

Neural Networks (Deep Learning)

In kernel method we made linear methods non linear by implicit (fixed) feature mapping $x \mapsto \phi(x)$ (which was possibly infinite dimensional)

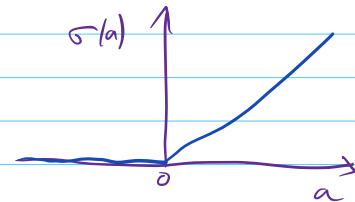
Neural networks, can think of the method as explicitly learning this feature map ϕ . Invented by Frank Rosenblatt in 1963 at Cornell!

A single neuron: $\sigma(U^T x + c)$
 $\sigma(a)$ is a non-linear activation function
 u are the weights of the neuron(s) we will learn

Eg. Activation functions:

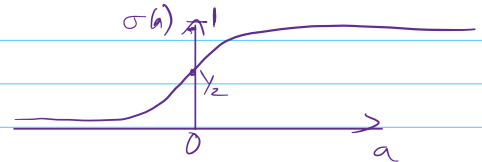
Relu: Rectified linear unit

$$\sigma(a) = \max\{a, 0\}$$



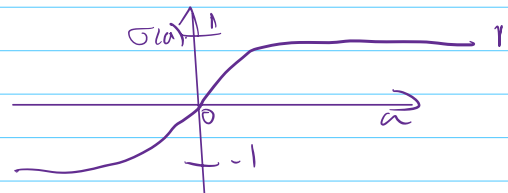
Sigmoid:

$$\sigma(a) = \frac{e^a}{1 + e^a}$$

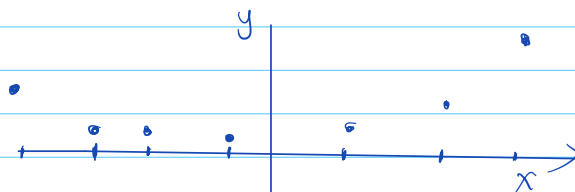


tanh:

$$\sigma(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$



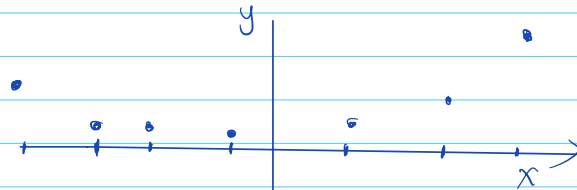
Regression with a single neuron:



A single hidden layer of neurons: (regression with k neurons)

$$\phi(x) = \sigma(U\vec{x} + \vec{c})$$

\vec{c} is a k dimensional vector
 U is a $k \times d$ matrix



e.g. ReLU

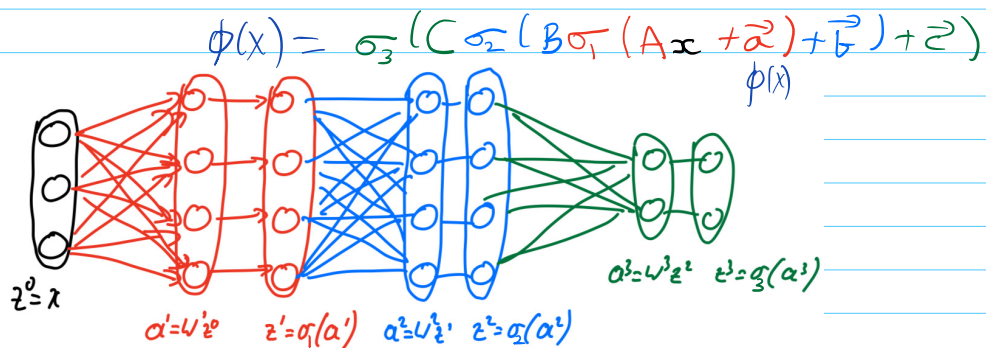
$$h(x) = w^T \phi(x) + b = w^T \sigma(Ux + \vec{c}) + b$$

$$= w^T \max(Ux + \vec{c}, \vec{0}) + b$$

What happens to the regression if $\sigma(a) = a$ (linear activation)?

Now lets try regression with RELU.

Deep Neural Network: We can stack multiple hidden layers one after the other, first layer gives linear boundaries, second layer combines corresponding non-linear functions, third layer combines non-linear functions output in second layer and so on. Depth quickly gives high expressibility.



In above example, how many neurons, how many parameters?

Forward Pass (prediction):

$$z_0 = x \quad (\text{Input Layer})$$

FOR $i = 1$ to L

$$a_i = W_i z_{i-1} + b_i$$

$$z_i = \sigma_i(a_i)$$

END

Return z_L

Typically σ is different from others

As usual bias can be aborted in each layer by inserting 1 to each layer as extra feature

We need to learn the weights and biases of each layer.
We will use gradient descent (or its cousin)!

How easy is it to compute gradients?

warmup with 1 dimensional example:

$$l(h(x), y) = \frac{1}{2} (h(x) - y)^2$$

$$h(x) = c \underbrace{\sigma(\underbrace{B \sigma(\underbrace{A x + a}_z) + b}_{z_2})}_{z_3} + c$$

$x \in \mathbb{R} \quad c, c, b, b, A, a \in \mathbb{R}$

l' derivative of loss
 σ' derivative of σ

gradients wrt.

$$\left. \begin{matrix} c \\ c' \\ c \end{matrix} \right\} \frac{\partial l(h(x), y)}{\partial c} = l'(h(x), y) \times \frac{\partial z_3}{\partial c} = l'(h(x), y) \times z_2$$

$$\frac{\partial l(h(x), y)}{\partial c} = l'(h(x), y) \times \frac{\partial z_3}{\partial c} = l'(h(x), y)$$

Already computed

$$\left. \begin{matrix} B, b \end{matrix} \right\} \frac{\partial l(h(x), y)}{\partial B} = l'(h(x), y) \times \frac{\partial z_3}{\partial B} = l'(h(x), y) \sigma'(a_2) \frac{\partial a_2}{\partial B}$$

$$= l'(h(x), y) \sigma'(a_2) z_1$$

$$\frac{\partial l(h(x), y)}{\partial b} = l'(h(x), y) \sigma'(a_2)$$

Backward pass (gradient step)

elementwise product

$$\vec{a} \odot \vec{f} = \begin{bmatrix} a_1 f_1 \\ a_2 f_2 \\ a_3 f_3 \\ \vdots \end{bmatrix}$$

$$\vec{\delta}_L = \nabla l(z_L, y) \odot \sigma'_L(a_L)$$

For $j = L-1$ to 1

$$W_j = W_j - \eta \vec{\delta}_j z_{j-1}^T$$

$$b_j = b_j - \eta \vec{\delta}_j$$

$$\vec{\delta}_{j-1} = \sigma'(a_{j-1}) \odot W_j^T \vec{\delta}_j$$

End

1. Compute a's and z's from forward pass (store them)
2. Run backward pass using a's and z's from forward pass (computes gradients)