

Machine Learning

Lecture 13
Intro to Neural Networks

Anna Kuzina

Learning Outcomes

After this lecture you should know:

- What is an artificial neuron, perceptron and artificial neural network
- What is a dense layer
- Which activation function are used
- How to train neural networks

Recap Discriminative models

Data:

$\{x_1, \dots, x_n\}$ - observed variables

$\{y_1, \dots, y_N\}$ - unobserved / target variables

Model with unknown parameters θ :

$$\theta^* = \arg \max_{\theta} \prod_{n=1}^N p(y_n | x_n, \theta)$$

↑

Maximum Likelihood $p(y | x, \theta)$ Bayesian (+ prior $p(\theta)$)

MAP $\theta^* = \arg \max_{\theta} \prod_{n=1}^N p(y_n | x_n, \theta) p(\theta)$ Full
↑

$p(\theta | x, \mathcal{F}) = \dots$

Recap Linear Regression

Data:

$\{x_1, \dots, x_n\}$ - observed variables

$\{y_1, \dots, y_N\}$ - unobserved / target variables

Model with unknown parameters θ :

$$p(y|x, \theta) = N(y | \omega^T x, \sigma^2)$$

ML: $\omega^* = \underset{\omega}{\operatorname{argmin}} \|X\omega - y\|^2$

Recap Logistic Regression

Data:

$\{x_1, \dots, x_n\}$ - observed variables

$\{y_1, \dots, y_N\}$ - unobserved / target variables $\in \{0, 1\}$

Model with unknown parameters θ :

$$\bullet p(y|x, \theta) = \text{Be}(y | \text{sigm}(\omega^T x))$$

$\text{sigm}(x) = \frac{1}{1 + \exp(-x)}$

$$P(y=1|x, \omega) = \text{sigm}(\omega^T x)$$

$$\omega^* = \arg \max_{\omega} \sum_{n=1}^N \log P(y_n | x_n, \omega)$$

Recap Discriminative models

Linear Regression:

$$p(y|x, \theta) = \mathcal{N}(y | w^T x, \sigma^2)$$

Logistic Regression:

$$p(y|x, \theta) = \mathcal{B}e(y | \text{sigm}(w^T x))$$

Can we parametrise distributions in other ways?

$$f_{\theta}(x)$$

Recap Discriminative models

Linear Regression:

$$p(y|x, \theta) = \mathcal{N}(y|w^T x, \sigma^2)$$

Logistic Regression:

$$p(y|x, \theta) = \mathcal{B}e(y|\text{sigm}(w^T x))$$

Can we parametrise distributions in other ways?

Regression:

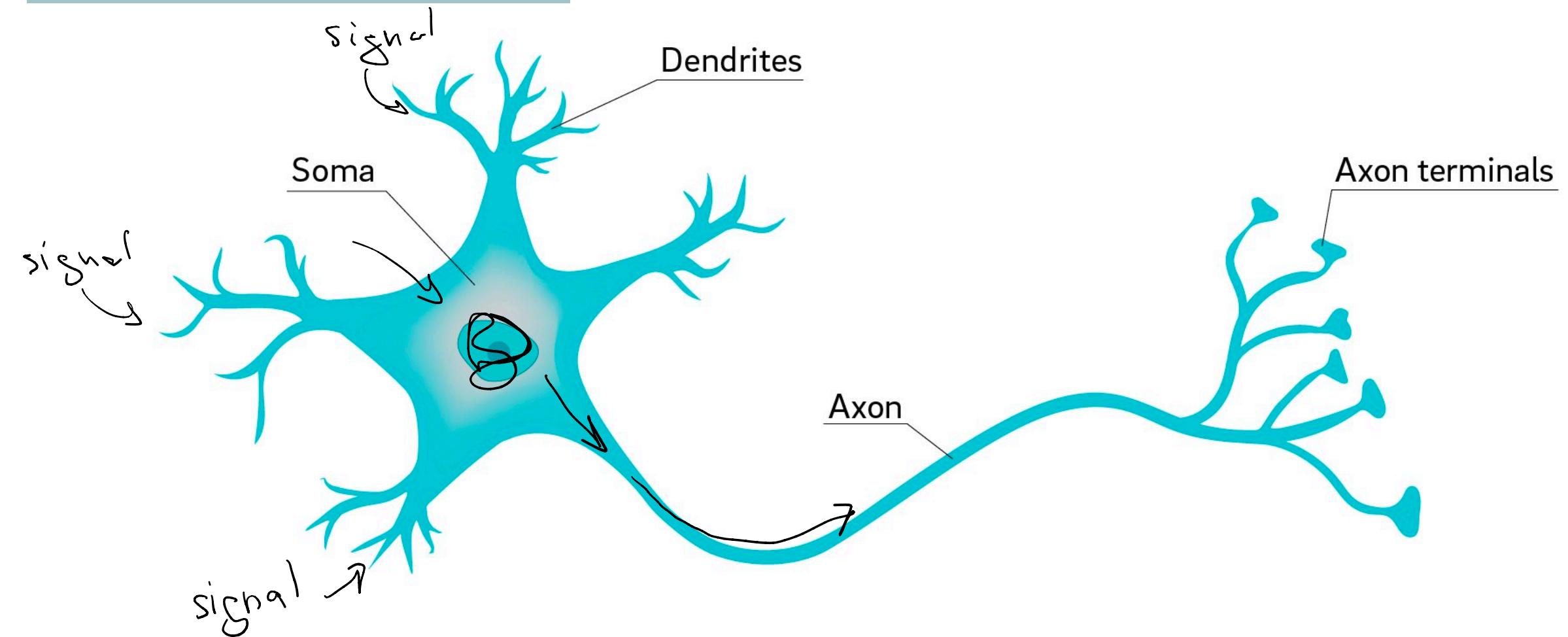
$$p(y|x, \theta) = \mathcal{N}(y|f(x), \sigma^2)$$

Binary Classification:

$$p(y|x, \theta) = \mathcal{B}e(y|f(x))$$

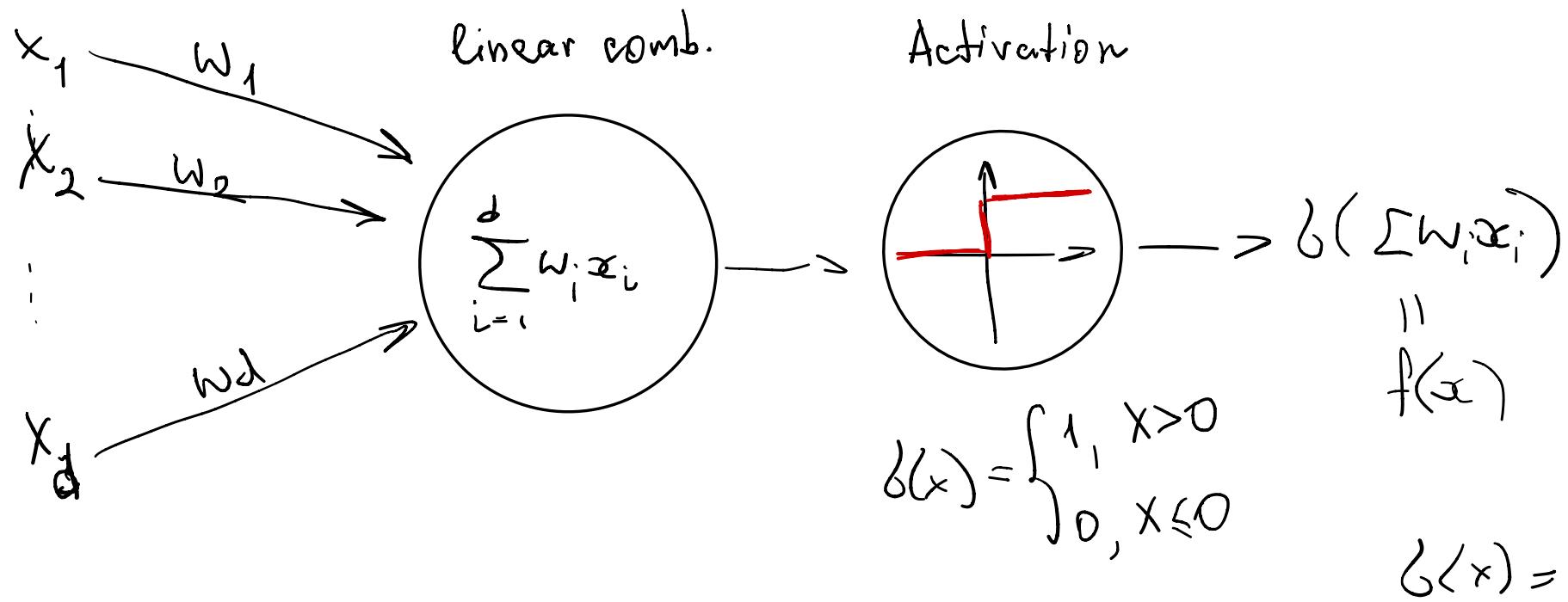
But we need to pay attention to the constraints!

Artificial Neuron: Inspiration



Artificial Neuron: Perceptron

$$x = (x_1, \dots, x_d)$$



$\delta(x) = x \Rightarrow$ Linear Regn

$\delta(x) = \text{sigm}(x) \Rightarrow$ Log. regn

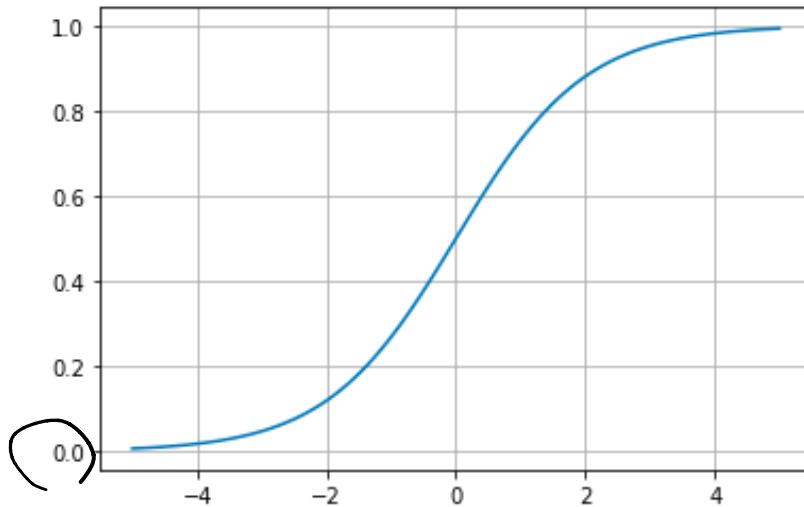
Artificial Neuron: Logistic Regression

$f(x) = \text{sigm}(\omega^T x) \Rightarrow \text{Artificial Neuron} \equiv \text{Log. Reg}$

Artificial Neuron: Activation Functions

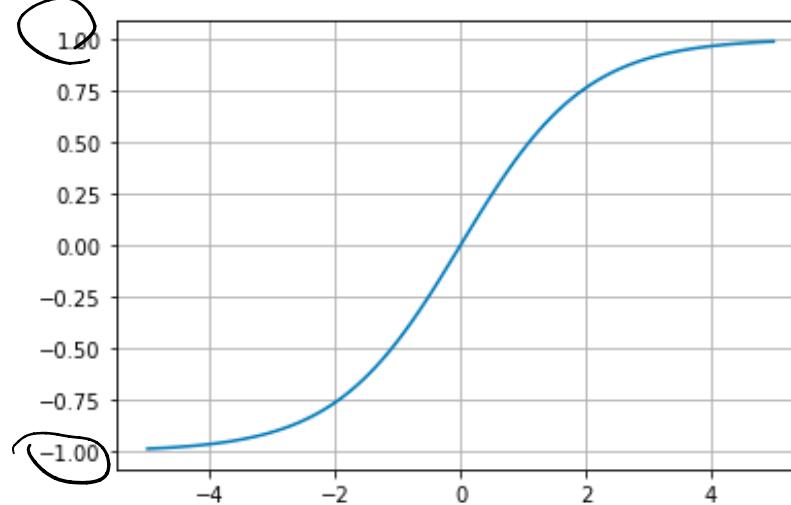
v

sigmoid



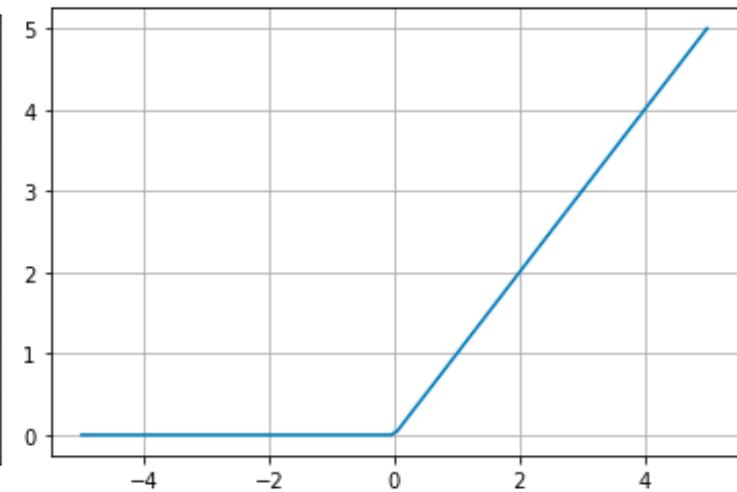
v

tanh

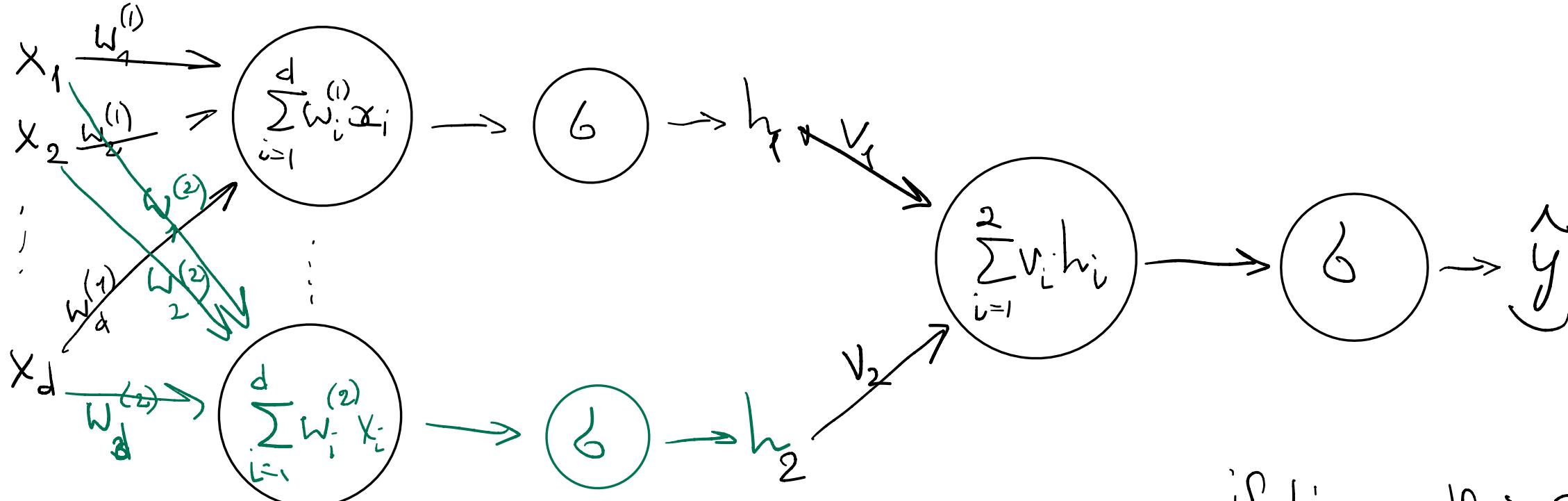


r

ReLU

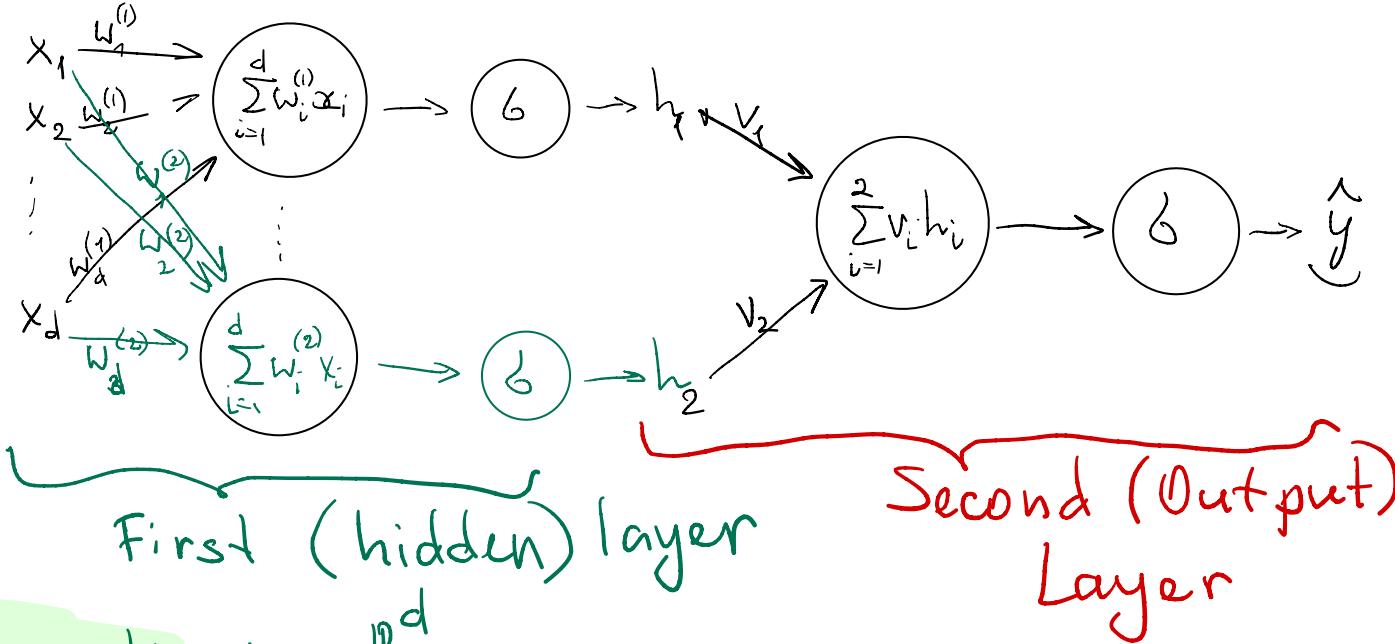


Artificial Neural Network (sceme)



if binary off \Rightarrow sigmoid
if regression $\Rightarrow f(x) = x$

Artificial Neural Network (formula)



Input: $x \in \mathbb{R}^d$

Weight: $w^{(1)} \in \mathbb{R}^d, w^{(2)} \in \mathbb{R}^d$

Act. func.: g

Output:

$$h_1 = g\left(\sum w_i^{(1)} x_i\right) = g(w^{(1)T} x)$$

$$h_2 = g\left(w^{(2)T} x\right)$$

Input: $h \in \mathbb{R}^2$
 Weight $V \in \mathbb{R}^2$
 Act. fun.: g

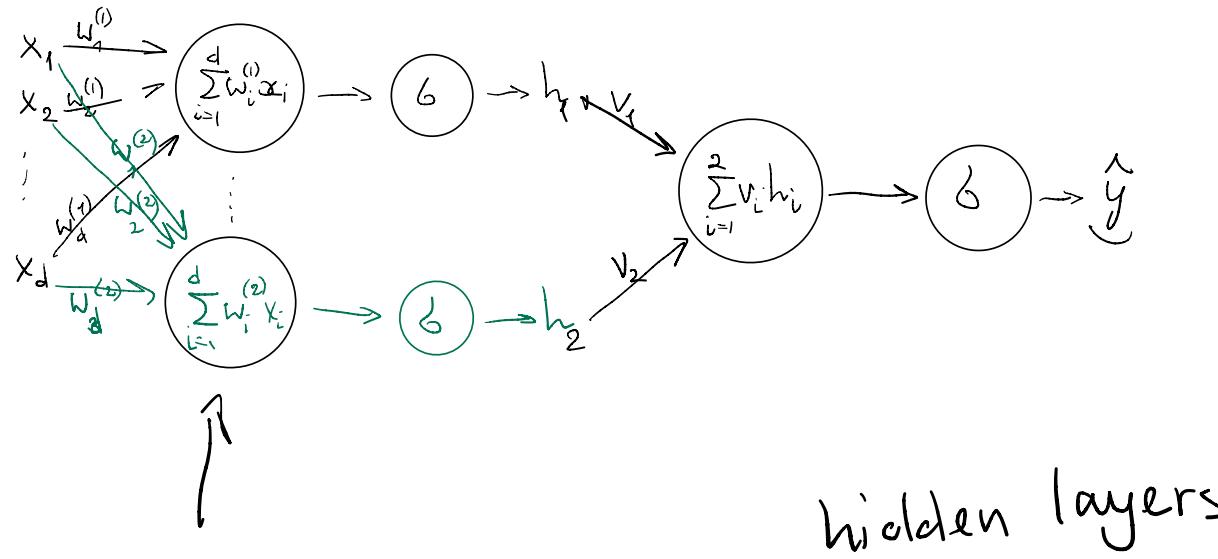
$$\hat{y} = g(V^T h)$$

$$h = \begin{bmatrix} g(w^{(1)T} x) \\ g(w^{(2)T} x) \end{bmatrix} = g\left(\begin{bmatrix} w^{(1)T} x \\ w^{(2)T} x \end{bmatrix}\right) = g(W^T x)$$

$$W = (w^{(1)T}, w^{(2)T}) \in \mathbb{R}^{d \times 2}$$

$$\hat{y} = g(W^T x)$$

Artificial Neural Network



$$\underbrace{f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1}_{\text{last "layer"} }(x)$$

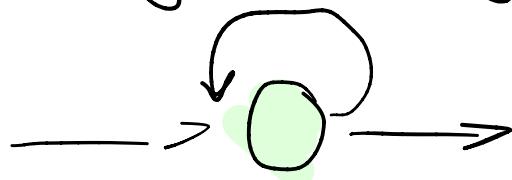
Layer 1: $f_1(x) = G(W^T x) \in \mathbb{R}^{\textcircled{2}}$

Layer 2: $f_2(x) = G(V^T x) \in \mathbb{R}^{\textcircled{1}}$

$$\hat{y} = f_2(f_1(x)) = \underbrace{f_2 \circ f_1}_{\text{composition}}(x)$$

composition
of functions

Layers in Artificial Neural Network

- Dense $f(x) = W^T x$, $W \in \mathbb{R}^{d \times h}$ \nwarrow hidden dimension
- Non-linearity
$$f(x) = \begin{bmatrix} g(x_1) \\ \vdots \\ g(x_d) \end{bmatrix}$$
- Convolutional $f(x) = W^T x$, but W should have special form
 \nwarrow Usually for images
- Recurrent 

How to train?

$$p(y|x, \theta) = p(y|f_\theta(x))$$

$$\hat{\theta}^* = \arg \max_{\theta} \prod_{n=1}^N p(y_n | f_\theta(x_n))$$

model param of some distr

or

$$\hat{\theta}^* = \arg \min_{\theta} L(y, x, \theta)$$

1) $X = \{x_1, \dots, x_N\}$, $Y = \{y_1, \dots, y_N\}$ are large

$N \gg$
Stochastic
g.d.

2) f_θ is a neural network
 \leftarrow Chain Rule
Backprop

SGD recap

X, Y - dataset with N observations
 $L(X, Y, \theta)$ ← loss

Initialize $\theta^{(0)}$

For i in max_iter :

- Sample mini-batch \tilde{X}, \tilde{Y} (m points, $m \ll N$)
- Compute $\nabla_{\theta} L(\tilde{X}, \tilde{Y}, \theta)$
- Update $\theta^{(i)} = \theta^{(i-1)} - \alpha \cdot \nabla_{\theta} L(\tilde{X}, \tilde{Y}, \theta)$

$$\left[\frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_{10000}} \right]$$

Chain Rule Recap

$$F(x) = \underbrace{f \circ g(x)}_{= f(g(x))}, \quad ? \quad \frac{\partial F(x)}{\partial x}$$

$$\frac{\partial F(x)}{\partial x} = \frac{\partial f(g(x))}{\partial g(x)} \cdot \frac{\partial g(x)}{\partial x}$$

$$\left. \begin{array}{l} h = g(x) \\ F(x) = f(h) \end{array} \right\} \Rightarrow \frac{\partial F(x)}{\partial x} = \frac{\partial f}{\partial h} \cdot \overbrace{\frac{\partial g(x)}{\partial x}}^{\text{L}}$$

Chain Rule in NN

$$y_L = f_L \circ f_{L-1} \circ \dots \circ f_1(x)$$

↑
 θ_L ↑
 θ_{L-1} ↑
 θ^1

• θ^L : $\frac{\partial L(y_L, x)}{\partial \theta^L} = \frac{\partial L}{\partial y_L} \cdot \frac{\partial y_L}{\partial \theta^L}$

• θ^{L-1} : $\frac{\partial L}{\partial \theta^{L-1}} = \frac{\partial L}{\partial y_L} \cdot \frac{\partial y_L}{\partial \theta^{L-1}} = \frac{\partial L}{\partial y_L} \cdot \frac{\partial y_L}{\partial y_{L-1}} \cdot \frac{\partial y_{L-1}}{\partial \theta^{L-1}}$

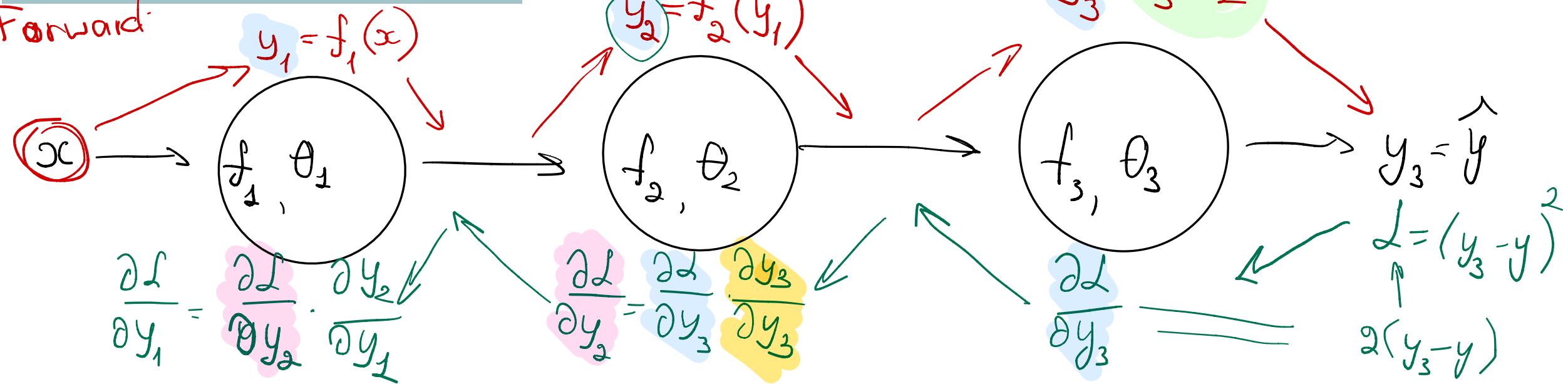
• θ^{L-2} : $\frac{\partial L}{\partial \theta^{L-1}} = \frac{\partial L}{\partial y_L} \cdot \frac{\partial y_L}{\partial y_{L-1}} \cdot \frac{\partial y_{L-1}}{\partial y_{L-2}} \cdot \frac{\partial y_{L-2}}{\partial \theta^{L-2}}$

$$\nabla_{\theta} L(y_L, x)$$

① $y_{L-1} = f_{L-1}(x, \theta_{L-1})$
 $y_L = f_L(y_{L-1})$

Backpropagation

Forward:



$$\frac{\partial L}{\partial \theta_1} = \frac{\partial L}{\partial y_1} \cdot \frac{\partial y_1}{\partial \theta_1}$$

$$\frac{\partial L}{\partial \theta_2} = \frac{\partial L}{\partial y_2} \cdot \frac{\partial y_2}{\partial \theta_2}$$

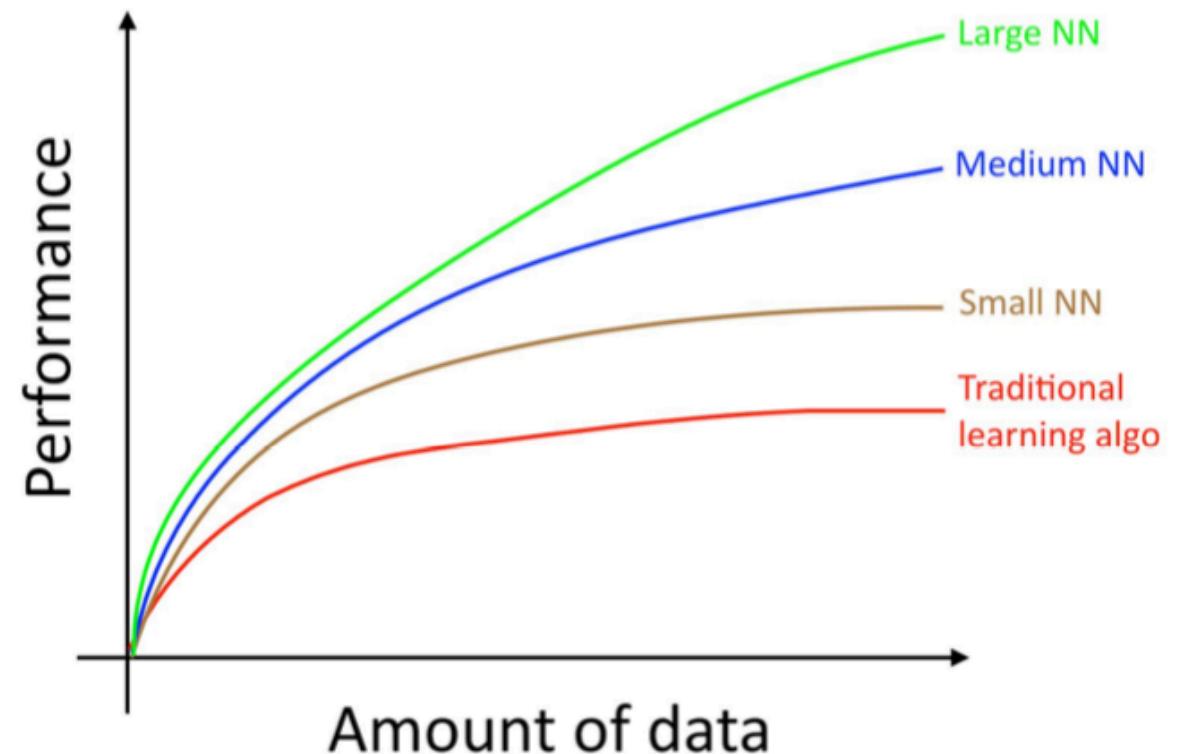
$$\frac{\partial L}{\partial \theta_3} = \frac{\partial L}{\partial y_3} \cdot \frac{\partial y_3}{\partial \theta_3}$$

SGD \Rightarrow Update θ

$\frac{\partial L}{\partial \theta_1}$
$\frac{\partial L}{\partial \theta_2}$
$\frac{\partial L}{\partial \theta_3}$

Deep Neural Networks

- In theory can approximate any arbitrary function
- Benefit from large datasets
- Can be implemented fast on GPU



Tasks

- Supervised: regression and classification
- Representation Learning
- Image / Text generation (next lecture)

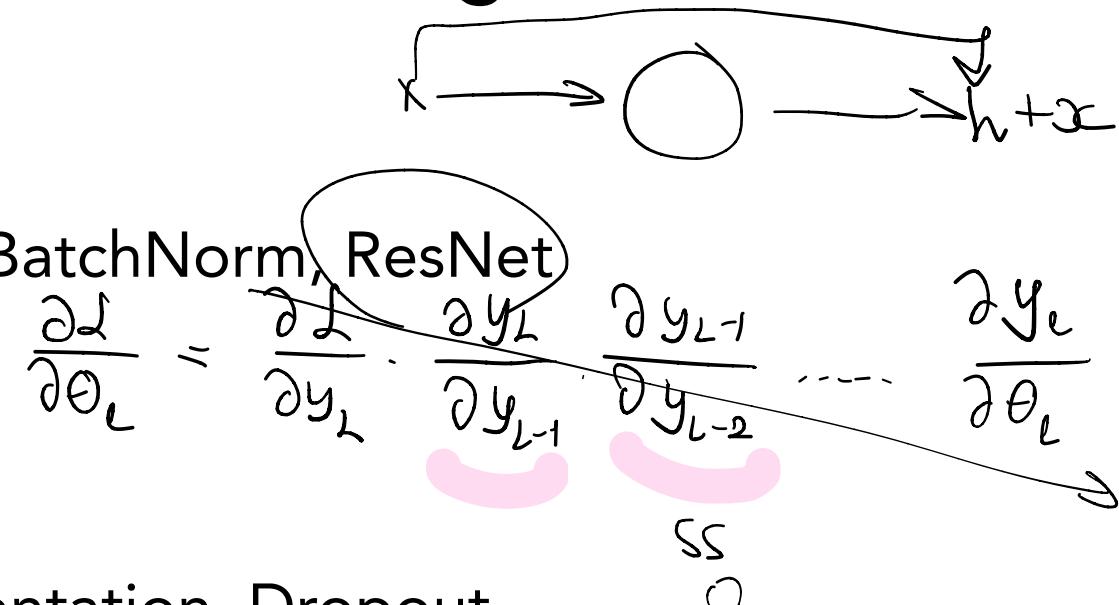
$$p(y|x)$$

$$p(x) \Rightarrow \tilde{x} \sim p(x)$$

Limitations and Challenges

- Vanishing Gradient

Other activation functions, BatchNorm, ResNet



- Overfitting

Regularization, Data Augmentation, Dropout

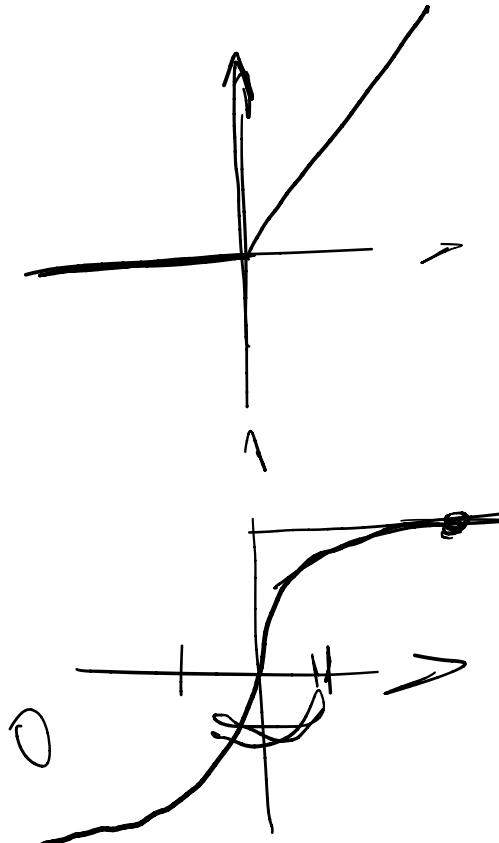
- Requires Large Dataset
Transfer Learning

D - # of params

$$D \gg N \gg 0$$

$$J + \frac{1}{2} \|\theta\|^2$$

↑ weight decay



Limitations and Challenges

- Vanishing Gradient
Other activation functions, BatchNorm, ResNet
- Overfitting
Regularization, Data Augmentation, Dropout
- Requires Large Dataset
Transfer Learning