

$$w = w_{joi}$$

$$\frac{\partial L}{\partial w} = \sum_k (t_k - y_k) (-1) \varphi'(\underbrace{\sum_j v_{kj} a_j}_{\text{net output of the 2nd layer}}) v_{kj} \varphi'(\underbrace{\sum_i w_{joi} x_i}_{\text{net output of the 1st layer}}) \cdot x_{io}$$

$$\text{new } w = \text{old } w - c \cdot \frac{\partial L}{\partial w}$$

e.g. $\varphi(x) = \frac{1}{1+e^{-x}}$

$$\varphi'(x) = \varphi(x) \cdot (1 - \varphi(x))$$

if x is away from zero, then φ' is small (very small e.g. for $|x| > 10$)

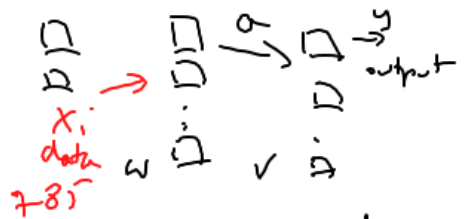
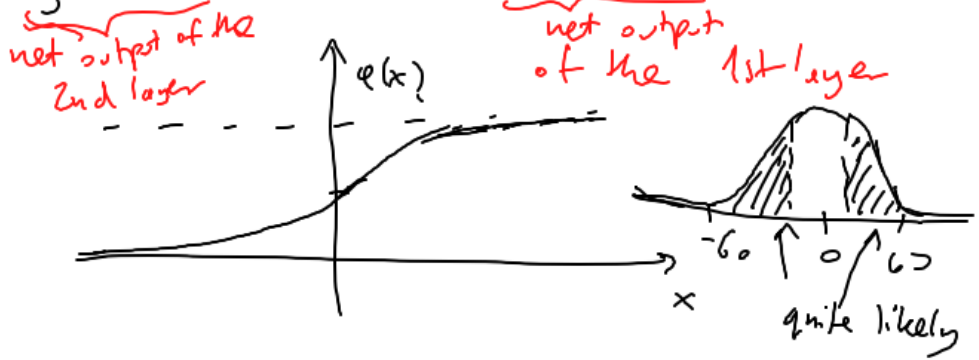
It's important to keep net outputs close to 0

- normalize the input
- initialize the weights in the right way

e.g. if $w \sim N(0, 1)$
 $(\frac{1}{\sqrt{785}})^2 \cdot 785$ inputs x

let's suppose half of x 's are 1, half are 0
 $\sum w_{joi} x_i \sim N(0, \boxed{\sqrt{392}^2})$

$\sigma \approx 20$
 $-60 \dots 60$

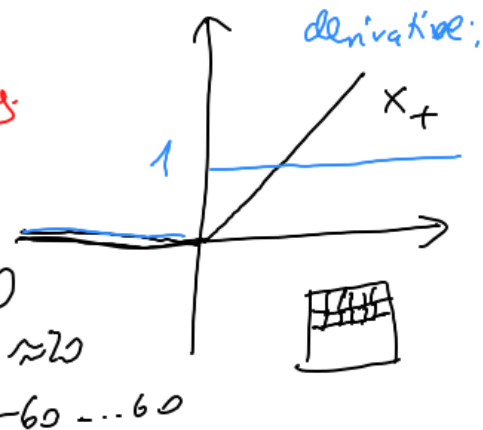


φ - activation function

$$a = \varphi(\sum_i w_{joi} x_i)$$

output of the 1st layer

we want e.g. $\frac{1}{2}$ here



{ Ridge regression (1st assignment):

we took all the training data and calculated the loss function based on that

↪ This method is too expensive in the use of neural networks.

• The simplest way: online learning: (example: 2nd assignment, images - classification)

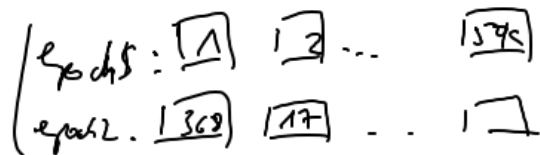
take some image, compute the output of the network, do the backprop. and weights' adjustment

↪ repeat for all images in the dataset



repeat some number of times (called epochs)

it's good to shuffle the training set between epochs (also before the first one)



loss function



Batch learning: take some number B of images
and take the average gradient for the adjustment

$$\text{new } w = \text{old } w - c \cdot \left(\frac{\partial L^{(1)}}{\partial w} + \dots + \frac{\partial L^{(B)}}{\partial w} \right) \cdot \frac{1}{B}$$

\uparrow
 derivative of the loss f.
 w/r to the weight w
 calculated using the image $w \perp$

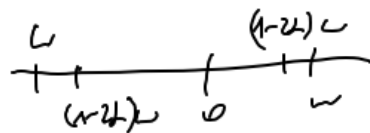
\swarrow one can incorporate it into c

regularization of weights:

~~add~~ take: (regularized loss function) = (usual loss function) + $\alpha \sum w^2$

$$\text{new } w = \text{old } w - c \left(\frac{\partial L}{\partial w} + 2\alpha w \right)$$

or
$$\text{new } w = \text{old } w \cdot (1 - 2\alpha c) - c \frac{\partial L}{\partial w}$$



w : weight
(but not for biases)
 $L \geq 0$ hyper parameter
small

Some history:

'40 first models

'50 first implementation (hardware)

'80 first applications - convolutional neural networks

2000 ...

Thm. (Cybenko)

Let σ - sigmoidal, cont. function

$\hookrightarrow \sigma(-\infty) = 0, \sigma(\infty) = 1, \sigma \nearrow$

Then the finite sums of the form

$$G(x) = \sum_{j=1}^m \alpha_j \sigma(y_j^T \cdot x + \theta_j)$$

\uparrow weights \uparrow input bias

are dense in $(C([0,1]^n), \|\cdot\|_{\infty})$.



$x \in [0,1]^n$

In other words, if $F: [0,1]^n \rightarrow \mathbb{R}$ continuous
then for any $\epsilon > 0$ one can find G as above
so that $|F(x) - G(x)| < \epsilon \quad \forall x \in [0,1]^n$