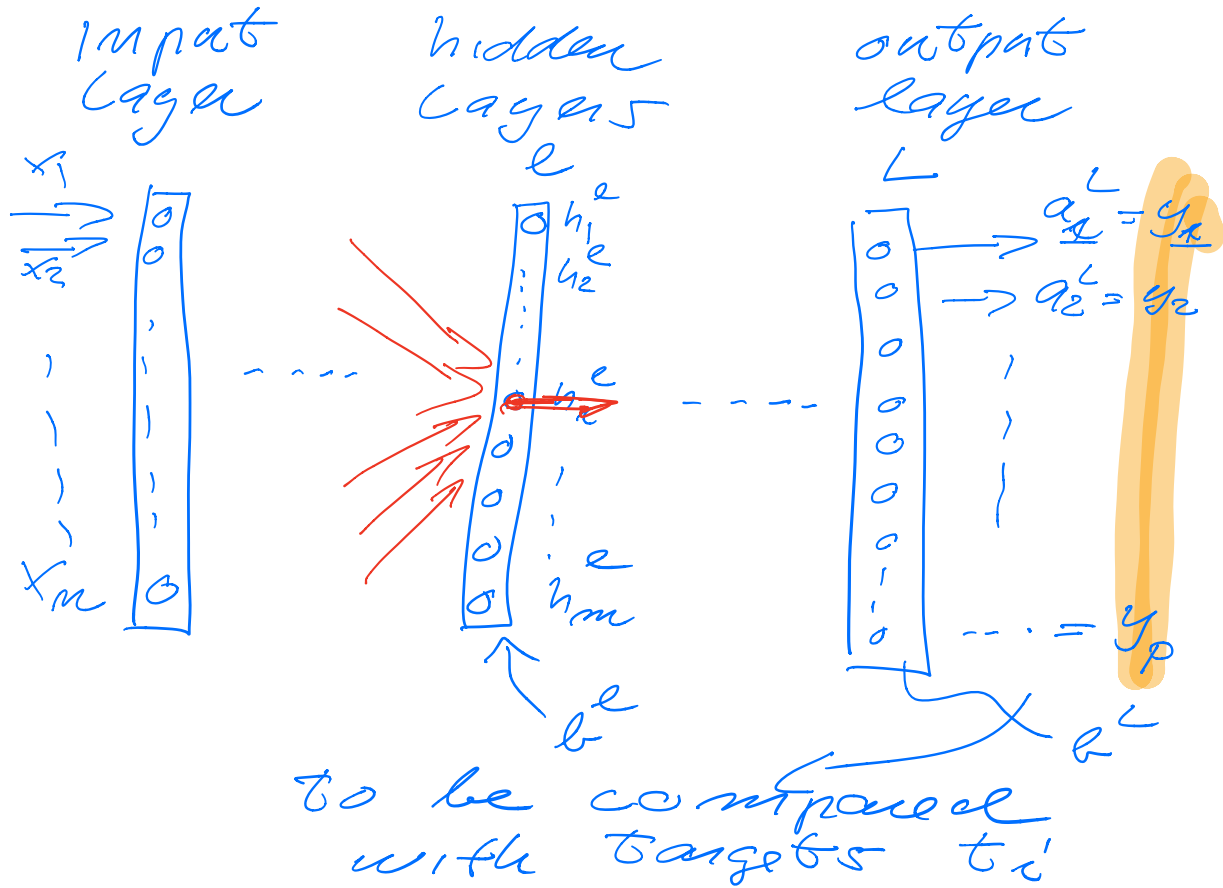


Lecture October 8

FFNN

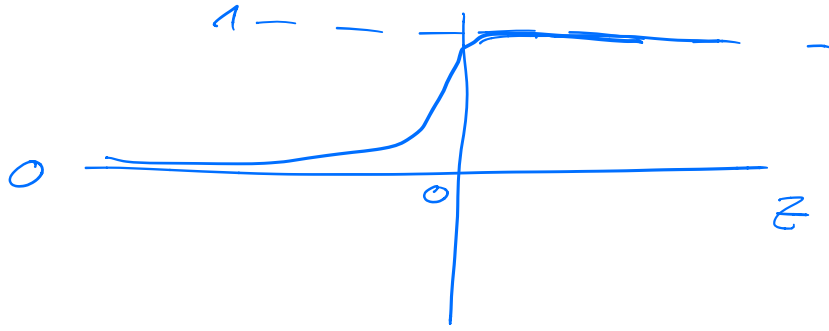


- activation functions

$$a_i^e = f(z_i^e) = f\left(\sum_j w_{ij}^e a_j^{e-1} + k_i^e\right)$$

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

$$f(z) = \tanh(z) = 2\sigma(2z) - 1$$



Rectified Linear Unit (ReLU)

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

output values only for $z > 0$,
that are different for zero.

Leaky ReLU

$$f(z_i, \alpha_i) = \max(0, z_i) + \alpha_i \min(0, z_i)$$

$$\alpha_i \sim 0.01$$

ELU

$$f(z) = \begin{cases} \alpha e^z - 1 & z < 0 \\ z & z \geq 0 \end{cases}$$

... many more

**Output layer with own
activation function + Cost.**

Cross entropy $\leftarrow \begin{cases} \text{Binary classification: } \sigma(z) \\ \text{More classes: Softmax} \end{cases}$

- Regression: MSE as cost function
 $\sigma(z)$ as activation.

$$\frac{1}{2} \sum_i (y_i^L - t_i)^2$$

architecture:

- # hidden layers (Depth)
- # nodes in each layer. (width)

Regularization parameters

- Ridge: add to cost function $\frac{\lambda}{2} W^T W$
- Lasso: $\lambda |W|_1$

SGD: epoch, learning, mini-batches

- standard
- Momentum SGD
- Adagrad
- RMS prop
- Adam

Feed Forward

$$a_i^{(1)} = f^{(1)}(W_i^{(1)} + b^{(1)})$$

$$a_i^{(2)} = f^{(2)}(w^{(2)} a_i^{(1)} + b^{(2)})$$

$$= f^{(2)}(w^{(2)} f^{(1)}(w^{(1)} x + b^{(1)}) + b^{(2)})$$

FFNN the basic operation is $w^{(l)(l-1)} a^{(l-1)} + b^{(l)}$

Backprop algo

$$\frac{\partial C(w^{(l)})}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} a_k^{(l-1)}$$

$$\delta_j^{(l)} = f'(\bar{z}_j^{(l)}) \frac{\partial C}{\partial a_j^{(l)}}$$

$$\delta_j^{(l)} = \frac{\partial C}{\partial f_j^{(l)}}$$

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} w_{kj}^{(l+1)} \cdot f'(\bar{z}_j^{(l)})$$

Gradient descent

$$w_{jk}^{(l)} \leftarrow w_{jk}^{(l)} - \gamma \delta_j^{(l)} a_k^{(l-1)}$$

$$\check{b}_j^{(e)} \leftarrow \check{b}_j^{(e)} - \delta \check{J}_j^{(e)}$$