

rows correspond to the neurons of the j -th layer

(we append -1 to encode the bias)

$$\text{net}^{(j)} = W^{(j)} \cdot X^{(j)} = \begin{bmatrix} w_{1j}^{(j)} \\ \vdots \end{bmatrix} \begin{bmatrix} x_1^{(j)} \\ \vdots \\ -1 \end{bmatrix} = \begin{bmatrix} \\ \end{bmatrix}$$

— net output, $j=1, 2, \dots, N$ of the j -th layer

$$Y^{(j)} = f[\text{net}^{(j)}] = \begin{bmatrix} f(\text{net}_1^{(j)}) \\ f(\text{net}_2^{(j)}) \\ \vdots \end{bmatrix}$$

— output of the j -th layer, $j=1, \dots, N$ (f — activation function)

$$X^{(j+1)} = \begin{bmatrix} Y^{(j)} \\ -1 \end{bmatrix}$$

— input of the $(j+1)$ -th layer ($j=1, 2, \dots, N-1$)

$$\rightarrow Y^{(N)} = \text{output of the network}$$

Back-propagation phase:

We have $X^{(j)}$, $\text{net}^{(j)}$, $Y^{(j)}$ from the forward pass.

the output was $Y^{(N)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$, let us suppose that we expected $t = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$

$$b^{(N)} = t - Y^{(N)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \text{ - error of the } N\text{th layer}$$

$$\delta^{(j)} = \begin{bmatrix} b_1^{(j)} \cdot f'(\text{net}_1^{(j)}) \\ b_2^{(j)} \cdot f'(\text{net}_2^{(j)}) \\ \dots \end{bmatrix} \quad \begin{array}{l} j = N, N-1, \dots, 1 \\ \text{ - } \text{ "delta signal" } \end{array}$$

$$\tilde{W}^{(j)} = W^{(j)} + c \delta^{(j)} \cdot (X^{(j)})^T \quad j = N, N-1, \dots, 1$$

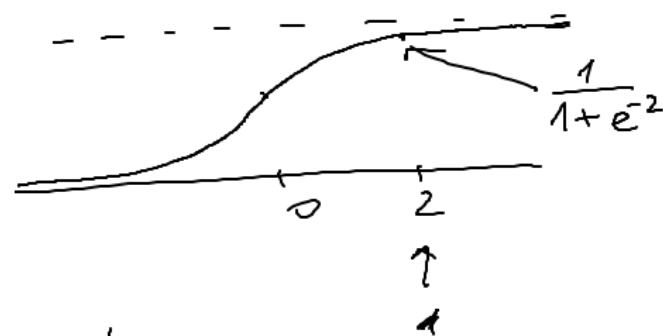
$\tilde{W}^{(j)}$ - weights after the adjustment

$$\tilde{b}^{(j)} = (W^{(j+1)})^T \cdot \delta^{(j+1)} \quad j = N-1, \dots, 1$$

$$b^{(j)} = \tilde{b}^{(j)} \text{ with last (bottom) entry removed, } j = N-1, \dots, 1$$

Say $f(x) = \sigma(x) = \frac{1}{1+e^{-x}}$

$\sigma'(x) = \sigma(x)(1-\sigma(x))$



Suppose $\text{net}_1^{(j)}$ has a large absolute value

then
$$g^{(j)} = \begin{bmatrix} b_1^{(j)} f'(\text{net}_1^{(j)}) \\ \vdots \\ \vdots \end{bmatrix}$$
close to 0

$$\tilde{W}^{(j)} = W^{(j)} + c \begin{bmatrix} b_1^{(j)} f'(\text{net}_1^{(j)}) \\ \vdots \\ \vdots \end{bmatrix} \cdot [x_1^{(j)} \ x_2^{(j)} \ \dots] =$$

$$= W^{(j)} + c \begin{bmatrix} \text{red oval} \\ \vdots \\ \vdots \end{bmatrix}$$
will be close to 0
these are the adjustments for the weights of the 1st neuron in the j-th layer

This is called 'saturation'.

the choice of initial weights

we know that we shouldn't choose the same weights for the neurons in the same layer

so let us take ^(i,j) random weights, e.g. $W_{ij}^{(1)} \sim N(0, \sigma^2)$

then the net output of the first neuron is

$$\text{net}_1^{(1)} = \sum_j W_{1j}^{(1)} x_j^{(1)}$$

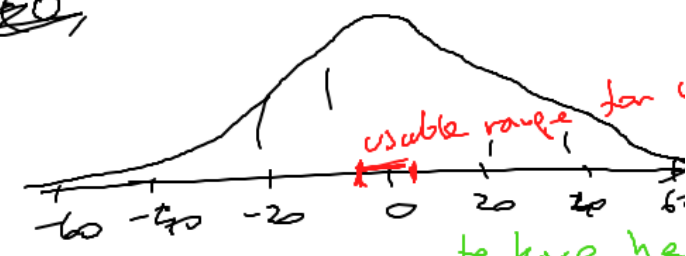
e.g. MNIST: $x_j^{(1)} \in [0, 1]$, $j = 1, \dots, 28^2 = 784$
 $x_{785}^{(1)} = -1$, $N = 785$

suppose half of inputs is 0 and half is 1

then $\text{net}_1^{(1)} = \text{sum of more-less } \frac{N}{2} \text{ variables } \sim N(0, \sigma^2)$
 $\sim N(0, \frac{N}{2} \sigma^2) \approx N(0, (20\sigma)^2)$

If we took $\sigma = 1$, then $\text{net}_1^{(1)} \sim N(0, 20^2)$

$\varphi'(\text{net}_1^{(1)})$ is likely

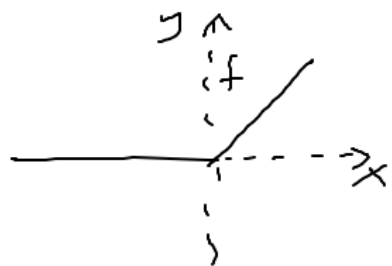


$\frac{N}{2} = 1$

... or take different activation function (apart from the last layer?)

$$f(x) = \max(x, 0)$$

(ReLU)



Training process:

• take one sample, compute output, compare with the target value,
use backpropagation to adjust the weight

• batch-learning:

take N samples, compute outputs, compare with the target values,

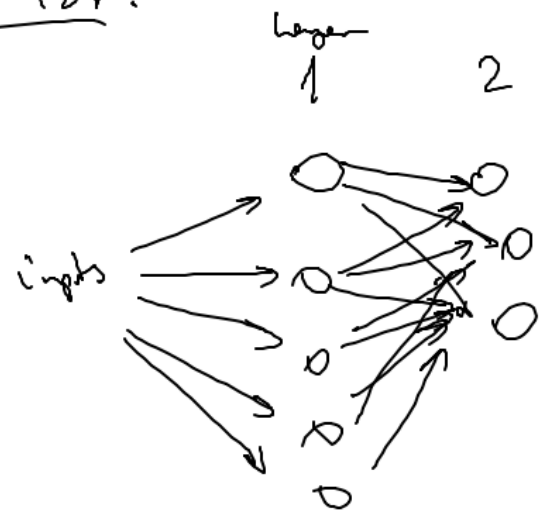
we backpropagation to compute the adjustments but

do not change the weights

afterwards adjust the weights by the sum/mean of the



So far:



dense
layers

~~neural net~~



Linear Regression

(supervised, continuous)

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f_{\theta}(x_i) = \begin{bmatrix} 1 & x_{i1} & \dots & x_{in} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} = \theta_0 + \sum_{j=1}^n \theta_j x_{ij}$$

← affine function ('linear')

$\theta_0, \dots, \theta_n$ - coefficients

model:

$$f(x_i) = f_{\theta}(x_i) + e_i, \quad i=1, \dots$$

↑

error,

e_i iid, mean 0

variance σ^2

given data (x_i, y_i) $i=1, \dots, n$
↑ inputs $\in \mathbb{R}^n$ ↑ output $\in \mathbb{R}$
 $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

we want to find "the best" coefficients θ ,

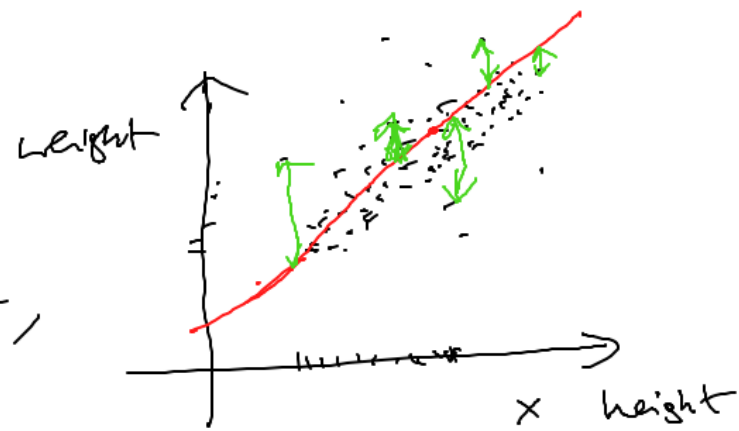
so that the formula

$$\hat{y} = f_{\theta}(x)$$

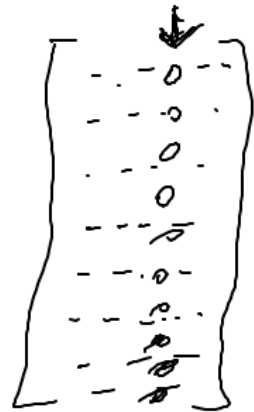
gives "the best" outputs.

one of the possible choices: minimise the square error:

$$L(y, \hat{y}) = \|y - \hat{y}\|_2^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



When $X = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ x_{21} & \dots & x_{2n} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{mn} \end{bmatrix}$ is of full rank



Then there is a unique \hat{y} minimising L :

$$\hat{y} = (X^T X)^{-1} X^T y$$

(note: there are better ways to compute \hat{y} numerically)

use scikit-learn.org (sklearn)

$$\begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

To make the model more flexible, one may add some terms,
e.g.

$$\underbrace{x_{ij} x_{ik}}_{\text{cross corr. terms}}, \quad x_{ij}^2, \quad x_{ij}^3, \dots$$

For one variable: numpy : polyfit

$$P_{\text{poly}}(x) = a_0 + a_1 \underbrace{x_{i,1}}_{\tilde{x}_{i,1}} + a_2 \underbrace{x_{i,1}^2}_{\tilde{x}_{i,2}} + \dots + a_n \underbrace{x_{i,1}^n}_{\tilde{x}_{i,n}}$$

