

Lecture 5

Nonlinear neural networks

We start this lecture with a remark from lecture 4.

Why the more general three-layer network still doesn't work as expected?

The answer is simple: mathematical trivial result

For any two consecutive weighted sums of the input, there exists a single weighted sum with identical behavior. So anything that the three-layer network can do, the two-layer network can also do.

-2

It is clear, that the main aim in training any neural network is to find correlations between the input and output. Only in this case we can expect that this neural network can predict results for new input data.

So we want that middle layer to *sometimes* correlate with an input layer, and *sometimes* not correlate.

The middle layer will have the opportunity not just always be correlated to some input.

It can be correlated to one input only when it wants to be and other times not be correlated at all (conditional correlation).

This approach is based on nonlinear neurons. Let us take a first, simple but very important example.

If the node's value dropped below 0, normally for linear neurons it would still have the same correlation to the input. But now we turn of the node setting it to 0, when it would be negative.

We have zero correlation to any inputs whenever they are negative.

Thus, for example, we can construct a nonlinear NN for which some node on level 1 (middle layer) is perfectly correlated to input node 1 (on layer level 0) ^{only} when another input node 2 is turned off.

We see that by turning off any middle node whenever it would be negative we allow the NN to sometimes to correlation from various inputs. This result is impossible for two-layer NN.

There are many kinds of nonlinearities. We started from one simple but very popular example called **relu** activation function (rectified Linear Unit)

```
def relu(x):
```

```
    return (x > 0) * x
```

This function sets all negative numbers to 0.

Also we will need the derivative of this function (in order to calculate the gradient)

```
def relu2deriv (input):  
    return input input > 0
```

It returns 1 for $\text{input} > 0$ and returns 0 otherwise.

Now we are ready to modify the three-level NN proposed in Lecture 3.

This linear deep NN not converged during training stage. Let us check if application of relu type nonlinear neurons can change this situation.

-7-

$$N=3, M=4$$

$$\text{hidden-size} = 4 \quad \# P=4$$

$$\text{weights}_{01} = \dots \quad \# \text{random } [N \times P]$$

$$\text{weights}_{12} = \dots \quad \# \text{random } [P \times 1]$$

for iteration in range(60):

$$\text{layer 2-error} = 0$$

for i in range(len(inputs)):

$$\text{layer 0} = \text{inputs}[i:i+1] \quad \# M$$

$$\text{layer 1} = \text{relu}(\text{np.dot}(\text{layer 0},$$

$$\text{weights}_{01}))$$

$$\text{layer 2} = \text{np.dot}(\text{layer 1},$$

$$\text{weights}_{12})$$

$$\text{layer 2-error} += \text{np.sum}(\text{layer 2-goals}[i:i+1]) \times 2$$

$$\text{layer 2-goals}[i:i+1]) \times 2$$

$$\text{layer 2-delta} = \text{layer 2-goals}[i:i+1]$$

Back propagation

$$\text{layer 1-delta} = \text{layer 2-delta} \cdot$$

$$\text{dot}(\text{weights}_{12}, T) * \text{relu 2 deriv}$$
$$(\text{layer}_{-1})$$

$$\text{weights}_{12} += \text{alpha} * \text{layer 1. T. dot}(\text{layer}_{-2-delta})$$

$$\text{weights}_{01} += \text{alpha} * \text{layer 0. T. dot}(\text{layer}_{-1-delta})$$

if (iteration % 10 == 9) :

print("Error:" + str(layer2-error))

Let us make experiments - the following results are obtained

- Error : 0.634231
- Error : 0.358384
- Error : 0.083018
- Error : 0.006467
- Error : 0.000329
- Error : 1.505e-05.

After inclusion of nonlinear activation function the three-level NN \times have new possibilities to predict results for ~~non~~ data sets not used during training. The iterative process is converging fastly.

The *reLU* deriv function defines the slope (derivative) of the *reLU* activation function and it regularizes update of weights. Remember, that weights w_1 are used to calculate a prediction on level 1 and *reLU* function is used in this update step.

~~6~~-10-

How we compute the deltas for layer 1.

- first we take the standard step for three-layer NN: multiply the output (Layer 2) deltas by the weights w_{12} .
- second, we should take into account that we applied relu activation function in computation of Layer 1 nodes, and set value 0 if the predicted value was negative. Thus now

we apply relu derivative to regularize layer 1-delta values, by setting the delta of that node to 0.