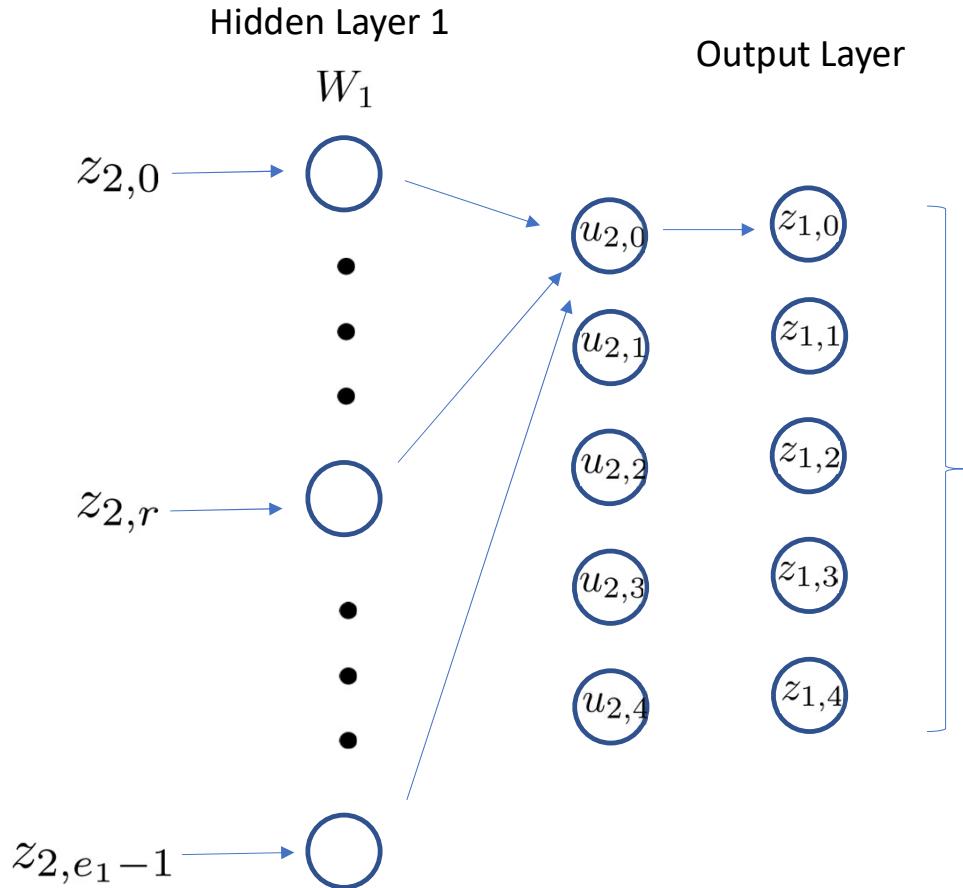


In general

$$u_{r+1} = z_{r+1} W_r + b_r$$

and

$$z_r = \phi_r(u_{r+1})$$



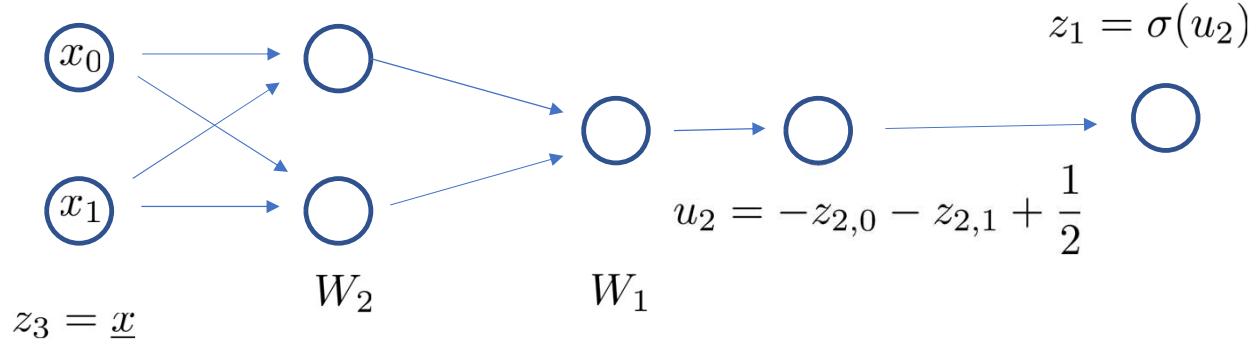
Final activation function , ϕ_1 is typically the softmax or the logistic function (classification) or the identity (regression)

$$\phi_1(\underline{z}_2 W_1 + \underline{b}_2) = \phi_1(\underline{u}_2) = \underline{z}_1$$

Let us consider an extremely simple network, with 1 hidden layer and input dimension 2. $W_2 = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$, $\underline{b}_2 = (0, 0)$ and $W_1 = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$, $\underline{b}_1 = \frac{1}{2}$. The activation functions ϕ_2 is the *ReLU* and ϕ_1 is the logistic function $\sigma(z) = \frac{1}{1 + e^{-z}}$

$$u_{3,0} = x_0 - x_1$$

$$z_{2,0} = \max(x_0 - x_1, 0)$$



$$u_{3,1} = -x_0 + x_1$$

$$z_{2,1} = \max(-x_0 + x_1, 0)$$

Consider the four input vectors $\underline{x}_1 = (0, 0), \underline{x}_2 = (1, 0), \underline{x}_3 = (0, 1), \underline{x}_4 = (1, 1)$

Then

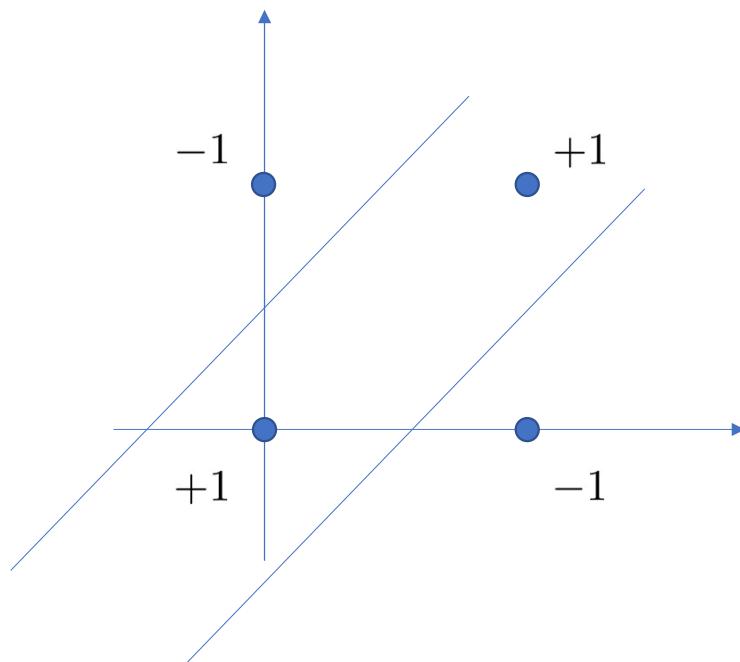
$$\underline{x}_1 \mapsto (0, 0) \mapsto \phi_1(0, 0) = (0, 0) \mapsto \frac{1}{2} \mapsto \sigma\left(\frac{1}{2}\right)$$

$$\underline{x}_2 \mapsto (1, -1) \mapsto \phi_1(1, -1) = (1, 0) \mapsto -1 + \frac{1}{2} \mapsto \sigma\left(-\frac{1}{2}\right)$$

$$\underline{x}_3 \mapsto (-1, 1) \mapsto \phi_1(-1, 1) = (0, 1) \mapsto -1 + \frac{1}{2} \mapsto \sigma\left(-\frac{1}{2}\right)$$

$$\underline{x}_4 \mapsto (0, 0) \mapsto \phi_1(0, 0) = (0, 0) \mapsto \frac{1}{2} \mapsto \sigma\left(\frac{1}{2}\right)$$

Thus this simple network separates the $\neg XOR$ data set



In a sense the graphic depiction of the ANN is more confusing than illuminating. The reality is that a (in this case *fully connected, feed forward Neural Network*) is just a composite function:

$$f : \underline{x} \mapsto \phi_1(\dots (\phi_{N-1}(\underline{x}W_{N-1} + \underline{b}_{N-1}) W_{N-2} + \underline{b}_{N-2}) \dots$$

or using the z and u variables

$$\begin{aligned} x &\mapsto \underline{z}_1 \\ \underline{z}_1 &= \phi_1(\underline{u}_2) \\ \underline{u}_2 &= \underline{z}_2 W_1 + \underline{b}_2 \\ \underline{z}_2 &= \phi_2(\underline{u}_3) \\ \underline{u}_3 &= \underline{z}_3 W_2 + \underline{b}_2 \\ &\vdots \end{aligned}$$

Let g_i be the function $g_i : \underline{z}_{i+1}, W_i, \underline{b}_i \mapsto \underline{z}_i = \phi_i(\underline{z}_{i+1}W_i + \underline{b}_i)$ and define the *feed-forward functions* $f_i, i = N - 1, \dots, 1$ by

$$f_i(\underline{x}, W_{N-1}, \dots, W_i, \underline{b}_N, \underline{b}_{N-1}, \dots, \underline{b}_i) = g_i \circ g_{i+1} \circ \dots \circ g_{N-1}$$

So

$$\begin{aligned} f_{N-1}(\underline{x}, W_{N-1}, \underline{b}_{N-1}) &= g_{N-1}(\underline{x}, W_{N-1}, \underline{b}_{N-1}) = \underline{z}_{N-1} \\ f_{N-2}(\underline{x}, W_{N-1}, W_{N-2}, \underline{b}_{N-1}, \underline{b}_{N-2}) &= g_{N-2}((g_{N-1}(\underline{z}_{N-1}, W_{N-1}, \underline{b}_{N-1}), W_{N-2}, \underline{b}_{N-2}) = \underline{z}_{N-2} \\ &\vdots \end{aligned}$$

The ANN itself is then the feed-forward function f_1 .

How do we train an ANN from a dataset i.e. how do we find parameters to fit the data set?

As usual we try to estimate the parameters to minimize a *loss function* \mathcal{L} , which can be for instance Cross Entropy in the classification case or Mean Squared Error (MSE) in the regression case.

The parameters of the ANN consists of the matrices W_{N-1}, \dots, W_1 and the bias vectors b_{N-1}, \dots, b_1 . Thus the total number of parameters is

$$\sum_{i=1} d_i e_i + \sum_{i=1} e_i = \sum_{i=1} d_{i-1}(d_i + 1)$$

Here d_{N-1} is the input dimension and d_1 is the output dimension.

For example if $N - 1 = 1$ and $\phi_1 = \sigma$, then the ANN is just the logistic regression $\underline{x} \mapsto \sigma(\underline{x}W + b)$, and there are $d_1 + 1$ parameters where $d_1 = \dim \underline{x}$.

We fix an input vector \underline{x} and the expected output y which is either a label or a real number and consider the loss

$$\mathcal{L}(y, f_1(\underline{x}, W_{N-1}, \dots, W_1, \underline{b}_{N-1}, \dots, \underline{b}_1))$$

We want to minimize the loss by doing gradient descent on the parameters.

This means computing the gradient with respect to all the parameters.
Fortunately this is not as difficult as it may seem.

We first compute the derivatives of a function of the form

$$g(\underline{z}, W, \underline{b}) = \phi(\underline{z}W + \underline{b})$$

where \underline{z} is an n -dimensional row vector, W is an $n \times m$ matrix, \underline{b} is m -dimensional and $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^m$.

Let $\underline{u} = (u_1, u_2, \dots, u_m) = \underline{z}W + \underline{b}$ then by the chain rule

$$\frac{\partial g}{\partial z_i} = \sum_k \frac{\partial g}{\partial u_k} \frac{\partial u_k}{\partial z_i}$$

Now

$$\frac{\partial g}{\partial u_k} = \left(\frac{\partial \phi^{(1)}}{\partial u_k}, \frac{\partial \phi^{(2)}}{\partial u_k}, \dots, \frac{\partial \phi^{(m)}}{\partial u_k} \right)$$

where $\phi = (\phi^{(1)}, \phi^{(2)}, \dots, \phi^{(m)})$ are the coordinate functions.

$$u_k = z_1 w_{1k} + z_2 w_{2k} + \dots + z_i w_{ik} + \dots + z_n w_{nk} + b_k$$

so

$$\frac{\partial u_k}{\partial z_i} = w_{ik}$$

and

$$\frac{\partial g}{\partial z_i} = \left(\sum_k w_{ik} \frac{\partial \phi^{(1)}}{\partial u_k}, \sum_k w_{ik} \frac{\partial \phi^{(2)}}{\partial u_k}, \dots, \sum_k w_{ik} \frac{\partial \phi^{(m)}}{\partial u_k} \right)$$

Hence $\frac{\partial g}{\partial z_i}$ is the i 'th row in the matrix product $\nabla \phi \cdot W^T$ where

$$\nabla \phi = \begin{pmatrix} \frac{\partial \phi^{(1)}}{\partial u_1} & \frac{\partial \phi^{(1)}}{\partial u_2} & \cdots & \frac{\partial \phi^{(1)}}{\partial u_m} \\ \frac{\partial \phi^{(2)}}{\partial u_1} & \frac{\partial \phi^{(2)}}{\partial u_2} & \cdots & \frac{\partial \phi^{(2)}}{\partial u_m} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial \phi^{(m)}}{\partial u_1} & \frac{\partial \phi^{(m)}}{\partial u_2} & \cdots & \frac{\partial \phi^{(m)}}{\partial u_m} \end{pmatrix}$$

is the Jacobian matrix of ϕ .

We write this as

$$\frac{\partial g}{\partial \underline{z}} = \nabla \phi \cdot W^T$$

In our case the activation function is just applying the *ReLU* function to each coordinate of the vector \underline{u} , so the Jacobian matrix is very simple

$$\begin{pmatrix} \frac{\partial \text{ReLU}}{\partial u_1} & 0 & 0 & \dots & 0 \\ 0 & \frac{\partial \text{ReLU}}{\partial u_2} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \frac{\partial \text{ReLU}}{\partial u_m} \end{pmatrix}$$

The derivative of the *ReLU* is the Heaviside function:

$$h(u) = \begin{cases} 0 & \text{if } u \leq 0 \\ 1 & \text{if } u > 0 \end{cases}$$

So the derivative

$$\frac{\partial g}{\partial \underline{z}} = h(\underline{u}) \cdot W^T$$

where

$$h(u) = \begin{pmatrix} h(u_1) & 0 & 0 & \dots & 0 \\ 0 & h(u_1) & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & h(u_m) \end{pmatrix}$$

Next we have

$$\frac{\partial g}{\partial w_{ij}} = \sum_k \frac{\partial g}{\partial u_k} \frac{\partial u_k}{\partial w_{ij}}$$

and

$$\frac{\partial u_k}{\partial w_{ij}} = \begin{cases} 0 & \text{if } j \neq k \\ z_i & \text{if } j = k \end{cases}$$

so

$$\frac{\partial g}{\partial w_{ij}} = \left(\frac{\partial \phi^{(1)}}{\partial u_j} z_i, \frac{\partial \phi^{(2)}}{\partial u_j} z_i, \dots, \frac{\partial \phi^{(m)}}{\partial u_j} z_i \right) = (0, 0, \dots, h(u_j)z_i, 0, \dots)$$

and we can write

$$\frac{\partial g}{\partial W} = \begin{pmatrix} \underline{z}^T \cdot (h(u_1), 0, \dots, 0) \\ \underline{z}^T \cdot (0, h(u_2), \dots, 0) \\ \vdots \\ \underline{z}^T \cdot (0, 0, \dots, h(u_m)) \end{pmatrix}$$

Remark that this is an array of matrices, each entry is the product of a $m \times 1$ by a $1 \times n$, hence an $m \times n$ matrix. We can view it as a *tensor* of dimension (m, n, m) .

Computing

$$\frac{\partial g}{\partial \underline{b}}$$

is very similar (but easier)

$$\frac{\partial g}{\partial b_i} = \sum_k \frac{\partial g}{\partial u_k} \frac{\partial u_k}{\partial b_i} = \sum_k \frac{\partial \phi}{\partial u_k} \frac{\partial (\underline{z} \cdot W_{|k} + b_k)}{\partial b_i} = \frac{\partial \phi}{\partial u_i}$$

Now assume we have a data point (y, \underline{x}) and assume we have a some loss function \mathcal{L} . We then want to compute $\frac{\partial \mathcal{L}}{\partial W_i}$ for $i = 1, 2, \dots, N - 1$. Using the chain rule we get

$$\frac{\partial \mathcal{L}}{\partial W_i} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \frac{\partial g_1}{\partial z_2} \cdot \frac{\partial g_2}{\partial z_3} \cdots \frac{\partial g_{i-1}}{\partial z_i} \cdot \frac{\partial g_i}{\partial W_i} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \nabla \phi_1 \cdot W_1^T \cdot \nabla \phi_2 \cdot W_2^T \cdots \nabla \phi_{i-1} \cdot W_{i-1}^T \cdot \frac{\partial g_i}{\partial W_i}$$

The gradients with respect to the bias vectors are computed similary:

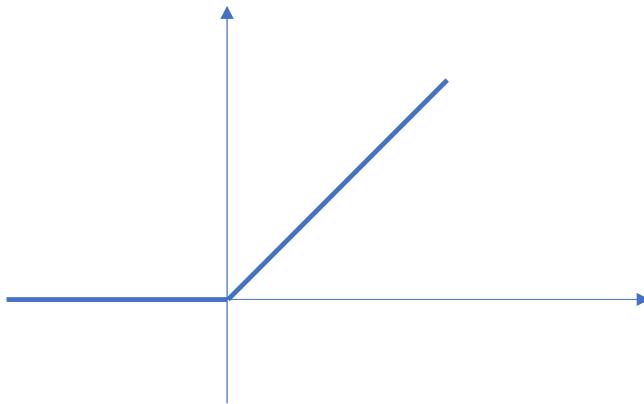
$$\frac{\partial \mathcal{L}}{\partial \underline{b}_i} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \frac{\partial g_1}{\partial z_1} \cdot \frac{\partial g_2}{\partial z_2} \cdots \frac{\partial g_{i-1}}{\partial z_{i-1}} \cdot \frac{\partial g_i}{\partial \underline{b}_i} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \nabla \phi_1 \cdot W_1^T \cdot \nabla \phi_2 \cdot W_2^T \cdots \nabla \phi_{i-1} \cdot W_{i-1}^T \cdot \frac{\partial g_i}{\partial \underline{b}_i}$$

Let's illustrate this with an example. We consider a 2-layer binary classification problem. The final activation function is given by sending a data point \underline{x} to $z_1 = \sigma(yf_1(\underline{x}))$. The other activation function ϕ_2 is the *ReLU* function. Thus the NN sends an input vector $\underline{x} = (x_1, x_2, \dots, x_n)$ to

$$z_1 = g_1(z_2, W_1, b_1) = \sigma(z_2 \cdot W_1 + b_1)$$

and

$$z_2 = g_2(\underline{x}, W_2, \underline{b}_2) = \text{ReLU}(\underline{x}W_2 + \underline{b}_2))$$



Here W_1 is an $m \times 1$ matrix, W_2 an $n \times m$ matrix, b_1 a scalar and \underline{b}_2 a $1 \times m$ vector.

Let $y = \pm 1$ be the label of \underline{x} then the loss function is

$$\mathcal{L}(z_1) = \begin{cases} -\log(z_1) & \text{if } y = 1 \\ -\log(1 - z_1) & \text{if } y = -1 \end{cases}$$

Hence we get

$$\frac{\partial \mathcal{L}}{\partial W_1} = \begin{cases} \frac{\partial \mathcal{L}}{\partial z_1} \cdot \frac{\partial g_1}{\partial W_1} = -\frac{1}{z_1} \cdot \underline{z}_2^T \cdot \sigma'(u_2) & \text{if } y = 1 \\ \frac{1}{1 - z_1} \cdot \underline{z}_2^T \cdot \sigma'(u_2) & \text{if } y = -1 \end{cases}$$

We get

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \frac{\partial g_1}{\partial z_2} \cdot \frac{\partial g_2}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial z_1} \sigma'(u_1) W_1^T \cdot \begin{pmatrix} \underline{x}^T \cdot (h(u_2^{(1)}), 0, 0, \dots, 0) \\ \underline{x}^T \cdot (0, h(u_2^{(2)}), 0, \dots, 0) \\ \vdots \\ \underline{x}^T \cdot (0, 0, \dots, 0, h(u_2^{(m)})) \end{pmatrix}$$

The product $W_1^T \cdot \begin{pmatrix} \underline{x}^T \cdot (h(u_2^{(1)}), 0, 0, \dots, 0) \\ \underline{x}^T \cdot (0, h(u_2^{(2)}), 0, \dots, 0) \\ \vdots \\ \underline{x}^T \cdot (0, 0, \dots, 0, h(u_2^{(m)})) \end{pmatrix}$ is a *tensor contraction*, it is the linear combination of $n \times m$ matrices

$$w_1 \left(\underline{x}^T \cdot (h(u_2^{(1)}), 0, 0, \dots, 0) \right) + \dots + w_m \left(\underline{x}^T \cdot (0, 0, \dots, 0, h(u_2^{(m)})) \right)$$

Adagrad

The Adagrad algorithm which means Adaptive Gradient algorithm works by adapting the learning rate at every step (like the Backtracking algorithm). The idea is to smoothen out spikes in the gradient by dividing each component of the gradient at step τ by the square root of the sum of squares of the component at the previous steps

Let $g_{\tau,j} = \sqrt{\sum_t^{\tau} ||\nabla Q_{i_t}(W_t)||^2}$ and let G_{τ} be the $k \times k$ diagonal matrix with the $g_{\tau,j}$ as the diagonal entries. Then the update rule is given by

$$W_{\tau+1} = W_{\tau} + \eta G_{\tau}^{-1} \nabla_W Q_{i_{\tau}}(W_{\tau})$$

Now we will use Pytorch to code a Neural Network

We start with some imports

```
1 import numpy as np
2 from sklearn.datasets import make_moons
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
```

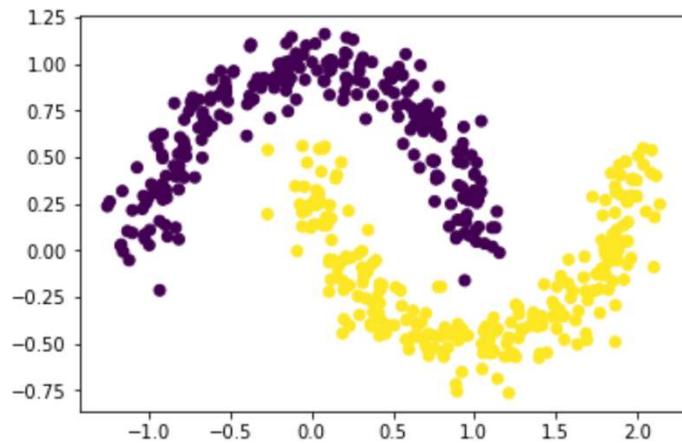
```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from torch.optim import SGD, Adagrad, RMSprop, Adam
```

If we want to use GPU computing (which Pytorch makes very simple) we want our data to be in a format that the GPU prefers. GPUs do calculations with 32 bit floats rather than the CPU which uses 64 bit floats. This is only a matter of precision in the very last decimals.

We convert our data to these formats

```
1 X,y = make_moons(500,noise=0.1)
2
3 y_ = [-1 if t == 0 else 1 for t in y]
```

```
1 plt.scatter(X[:,0],X[:,1],c=y_);
```



```
1 X = X.astype('float32')
2 y = y.astype('float32')
```

Our neural network is again a class, now it derives from the Pytorch class nn.Module

```
class Net(nn.Module):

    def __init__(self):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(2, 20)
        self.fc2 = nn.Linear(20, 40)
        self.fc3 = nn.Linear(40, 20)
        self.fc4 = nn.Linear(20, 2)
        self.ReLU = nn.ReLU()
```

In the constructor we first initialize the base class

```
def __init__(self):
    super(Net, self).__init__()
```

Here each layer has an input and an output dimension

fc means fully connected, but one can call it anything

Input dim

Output dim

```
self.fc1 = nn.Linear(2, 20)
self.fc2 = nn.Linear(20, 40)
self.fc3 = nn.Linear(40, 20)
self.fc4 = nn.Linear(20, 2)
```

Input dimension of the next layer matches the output dimension of the previous layer

These are all Linear layers with a matrix of parameters and bias vector

The class has two methods: the forward method is mandatory

Each layer is actually a linear function, u is the output

z is the output of applying the activation function, ReLU to u

Input data point

```
def forward(self,x):
    u = self.fc1(x)
    z = self.ReLU(u)
    u = self.fc2(z)
    z = self.ReLU(u)
    u = self.fc3(z)
    z = self.ReLU(u)
    u = self.fc4(z)
    z = torch.sigmoid(u)
```

```
return z
```

The `forward` method is automatically invoked from an instance of the class

```
ann = Net()
```

```
ann(x)      =      ann.forward(x)
```

We can also write a `predict` method

```
def predict(self,A):
    with torch.no_grad():
        output = self.forward(A)
        output.cpu()
        output_ = output.detach().numpy()
    return np.array([0 if x<1/2 else 1 for x in output_])
```

Feeding the data into the gradient descent loop is done by a special object `DataLoader`. To instantiate a `DataLoader` object we need to input another object of class `TensorDataset`. It is important that Pytorch deals exclusively with Tensor objects. Luckily it is easy to turn np.arrays into Tensors

```
X_p = torch.from_numpy(X)
y_p = torch.from_numpy(y).reshape(-1,1)
```

```
1 X_p
```

executed in 26ms, finished 18:50:34 2019-05-26

```
tensor([[ 1.6621e+00, -4.0047e-01],
       [ 9.4274e-01, -4.4460e-01],
       [ 1.0342e+00, -5.3469e-01],
       [ 2.5350e-01, -1.5420e-01],
       [-3.5460e-02,  9.2338e-01],
       [-2.2555e-01,  8.6350e-01],
       [ 3.1332e-02, -6.3017e-02],
       [ 2.5189e-01, -3.3939e-01],
       [ 2.0539e+00,  3.2519e-01],
       [ 2.9052e-01,  1.0639e+00],
```

```
1 X_p.shape
```

executed in 5ms, finished 18:51:43 2019-05-26

```
torch.Size([500, 2])
```

The most common error in using Pytorch is to not have tensors with the correct dimensions

```
1 | X_p[0:10].size()
```

executed in 5ms, finished 08:59:12 2019-05-27

```
torch.Size([10, 2])
```

```
1 | ann(X_p[0:10]).size()
```

executed in 7ms, finished 08:59:52 2019-05-27

```
torch.Size([10, 1])
```

We use the `nn.BCELoss` loss function.

The loss function computes

$$-y \log(\sigma(u)) - (1 - y) \log(1 - \sigma(u))$$

In order for this to make sense the tensors y and $\log(\sigma(u))$ must have the same size.

```
1 X_p[0:10].size()
```

executed in 5ms, finished 08:59:12 2019-05-27

```
torch.Size([10, 2])
```

```
1 torch.log(torch.sigmoid(ann(X_p[0:10]))).size()
```

executed in 10ms, finished 09:11:11 2019-05-27

```
torch.Size([10, 1])
```

If we just take `y_p = torch.from_numpy(y)`

we get

```
1 | y_p[0:10].size()
```

executed in 4ms, finished 09:14:07 2019-05-27

```
torch.Size([10])
```

and so does not have the shape we want (each label is a 0-dimensional tensor and we want a 1-dimensional tensor).

We fix this by reshaping

```
y_p = torch.from_numpy(y).reshape(-1, 1)
```

The `reshape (-1, 1)` command means to keep the first dimension (=length of the data set) and make the second dimension = 1

```
5 | y_p = torch.from_numpy(y).reshape(-1, 1)
```

executed in 5ms, finished 09:20:10 2019-05-27

```
1 | y_p[0:10].size()
```

executed in 5ms, finished 09:20:12 2019-05-27

```
: torch.Size([10, 1])
```

We can now instantiate our TensorDataset and DataLoader objects

```
from torch.utils.data import TensorDataset, DataLoader  
  
dataset = TensorDataset(X_p,y_p)  
training = DataLoader(dataset = dataset,batch_size=10,  
                      shuffle=True)
```

We also have to select a loss function, we choose BCEWithLogitsLoss = Binary Cross Entropy With Logits is our loss function from before

For each epoch we loop through the `DataLoader` object `training`, which emits a mini-batch of data points and corresponding labels

```
--  
for points, labels in training:
```

We send the data through the forward method to produce the logit `u` (the input to the sigmoid function)

```
output = ann(points)
```

and we can then compute the loss

```
loss = loss_fn(output, labels)
```

Next we have to select an optimizer.

We first try a simple SGD

```
from torch.optim import SGD
```

We have to write our training loop:

```
losses = []

for epoch in range(100):
    aggr_loss = 0.0
    for points, labels in training:
```

Now the optimizer comes in to minimize the loss:

first we zero all the gradients (if we had run a batch previously all the gradients of the loss function would have been computed and stored in the optimizer object)

```
optimizer1.zero_grad()
```

Then we compute all the gradients of the loss function with the data in the mini-batch. This is done with a single command

```
loss.backward()
```

which computes all the matrix products we had to compute in our hand coded model.

Finally the optimizer computes the step i.e. subtracting the learning rate multiple of the gradient from the parameters

```
optimizer1.step()
```

Here is the complete training loop

```
losses = []

for epoch in range(100):
    aggr_loss = 0.0
    for points,labels in training:

        output = ann(points)

        loss = loss_fn(output,labels)

        optimizer1.zero_grad()

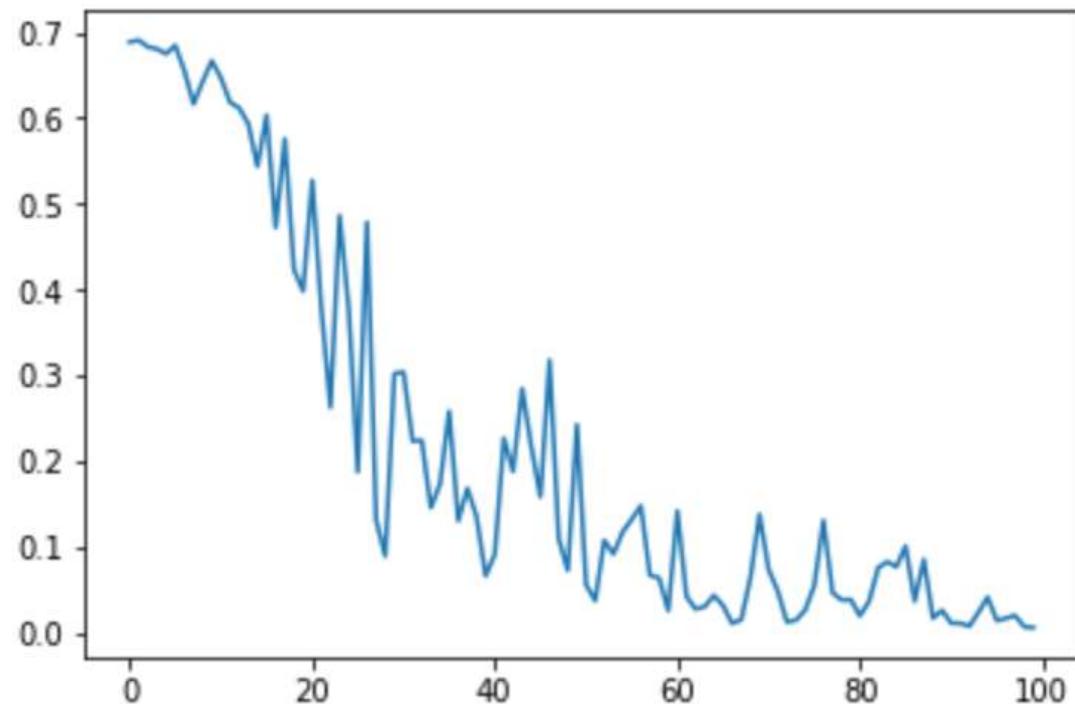
        loss.backward()

        optimizer1.step()

    losses.append(loss)
```

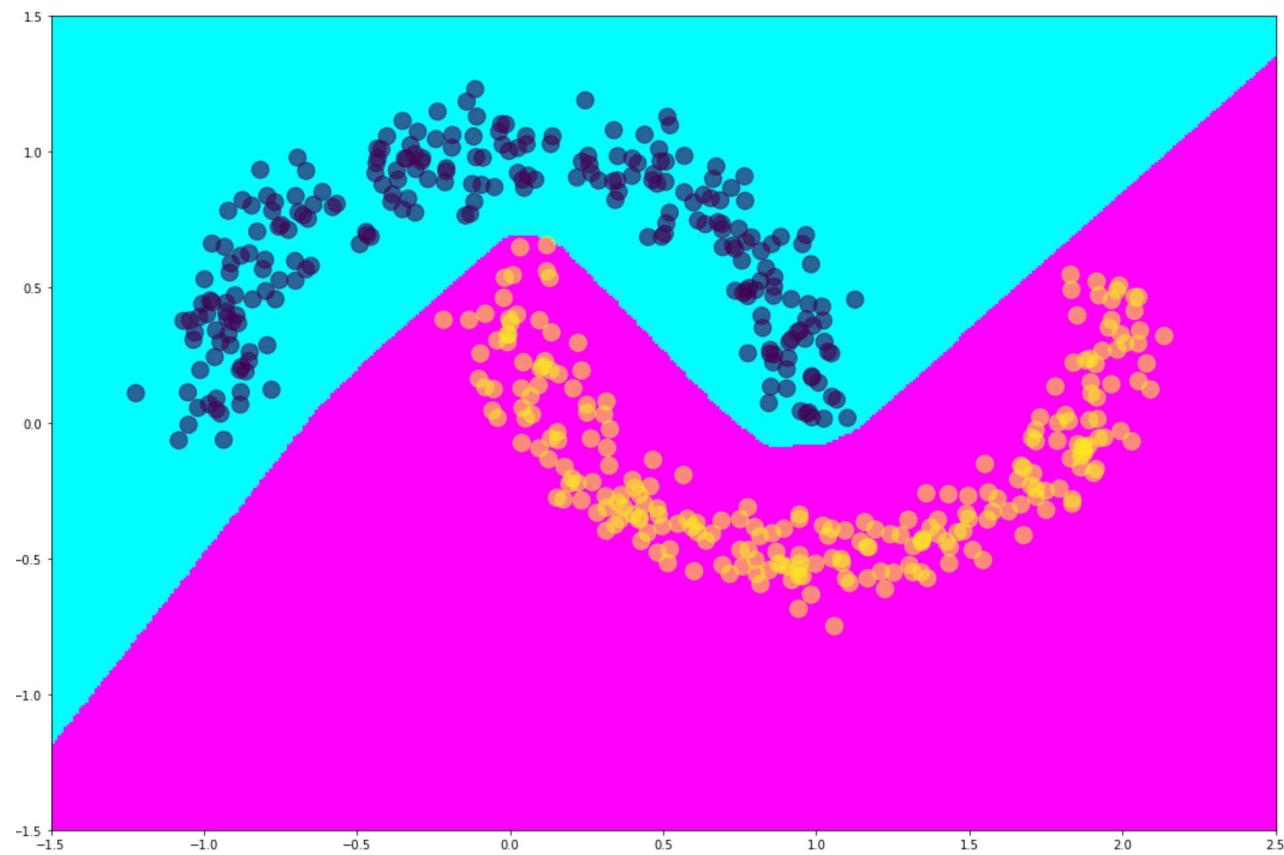
```
1 plt.plot(losses);
```

executed in 118ms, finished 09:39:52 2019-05-27



1	losses[-1]
executed in 22ms, finished 09:40:52 2019-05-27	

tensor(0.0064, grad_fn=<BinaryCrossEntropyWithLogitsBackward>)



SGD with momentum

The idea of the Momentum algorithm comes from physics. Viewing the parameter vector as a particle moving in space. Each update corresponds to giving the particle a push in a certain direction. But the particle was pushed at the previously and so already is moving in some direction i.e. it has momentum. This will change how the stochastic gradient update affects the parameter.

The idea is that this will speed up the rate of convergence.

Here is the algorithm

At epoch $t - 1$ the particle received a push of ΔW_{t-1} . Now at epoch t it receives a push of $-\eta \nabla_W Q_j(W_t)$ where Q_j is the part of the loss-function that comes from the training points in batch j .

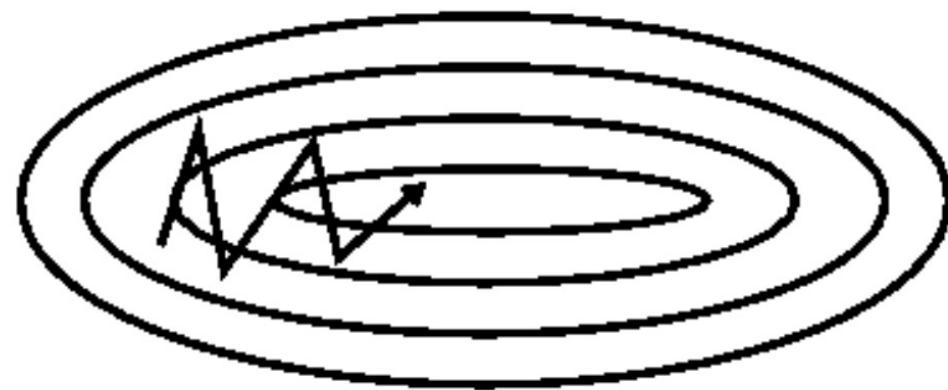
The new push is the convex combination of the previous momentum and the stochastic gradient:

$$\Delta W_t = (1 - \eta) \Delta W_{t-1} + \eta \nabla_W(Q_i(W_t))$$

so the update rule is:

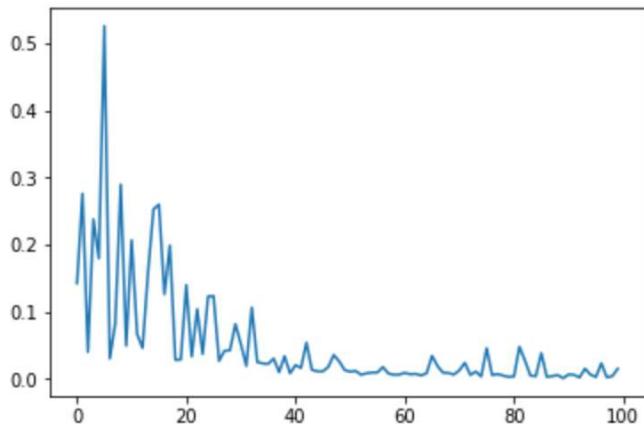
$$W_{t+1} = W_t - \alpha \Delta W_t$$

where α is a learning rate parameter. Thus the momentum optimizer requires two parameters α and η

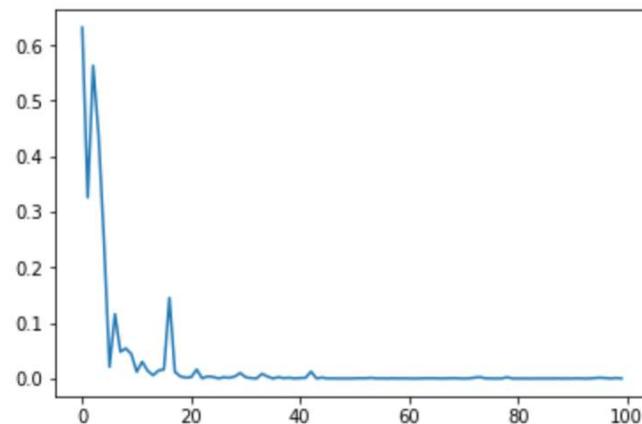


We can try other optimizers

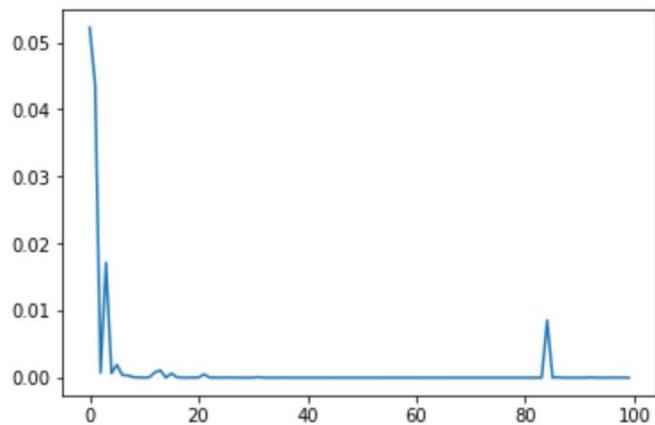
Adagrad



SGD with momentum



RMSProp



Adam

