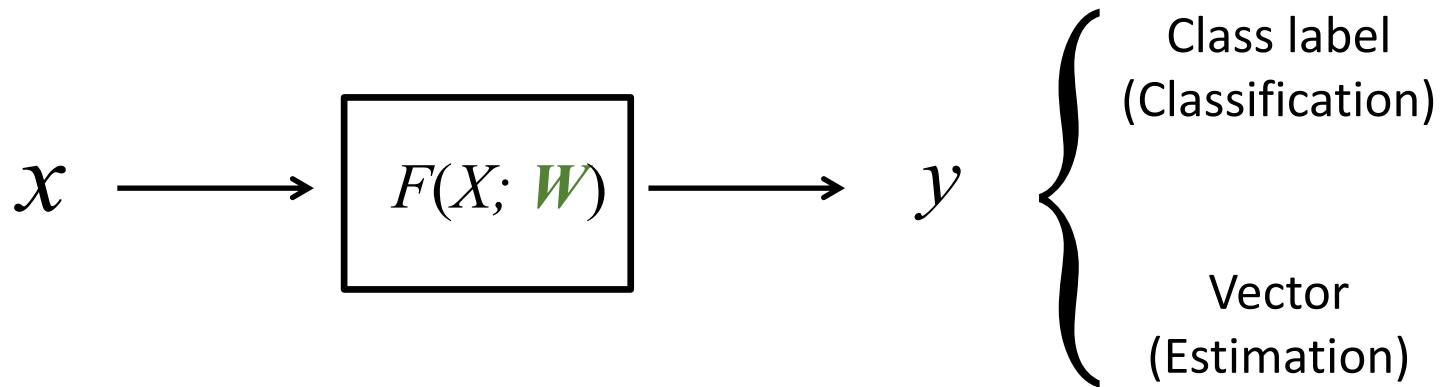


Basics of neural network

Outline

- Basic operations of deep model
- Basics of optimization
- Optimization for deep models

Machine Learning



Object recognition

{dog, cat, horse, flower, ...}



Super resolution

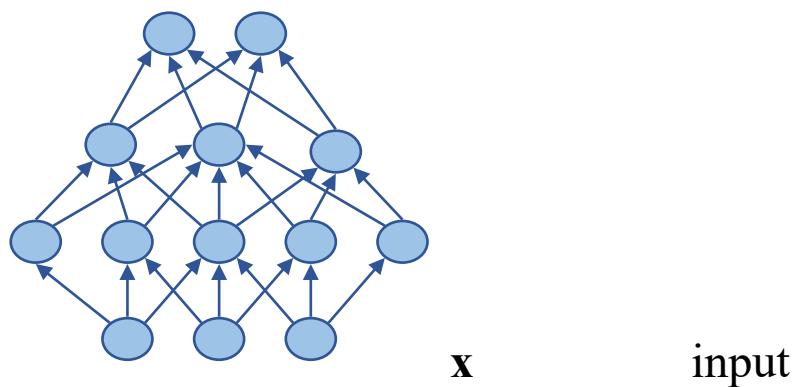


High-resolution
image

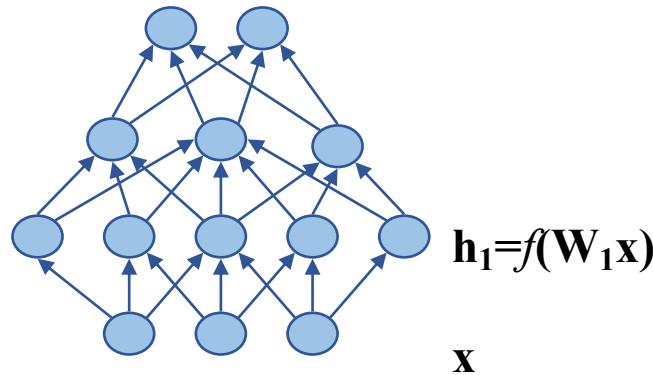
Low-resolution image

Neural network

Multilayer Perception (MLP)



Neural network

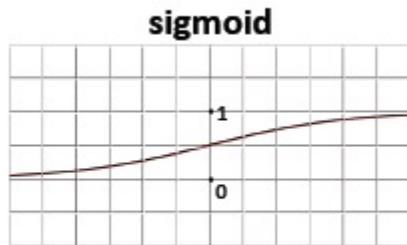


f : Non-linear/activation function

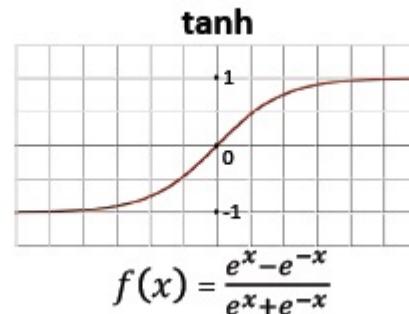
sigmoid: The largest response is 1, the smallest response is 0.

tanh: The largest response is 1, the smallest response is -1

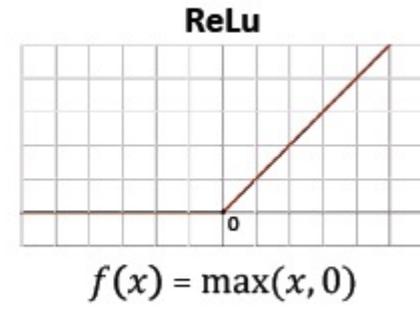
Relu (rectified linear unit): Thresholding, when smaller than 0, set to 0.



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



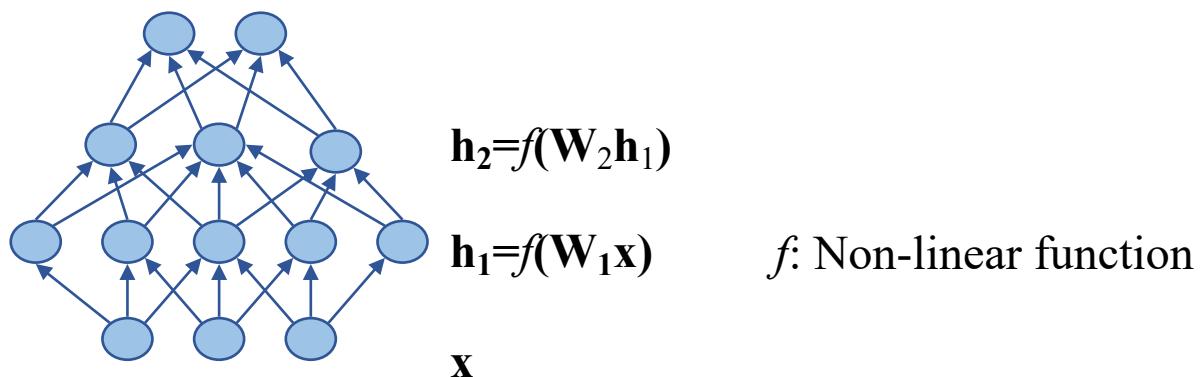
$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



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

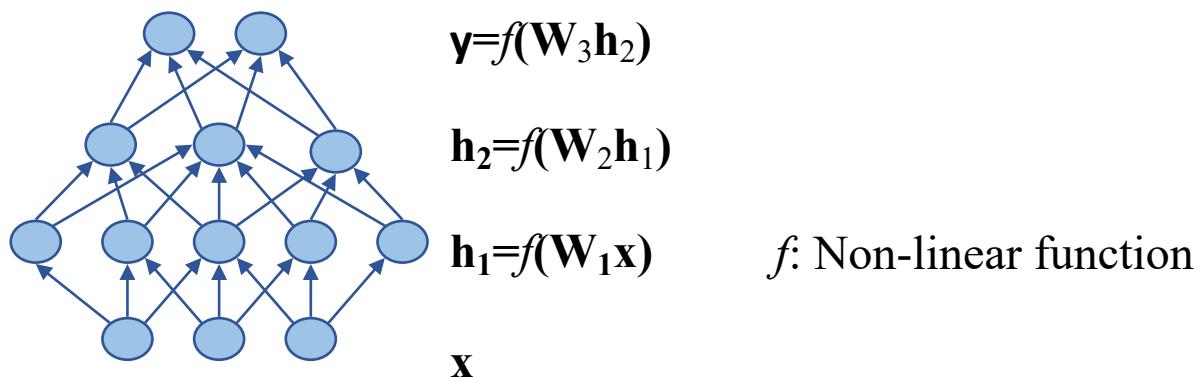
Neural network

Multilayer Perception (MLP)



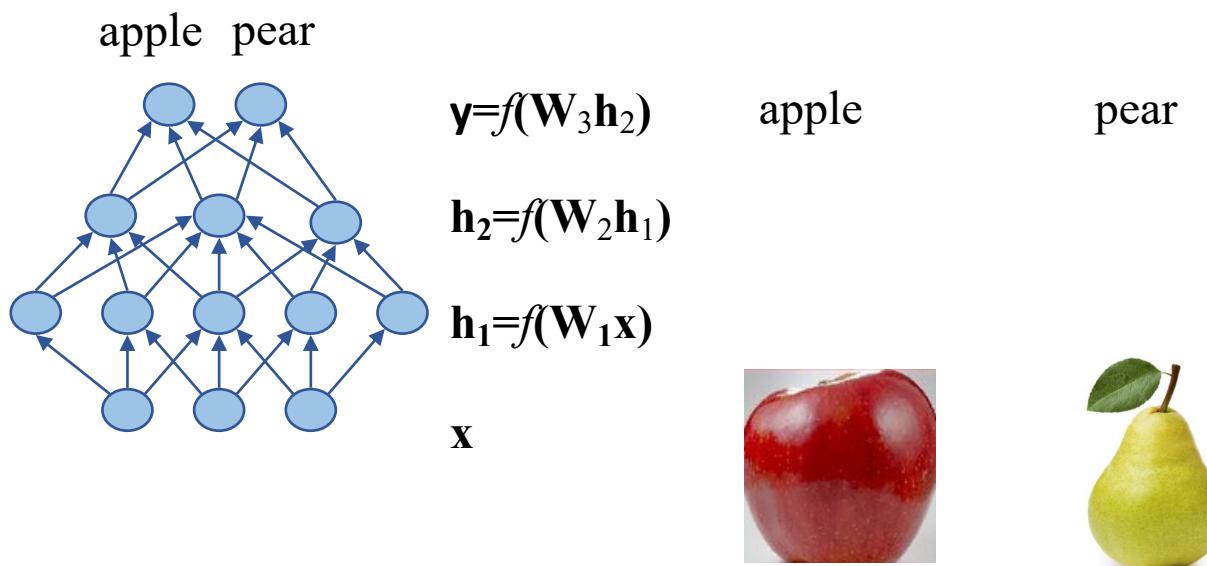
Neural network

Multilayer Perception (MLP)



Neural network

Multilayer Perception (MLP)



Neural
network

Back
propagation



Geoffrey Hinton



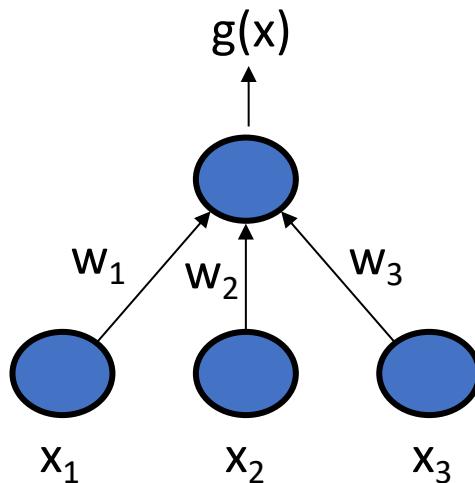
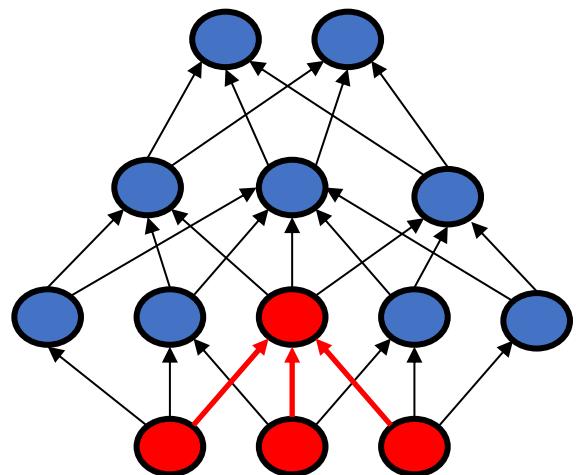
Nature

1940s

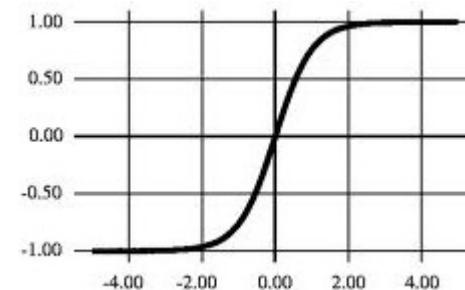
1986



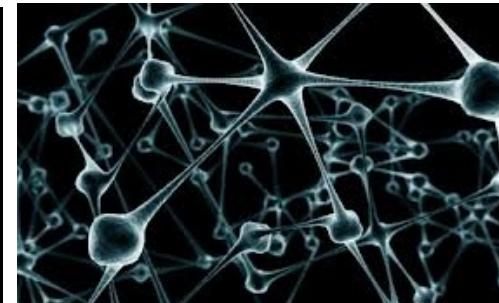
Lecun Yann joined Geoffrey Hinton's group in 1987

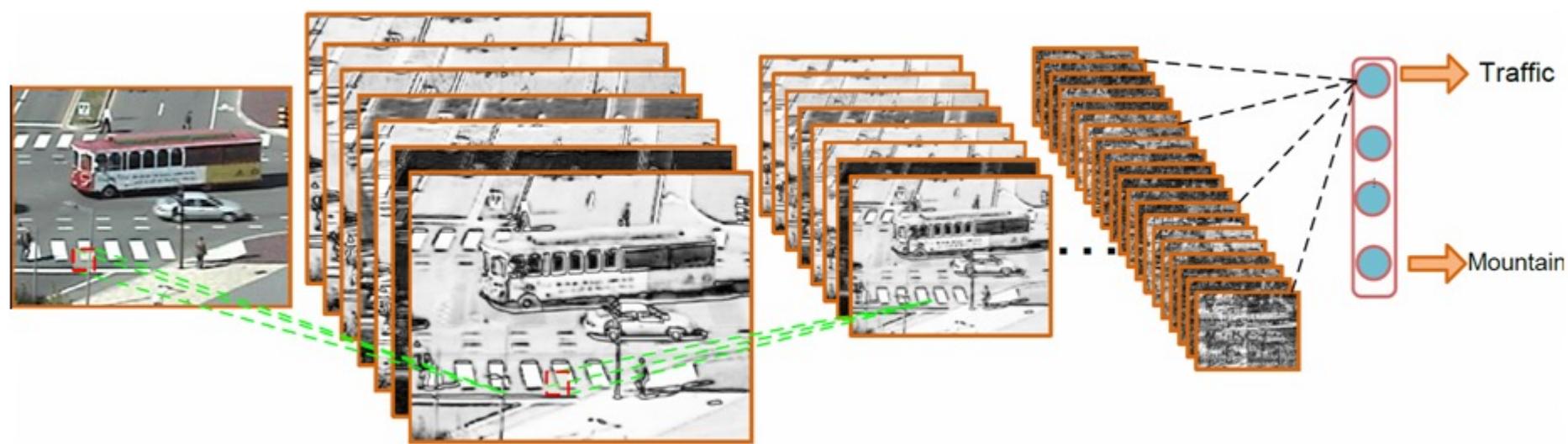
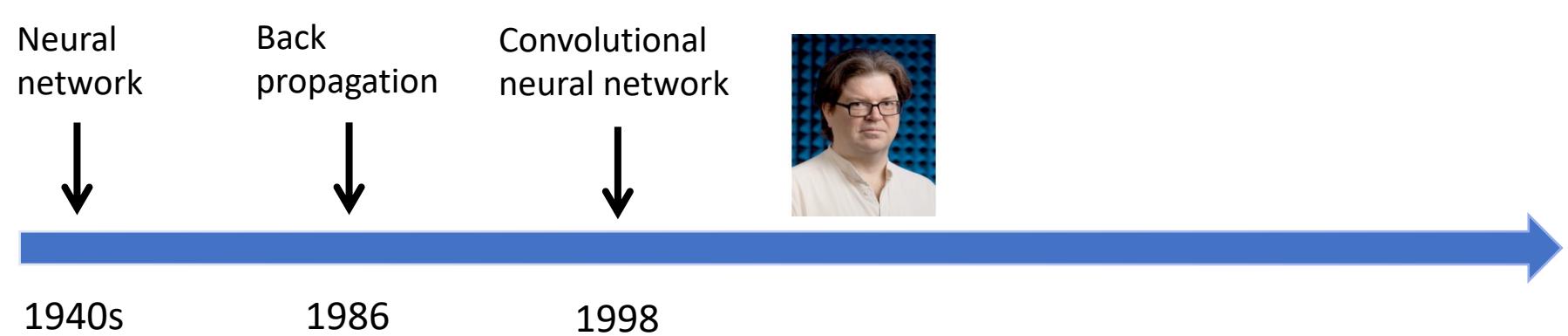


$$g(\mathbf{x}) = f\left(\sum_{i=1}^d x_i w_i + w_0\right) = f(\mathbf{w}^t \mathbf{x})$$



$f(\text{net})$





Neural
network

Back
propagation

Convolutional
neural network

Deep
belief net

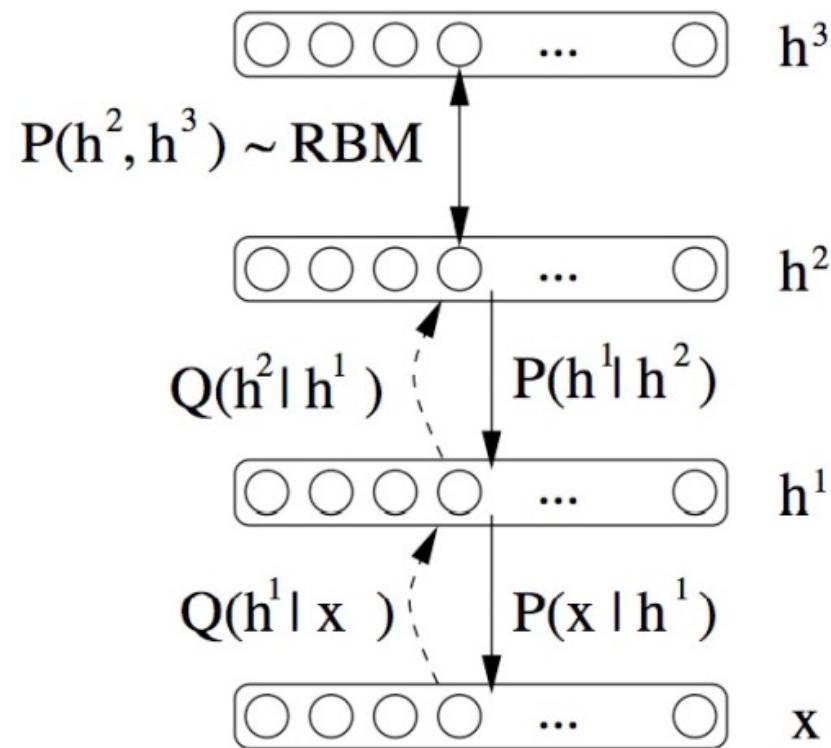


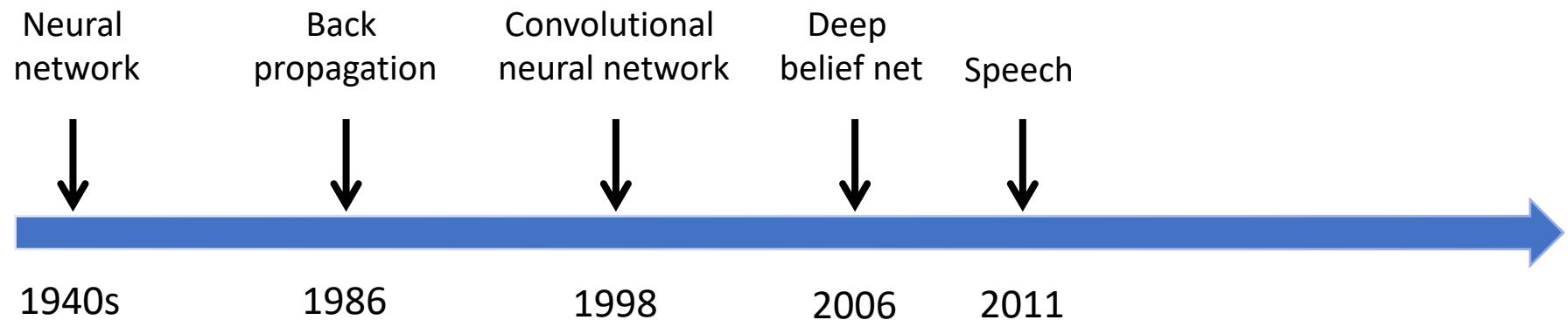
1940s

1986

1998

2006





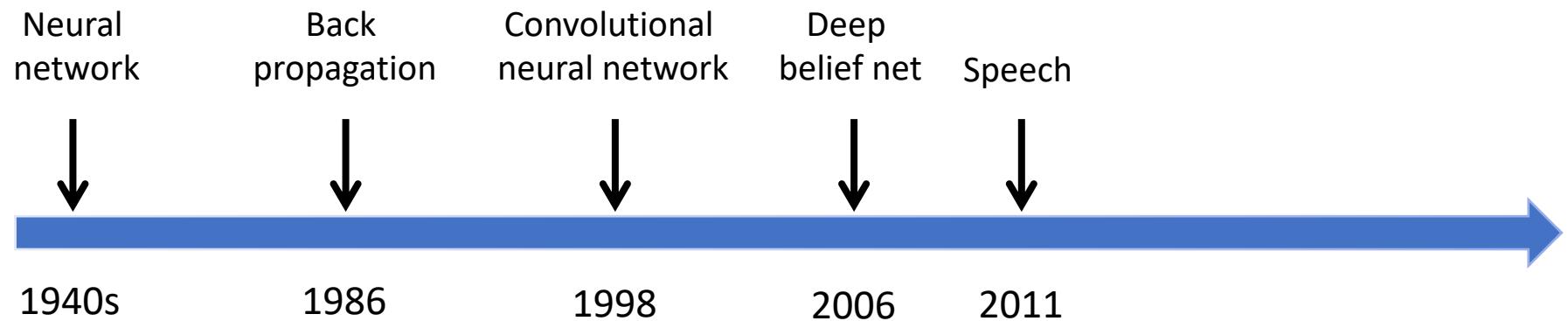
deep learning results

task	hours of training data	DNN-HMM	GMM-HMM with same data
Switchboard (test set 1)	309	18.5	27.4
Switchboard (test set 2)	309	16.1	23.6
English Broadcast News	50	17.5	18.8
Bing Voice Search (Sentence error rates)	24	30.4	36.2
Google Voice Input	5,870	12.3	
Youtube	1,400	47.6	52.3

Deep Networks Advance State of Art in Speech

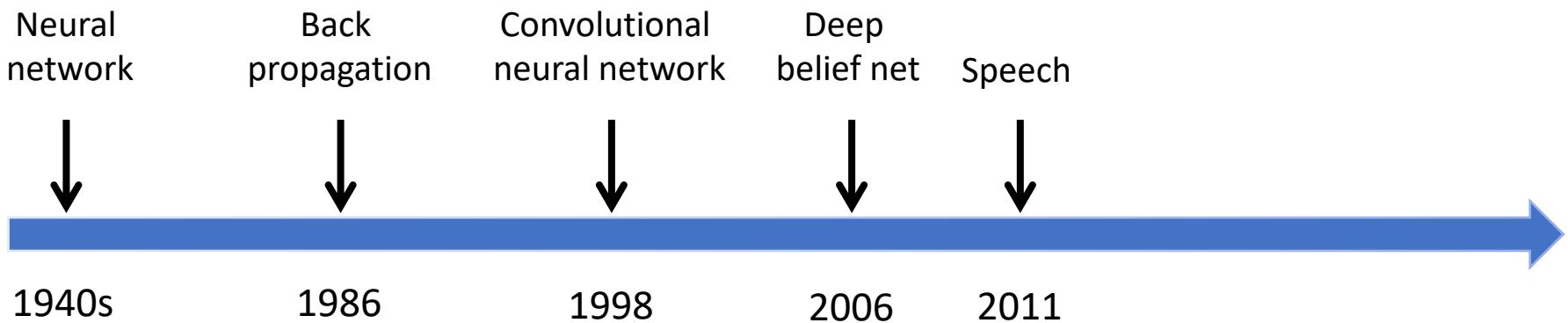
Deep Learning leads to breakthrough in speech recognition at MSR.





Not well accepted by the vision community 😞





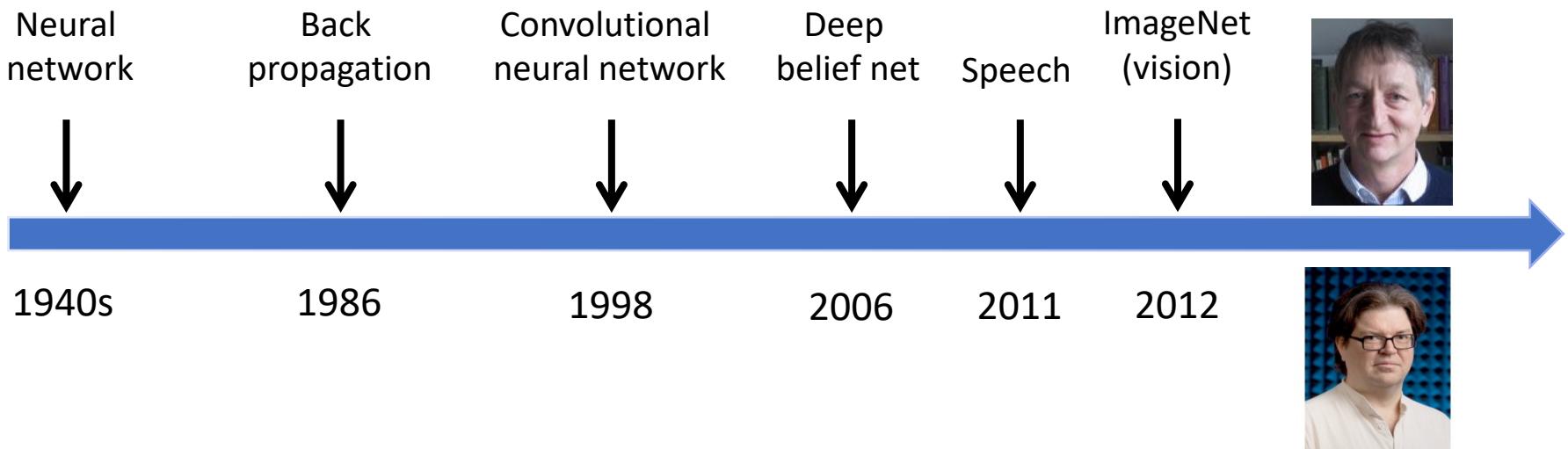
LeCun's open letter in CVPR 2012



So, I'm giving up on submitting to computer vision conferences altogether. CV reviewers are just too likely to be clueless or hostile towards our brand of methods. Submitting our papers is just a waste of everyone's time (and incredibly demoralizing to my lab members)

I might come back in a few years, if at least two things change:

- Enough people in CV become interested in feature learning that the probability of getting a non-clueless and non-hostile reviewer is more than 50% (hopefully [Computer Vision Researcher]'s tutorial on the topic at CVPR will have some positive effect).
- CV conference proceedings become open access.



Rank	Name	Error rate	Description
1	U. Toronto	0.15315	Deep learning
2	U. Tokyo	0.26172	Hand-crafted features and learning models. Bottleneck.
3	U. Oxford	0.26979	
4	Xerox/INRIA	0.27058	

Object recognition over 1,000,000 images and 1,000 categories (2 GPU)

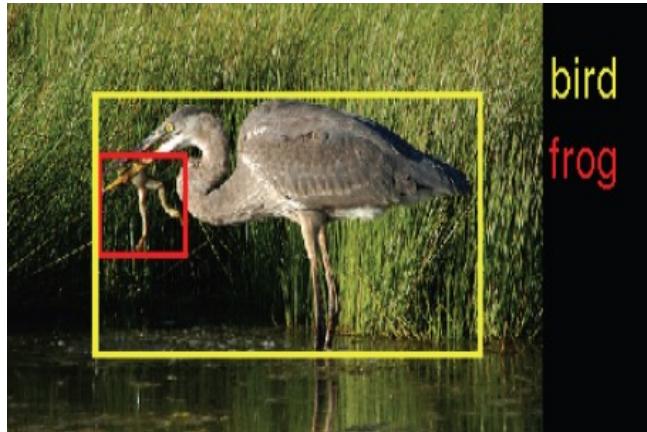
Current best result < 0.03

ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

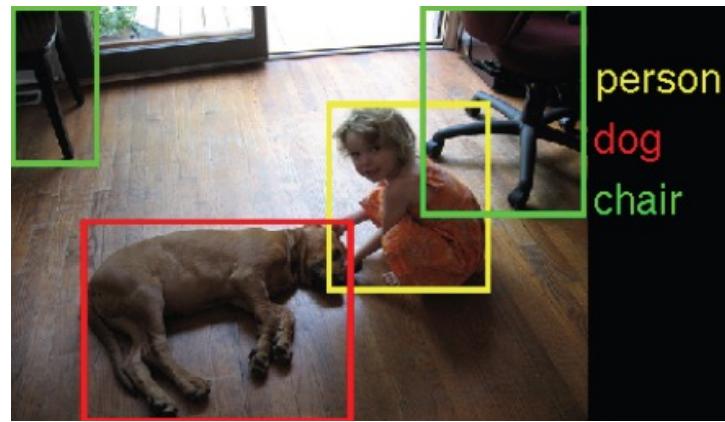


ImageNet Object Detection Task

- 200 object classes
- 60,000 test images



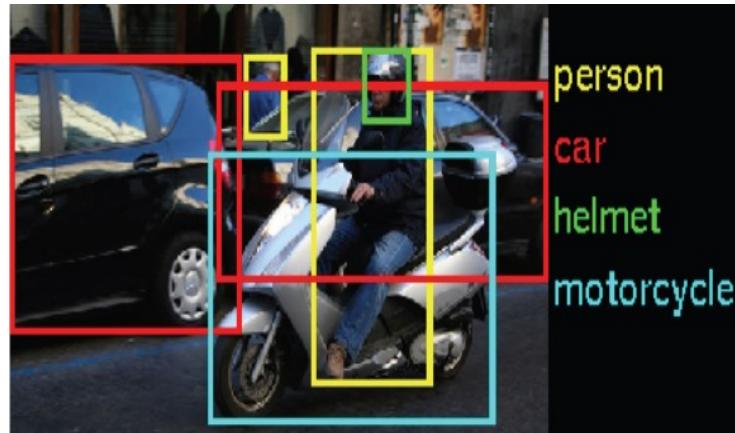
bird
frog



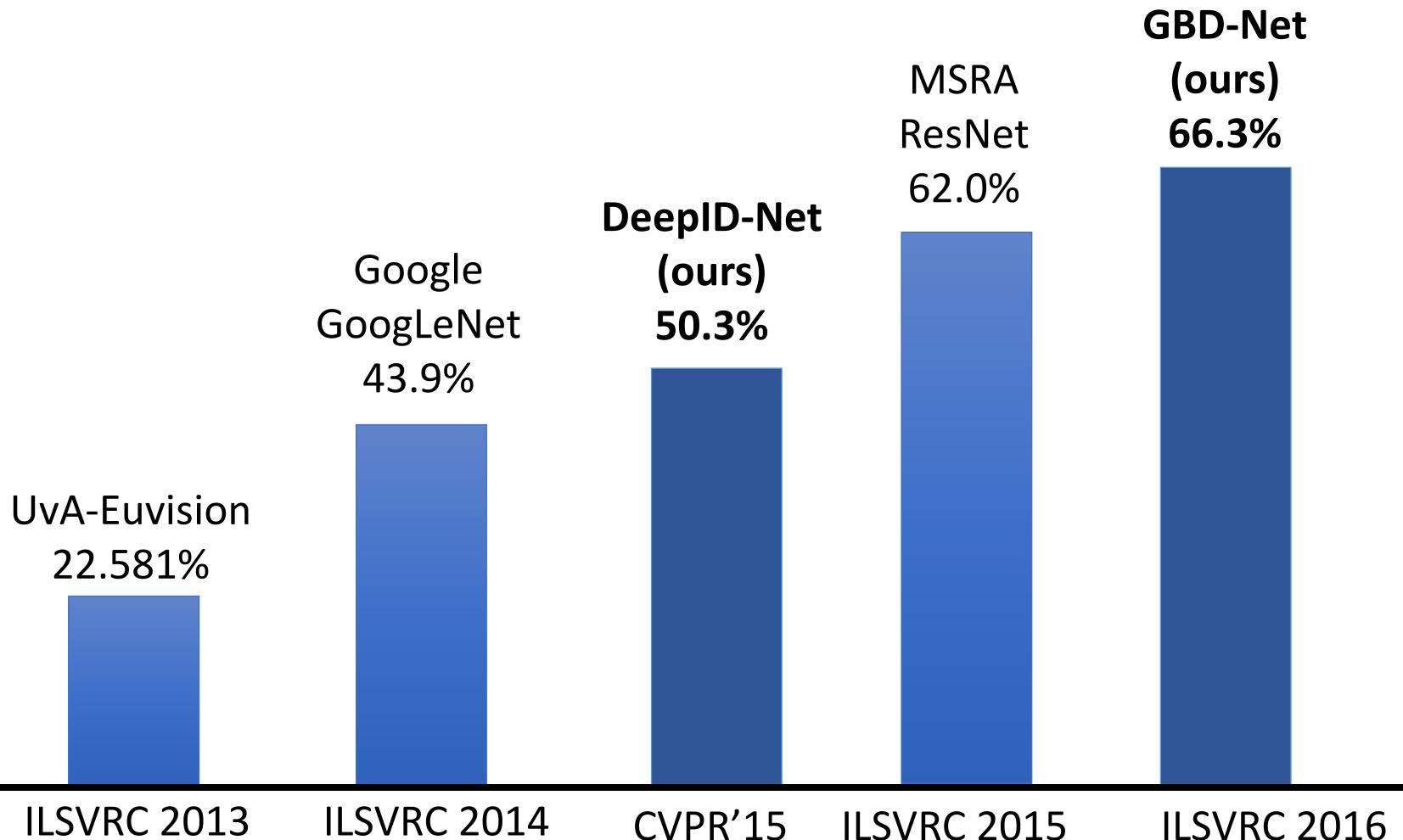
person
dog
chair



person
hammer
flower pot
power drill



person
car
helmet
motorcycle

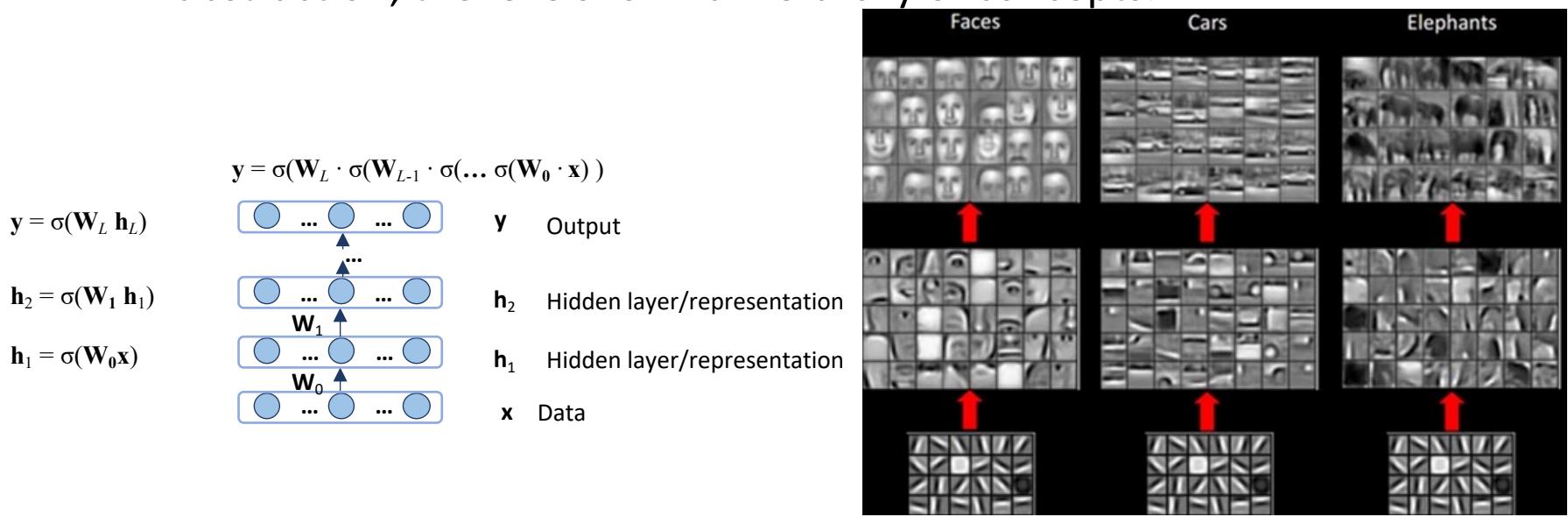


W. Ouyang and X. Wang, et al. "DeepID-Net: Deformable Deep Convolutional Neural Networks for Object Detection," CVPR 2015/TPAMI17

X. Zeng, W. Ouyang, B. Yang, J. Yan, and X. Wang, "Gated Bi-directional CNN for Object Detection," ECCV 2016 / TPAMI17

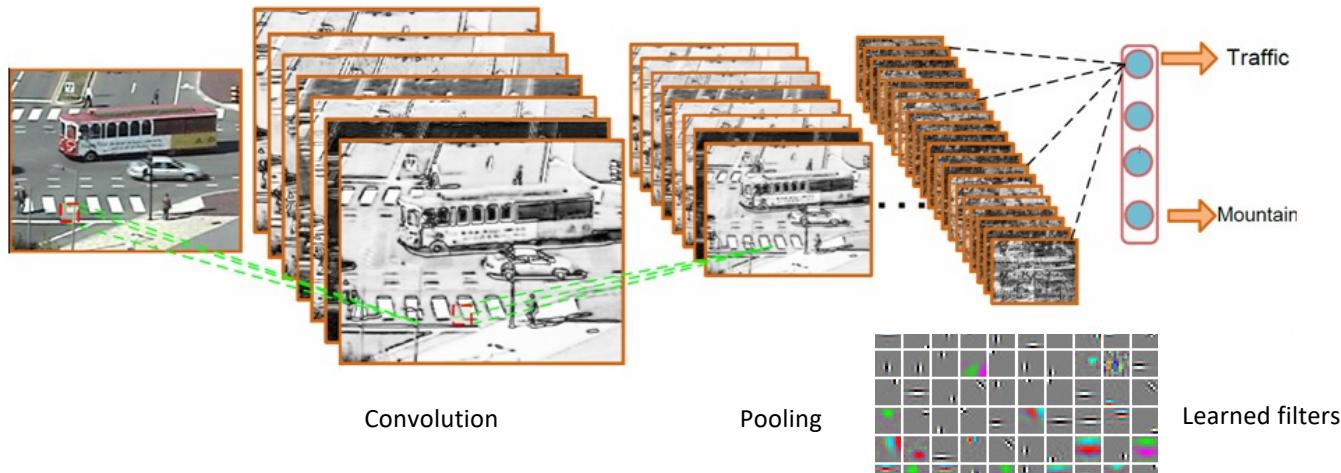
Deep learning

- Learning of multi-layer non-linear transformations
- Multiple layers have multiple levels of representations with multiple non-linear transformations.
- Higher level representations are derived from lower level representations to form a hierarchical representation.
- Multiple levels of representations correspond to different levels of abstraction; the levels form a hierarchy of concepts.



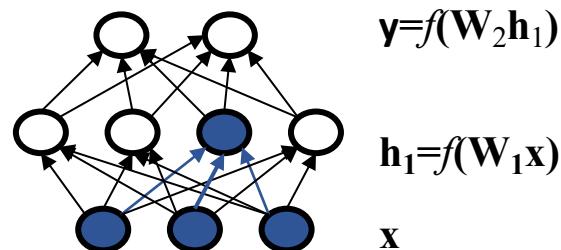
Classical Deep Models

- Convolutional Neural Networks (CNN)
 - First proposed by Fukushima in 1980
 - Improved by LeCun, Bottou, Bengio and Haffner in 1998



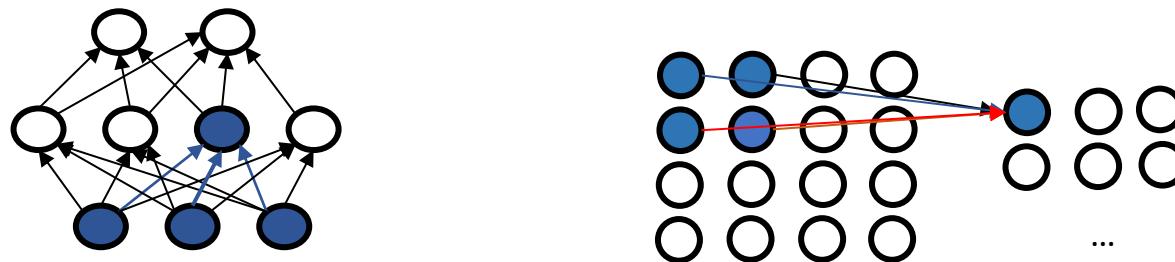
Comparison of fully connected network and convolutional neural network

- Fully connected network (FC)
 - Each neuron at the current layer is connected with all the neurons in the previous layer



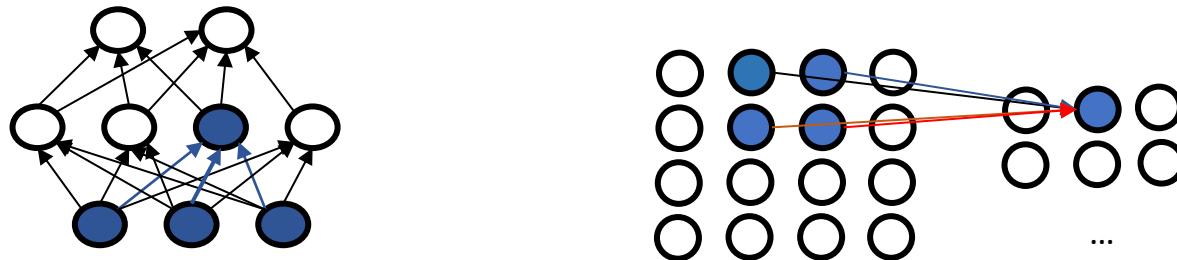
Comparison of fully connected network and convolutional neural network

- Fully connected network (FC)
 - Every output unit interacts with every input unit
- Convolutional neural network (CNN)
 - Every output unit interacts with a small set of input units, depending on location
 - Parameters are shared



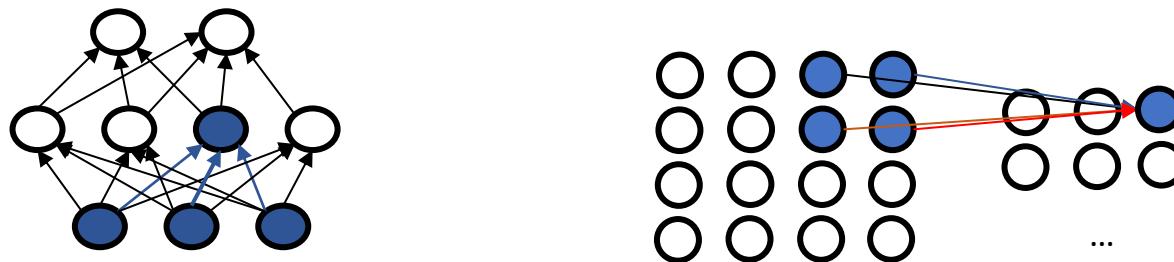
Comparison of fully connected network and convolutional neural network

- Fully connected network (FC)
 - Every output unit interacts with every input unit
- Convolutional neural network (CNN)
 - Every output unit interacts with a small set of input units, depending on location
 - Parameters are shared



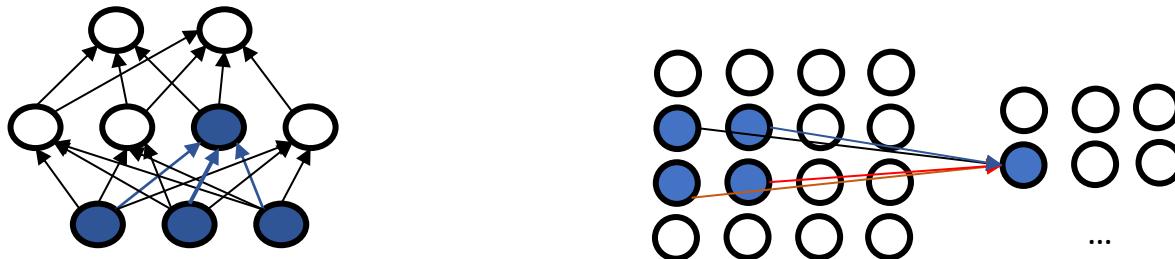
Comparison of fully connected network and convolutional neural network

- Fully connected network (FC)
 - Every output unit interacts with every input unit
- Convolutional neural network (CNN)
 - Every output unit interacts with a small set of input units, depending on location
 - Parameters are shared

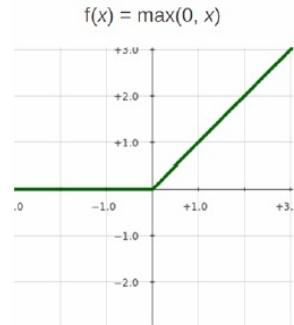


Comparison of fully connected network and convolutional neural network

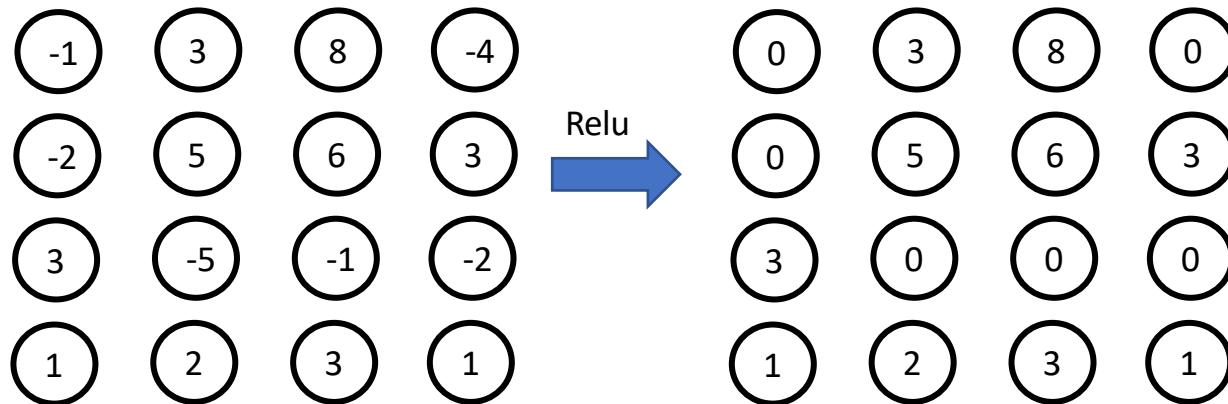
- Fully connected network (FC)
 - Every output unit interacts with every input unit
- Convolutional neural network (CNN)
 - Every output unit interacts with a small set of input units, depending on location
 - Parameters are shared



Non-linear function Relu

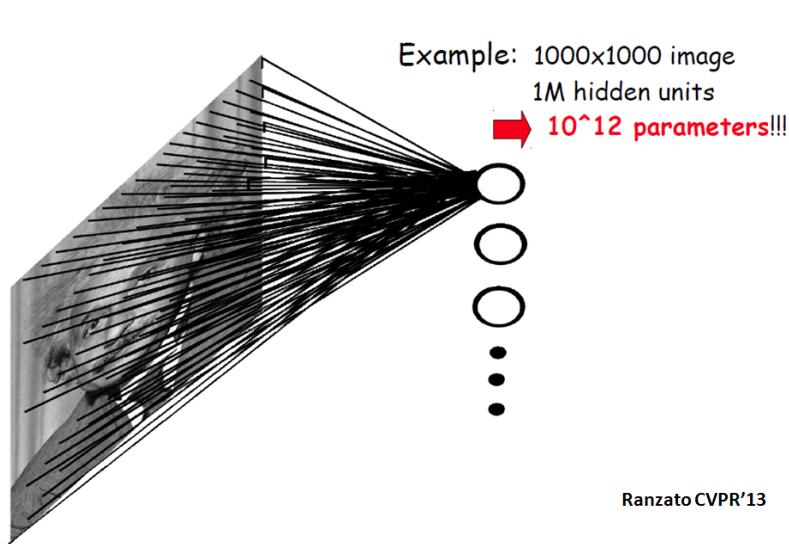


- Rectified linear unit (Relu)
- $f(x)= \max(x, 0);$



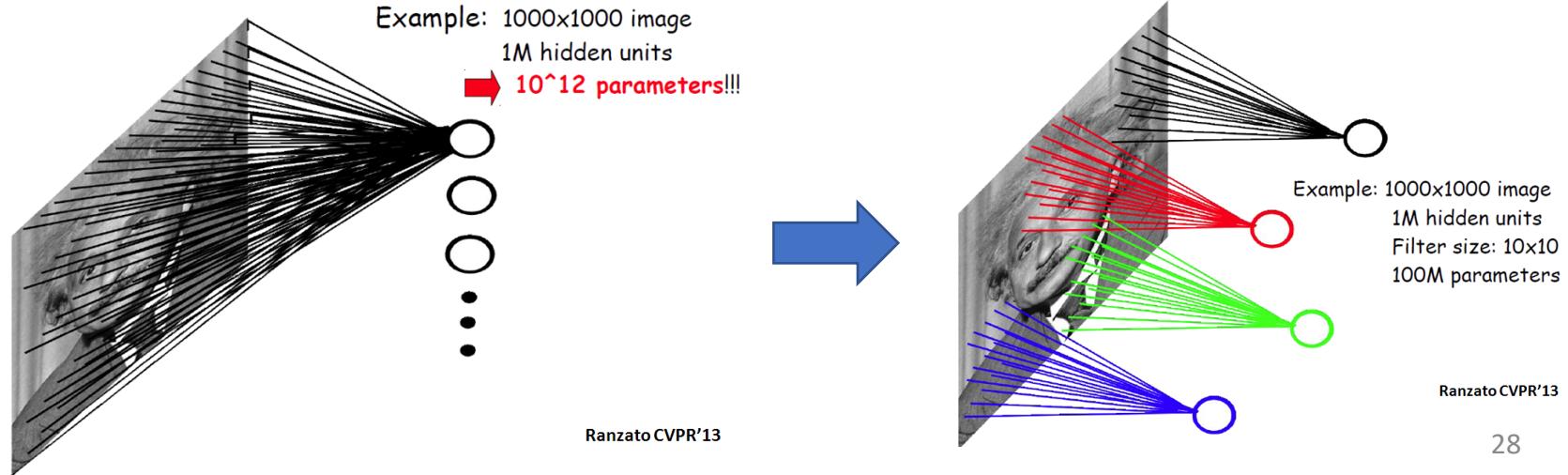
Problems of fully connected neural networks

- Every output unit interacts with every input unit
- The number of weights grows largely with the size of the input image
- Pixels in distance are less correlated



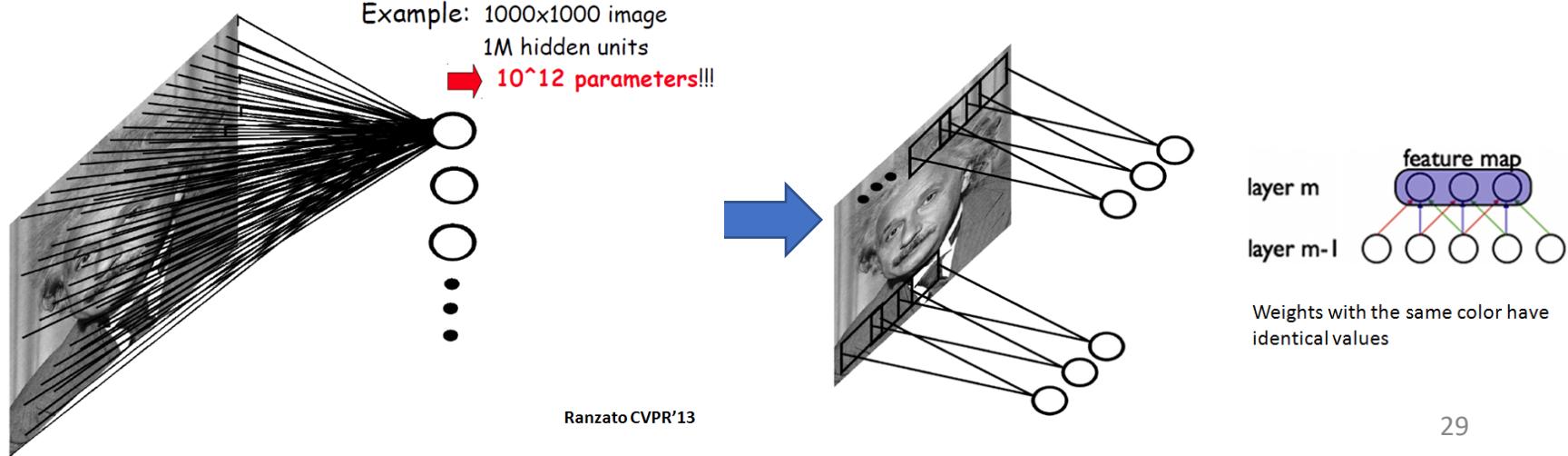
Locally connected neural networks

- Sparse connectivity: a hidden unit is only connected to a local patch (weights connected to the patch are called filter or kernel)
- It is inspired by biological systems, where a cell is sensitive to a small sub-region of the input space, called a receptive field. Many cells are tiled to cover the entire visual field.
- The design of such sparse connectivity is based on domain knowledge.



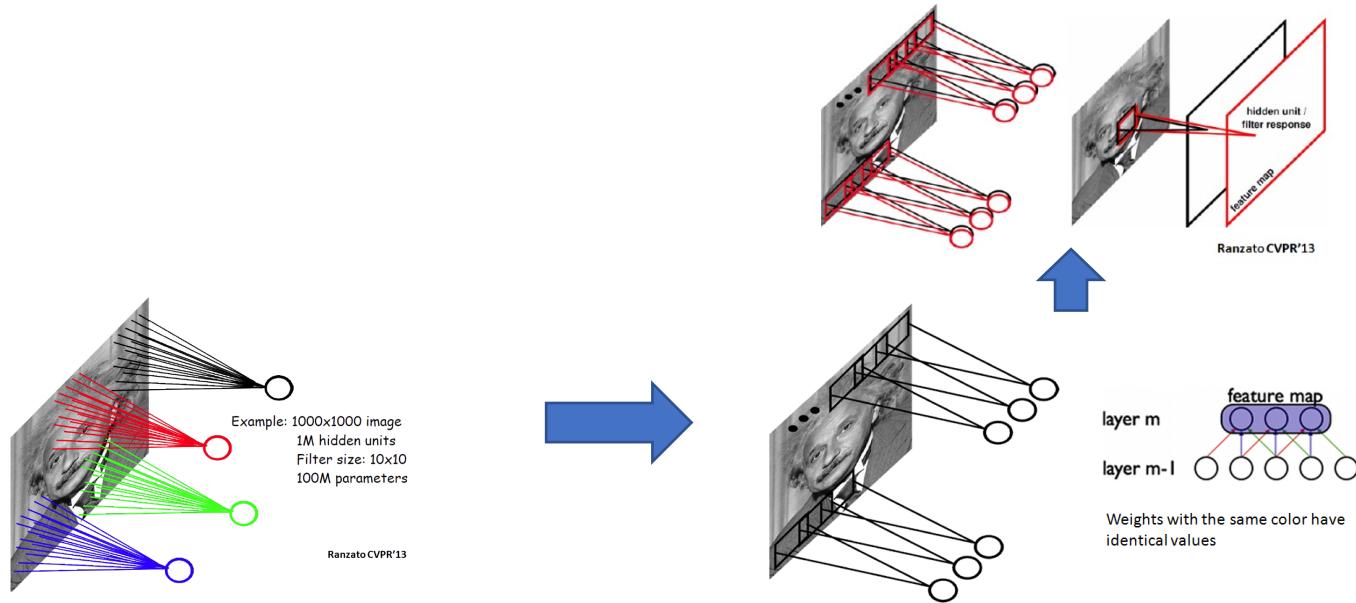
Locally connected neural networks with shared weights

- Translation invariance: capture statistics in local patches and they are independent of location
 - Similar edges may appear at different locations
- Hidden nodes at different locations share the same weights. It greatly reduces the number of parameters to learn
- In some applications (especially images with regular structures), we may only locally share weights or not share weights at top layers



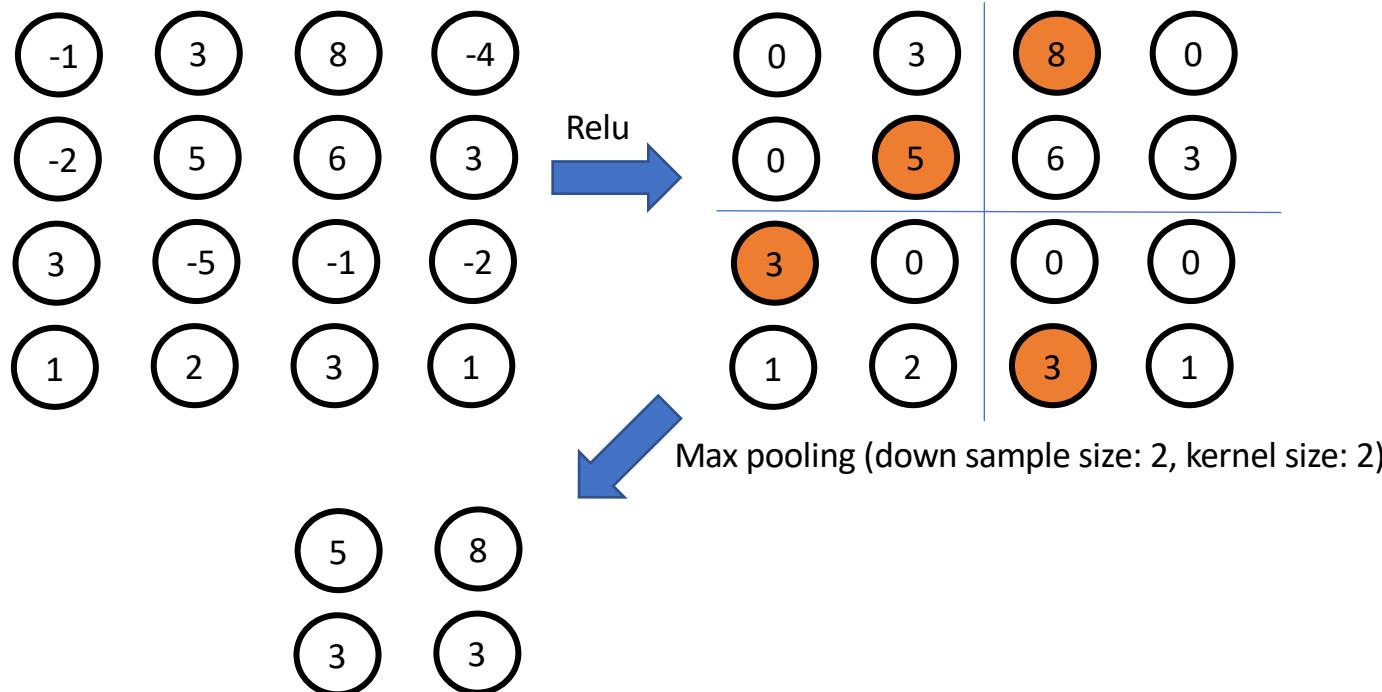
Multiple filters

- Multiple filters generate multiple feature maps
- Detect the spatial distributions of multiple visual patterns



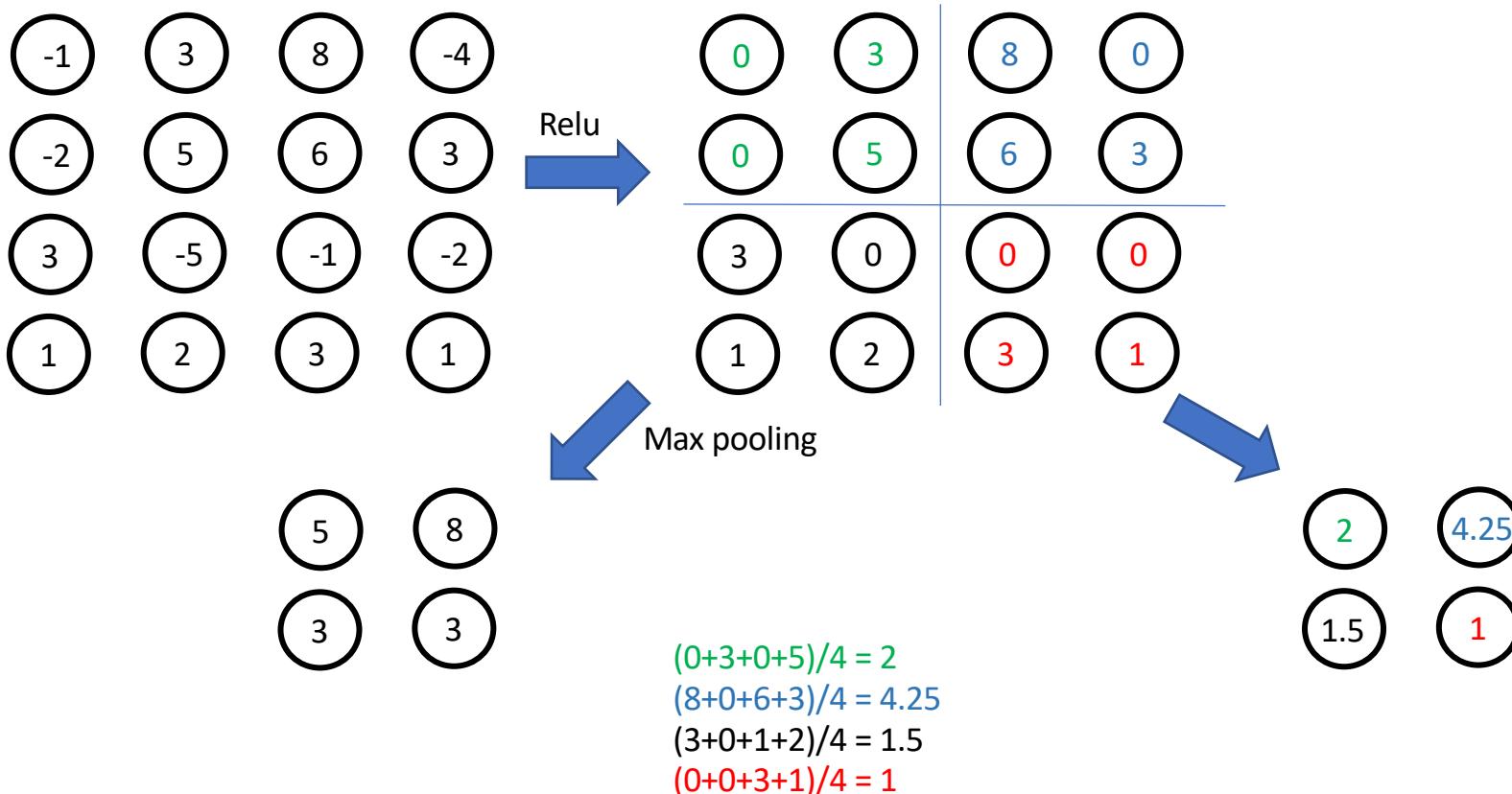
Pooling

- Max pooling



Pooling

- Max pooling and average pooling



Convolutional Neural Networks (CNN)

- Visualization results
 - <https://www.youtube.com/watch?v=AgkfIQ4IGaM>
 - <https://www.youtube.com/watch?v=xF5WNpTzCA>

Outline

- Basic operations of deep model
- **Basics of optimization**
- Optimization for deep models

Optimization basics

- General form

$$\begin{aligned} & \min f_0(x) \\ \text{s.t. } & f_i(x) \leq b_i, \quad i = 1, \dots, m \end{aligned}$$

$x = [x_1, \dots, x_n]^T \in \mathbf{R}^n$ optimization variables

$f_0 : \mathbf{R}^n \rightarrow \mathbf{R}$ objective function

$f_i : \mathbf{R}^n \rightarrow \mathbf{R}, i = 1, \dots, m$ constraint functions

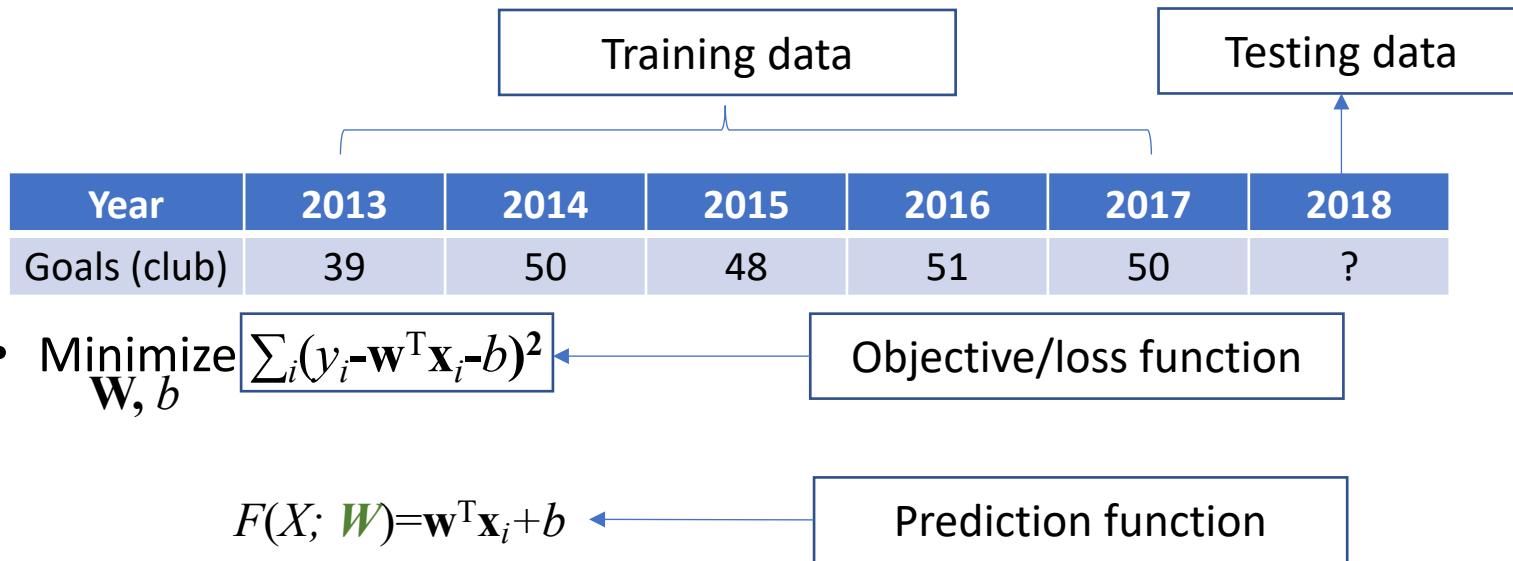
Our aim is to find an **optimal solution** x^* that minimizes f_0 whilst satisfying all the constraints.

Example:

maximize	$\mathbf{c}^T \mathbf{x}$
subject to	$A\mathbf{x} \leq \mathbf{b}$
and	$\mathbf{x} \geq \mathbf{0}$



Goals of Messi



Picture from <http://www.goal.com/en-gb/news/how-messi-became-the-worlds-deadliest-goalscorer/f1llyv71yqs91nzapr8yvo5td>

Goals from <http://messivsronaldo.net/calendar-year-stats/>

Terms in machine learning

Training data						Testing data
Year	2013	2014	2015	2016	2017	2018
Goals (club)	39	50	48	51	50	?

- Minimize $\sum_i(y_i - \mathbf{w}^T \mathbf{x}_i - b)^2$

Objective/loss function

- Training data:** The data used for training the model, i.e. learning the parameters.
- $(x_1, y_1) = (2013, 39), (x_2, y_2) = (2014, 50) \dots$
- \mathbf{W}, b : parameters to be learned from training data
- Testing data:** The data used for evaluating the model. How effective is the learned model?
- Minimize: $(39-w*2013-b)^2 + (50-w*2014-b)^2 + (48-w*2015-b)^2 + (51-w*2016-b)^2 + (50-w*2017-b)^2$

Linear Regression

Year	2013	2014	2015	2016	2017	2018
Goals (club)	39	50	48	51	50	?

- Minimize $\sum_i(y_i - \mathbf{w}^T \mathbf{x}_i - b)^2$

$x = [2013 \ 2014 \ 2015 \ 2016 \ 2017];$

$y = [39 \ 50 \ 48 \ 51 \ 50];$

`mdl = LinearModel.fit(x,y);`

`b = mdl.Coefficients.Estimate(1);`

`a = mdl.Coefficients.Estimate(2);`

$$y2 = a*x + b;$$

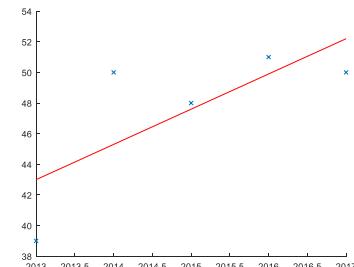
hold on;

`plot(x,y,'x');`

`plot(x,y2,'r');`

hold off;

$$a*2018+b$$



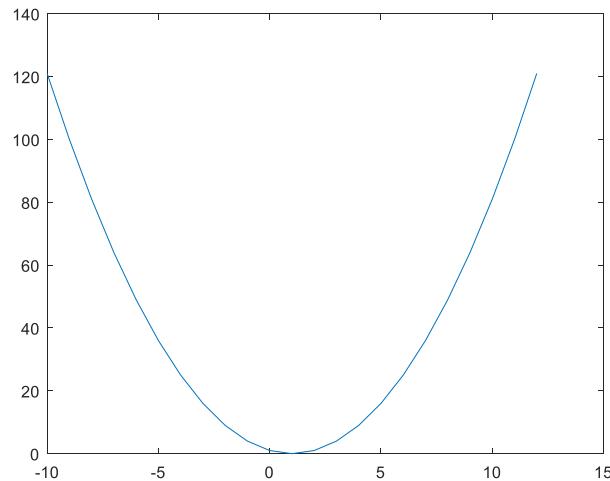
>> a*2018+b

ans =

$$54.5000$$

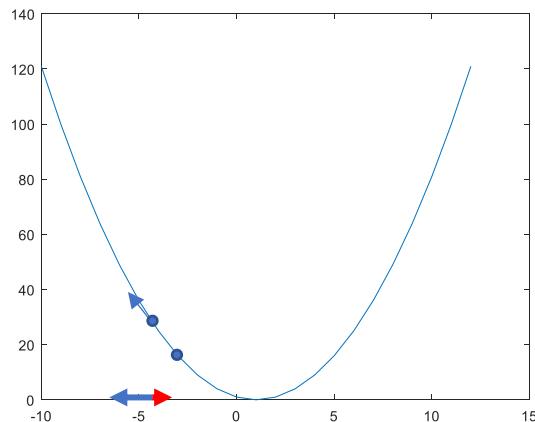
Linear Regression

- Minimize $\sum_i (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2$
- Minimize $(x-1)^2$



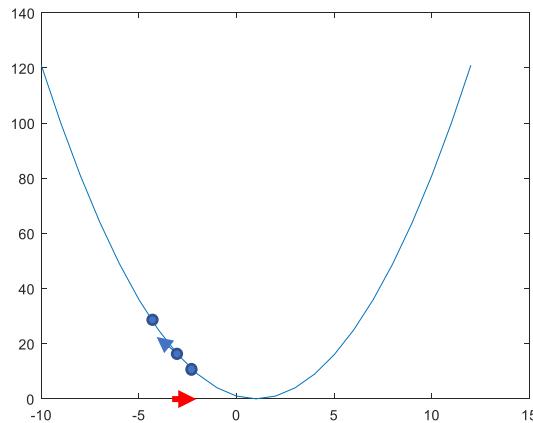
Gradient descent for linear regression

- Gradient descent
 - Follow the inverse direction of gradient for minimizing the objective function
- $x_t = x_{t-1} - \alpha J(x_t)$
- Minimize $(x-1)^2$
 - $\alpha = 0.1$
 - $x_0 = -4$
 - $x_1 = x_0 - \alpha f(x_0) = -4 - \alpha * (-10) = -3$
 - $f(-4) = 2 * (-4-1) = -10$
- $J(x) = (x-1)^2$
- $J(x)' = 2(x-1)$



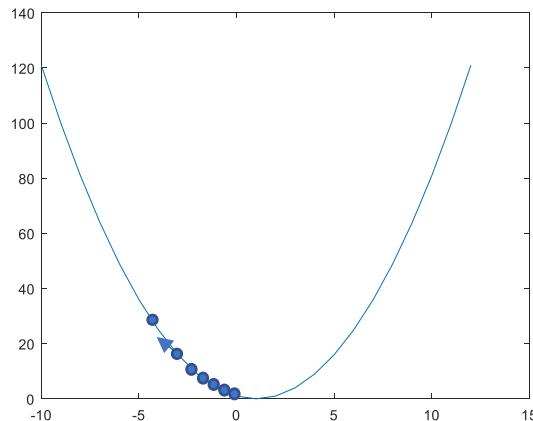
Gradient descent for linear regression

- Gradient descent
 - Follow the inverse direction of gradient for minimizing the objective function
- $x_t = x_{t-1} - \alpha J(x_{t-1})'$ $\alpha = 0.1$
 $x_0 = -4$
- Minimize $(x-1)^2$ $x_1 = x_0 - \alpha f(x_0)' = -4 - \alpha * (-10) = -3$ $J(-4)' = 2 * (-4-1) = -10$
 $x_2 = x_1 - \alpha f(x_1)' = -3 - \alpha * (-10) = -2.2$ $J(-3)' = 2 * (-3-1) = -8$
- $J(x) = (x-1)^2$
- $J(x)' = 2(x-1)$



Gradient descent for linear regression

- Gradient descent
 - Follow the inverse direction of gradient for minimizing the objective function
- $x_t = x_{t-1} - \alpha J(x_{t-1})'$ $\alpha = 0.1$
 $x_0 = -4$
- Minimize $(x-1)^2$ $x_1 = x_0 - \alpha f(x_0)' = -4 - \alpha * (-10) = -3$ $J(-4)' = 2 * (-4-1) = -10$
 $x_2 = x_1 - \alpha f(x_1)' = -3 - \alpha * (-10) = -2.2$ $J(-3)' = 2 * (-3-1) = -8$
- $J(x) = (x-1)^2$...
- $J(x)' = 2(x-1)$...



Loss function

$0.5\lambda\|\mathbf{w}\|^2$: Regularization term, the magnitude of \mathbf{w} should not be too large.

- Prediction $p_i : p_i = \mathbf{w}^T \mathbf{x}_i - b$ y_i can be any continuous value
- $J(\mathbf{w}) = \sum_i (y_i - p_i)^2 + 0.5\lambda\|\mathbf{w}\|^2$
 - Least square loss for linear regression, pixel value?
- $J(\mathbf{w}) = \sum_i \max(0, 1 - y_i * p_i) + 0.5\lambda\|\mathbf{w}\|^2$ $y_i = -1 \text{ or } 1$
 - Soft-margin loss for linear support vector machine, dog or not?
- $J(\mathbf{w}) = -\sum_i \sum_c y_{i,c} \log p_{i,c} + 0.5\lambda\|\mathbf{w}\|^2$ $y_{i,c} = 1 \text{ or } 0 \quad \sum_c y_{i,c} = 1$
 - Cross entropy loss for multi-class classification, dog, cat, or flower?



$$[y_{i,1}, y_{i,2}, y_{i,3}]$$

$$[1, 0, 0]$$



$$[0, 1, 0]$$



$$[0, 0, 1]$$

$0.5\lambda\|\mathbf{w}\|^2$: Regularization term, the magnitude of \mathbf{w} should not be too large.

Gradients for loss functions

$$p_i = \mathbf{w}^T \mathbf{x}_i - b$$

- Least square: $J(\mathbf{w}) = \sum_i (p_i - y_i)^2 + 0.5\lambda\|\mathbf{w}\|^2$ y_i can be any continuous value

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_i 2(p_i - y_i) \frac{\partial p_i}{\partial \mathbf{w}} + \lambda \mathbf{w}$$

- Soft-margin loss: $J(\mathbf{w}) = \sum_i \max(0, 1 - y_i * p_i) + 0.5\lambda\|\mathbf{w}\|^2$ $y_i = -1 \text{ or } 1$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_i (-y_i) u(1 - y_i * p_i) \frac{\partial p_i}{\partial \mathbf{w}} + \lambda \mathbf{w}$$

$u(x) = 1, \text{ if } x > 0,$
 $u(x) = 0, \text{ if } x \leq 0.$

- Cross entropy loss: $J(\mathbf{w}) = -\sum_i \sum_c y_{i,c} \log p_{i,c} + 0.5\lambda\|\mathbf{w}\|^2$ i th sample, c th class

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = -\sum_i \sum_c y_{i,c} \frac{1}{p_{i,c}} \frac{\partial p_{i,c}}{\partial \mathbf{w}} + \lambda \mathbf{w}$$

$y_{i,c} = 1 \text{ or } 0$ $\sum_c y_{i,c} = 1$



$$y_{i,1} = 1$$

$$[y_{i,1}, y_{i,2}, y_{i,3}]$$



$$y_{i,2} = 1$$

$$[0, 1, 0]$$



$$y_{i,3} = 1$$

$$[0, 0, 1]$$

$$J(\mathbf{w}) = \sum_i \max(0, 1 - y_i^* p_i) + 0.5\lambda \|\mathbf{w}\|^2$$

$$\begin{aligned} \max(0, 1 - y_i^* p_i) &= 1 - y_i^* p_i \quad \text{if } 1 - y_i^* p_i > 0 \\ &\text{or } = 0 \quad \quad \text{if } 1 - y_i^* p_i \leq 0 \end{aligned}$$

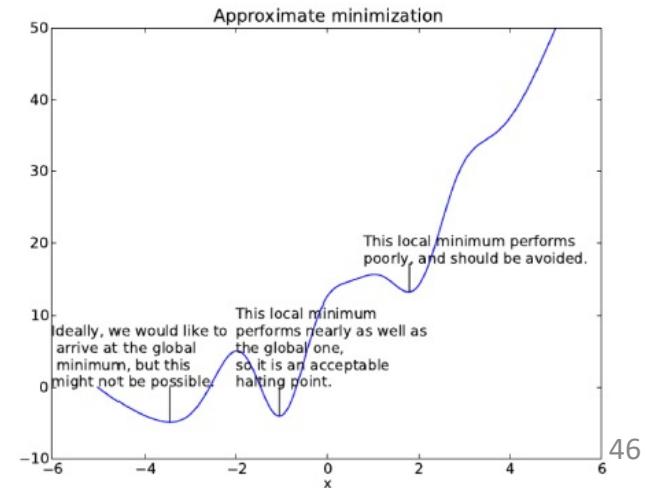
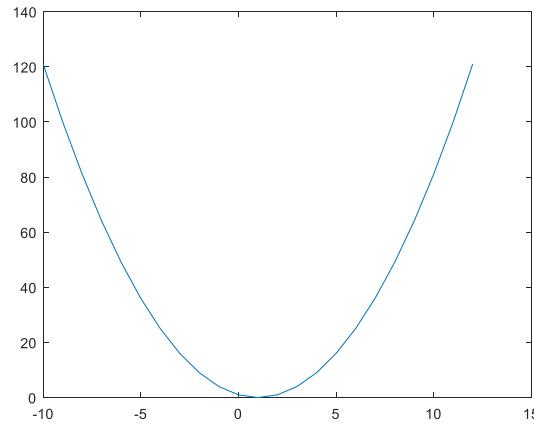
$$\begin{aligned} \frac{\partial J(p)}{\partial p} &= -y_i && \text{if } 1 - y_i^* p_i > 0 \\ &\text{or } = 0 && \text{if } 1 - y_i^* p_i \leq 0 \end{aligned}$$

$$\begin{aligned} u(x) &= 1, \text{ if } x > 0, \\ u(x) &= 0, \text{ if } x \leq 0. \end{aligned}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_i (-y_i) u(1 - y_i^* p_i) \frac{\partial p_i}{\partial \mathbf{w}} + \lambda \mathbf{w}$$

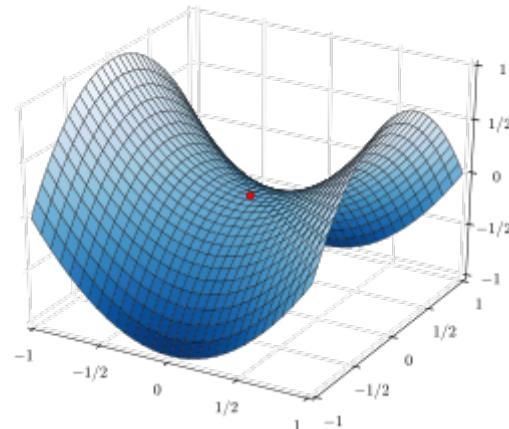
Local minimum and saddle points

- Minima
 - Smallest value of the function, either within a given range (**the local** minima) or on the entire domain of a function (**the global** minima)
- Local minima = global minima => Convex
- Local minima \neq global minima => Non-Convex



Local minimum and saddle points

- Minima
- Saddle point:
 - a point on the surface of the graph of a function where the slopes (derivatives) in orthogonal directions are both zero, but which is not a **local** minima/maxima of the function.



Outline

- Basic operations of deep model
- Basics of optimization
- Optimization for deep models

Back-propagation

- A chaining rule for obtaining gradient
- $p=f_1(h)$
- $J=f_2(p)$

$$\frac{\partial J}{\partial h}?$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h} = f_2(p)' f_1(h)'$$

$$h=f_1(h), p = f_2(p)$$
$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient
- $p=f_1(h)=2h$
- $J=f_2(p)=(p-1)^2$

$$\frac{\partial J}{\partial h}?$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h} = f_2(p)'f_1(h)'=2(p-1)\cdot 2=2(2h-1)\cdot 2=4\cdot(2h-1)$$

$$J = f_2(f_1(h)) = (2h - 1)^2$$

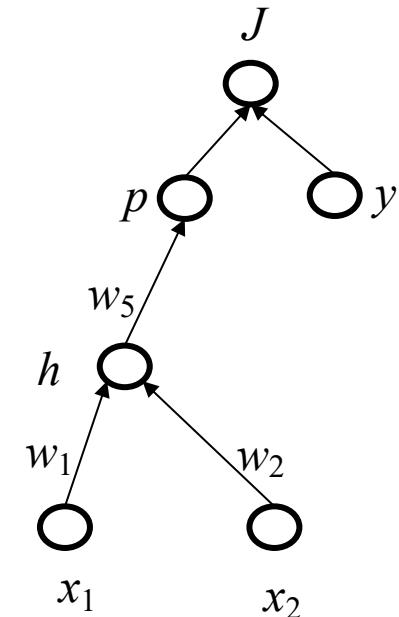
$$\frac{\partial J}{\partial h} = 4 \cdot (2h - 1)$$

$$h=f_1(h), p = f_2(p)$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5 h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$



$$h=f_1(h), p = f_2(p)$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient

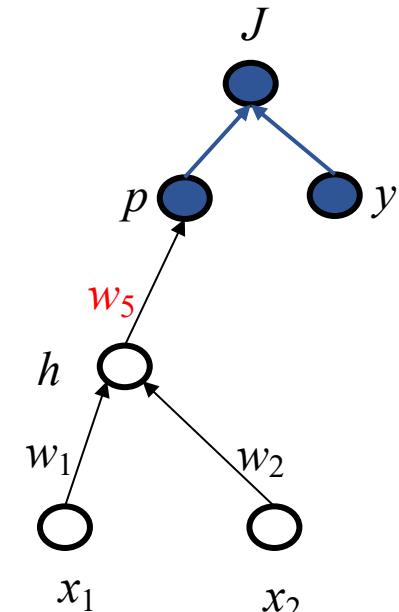
$$h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)$$

$$p_i = f_2(h_i; \mathbf{w}_5) = \mathbf{w}_5 h_i$$

- Loss:

$$J(\mathbf{p}) = \sum_i (p_i - y_i)^2$$

$$\frac{\partial J}{\partial w_5} = \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial w_5}$$

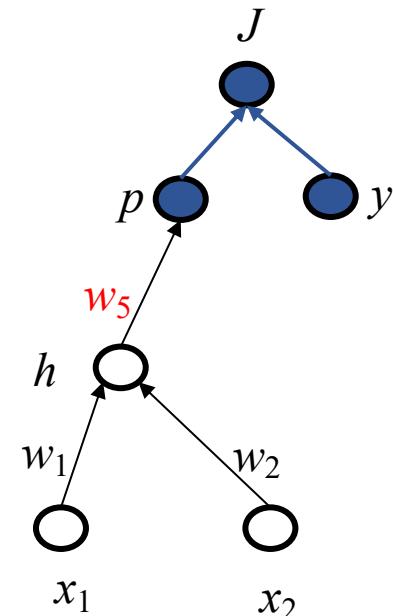


$$h=f_1(h), p = f_2(p)$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5 h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$



$$\begin{aligned}\frac{\partial J}{\partial w_5} &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial w_5} \\ &= \sum_i 2(p_i - y_i) \frac{\partial p_i}{\partial w_5}\end{aligned}$$

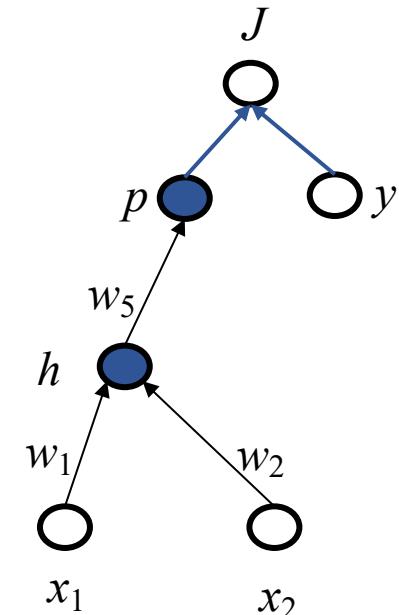
$$h=f_1(h), p = f_2(p)$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1x_{1,i} + w_2x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$

$$\begin{aligned}\frac{\partial J}{\partial w_5} &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial w_5} \\ &= \sum_i 2(p_i - y_i) \frac{\partial p_i}{\partial w_5} \\ &= \sum_i 2(p_i - y_i)h_i\end{aligned}\quad \Delta p_i = \frac{\partial J}{\partial p_i} = 2(p_i - y_i)$$



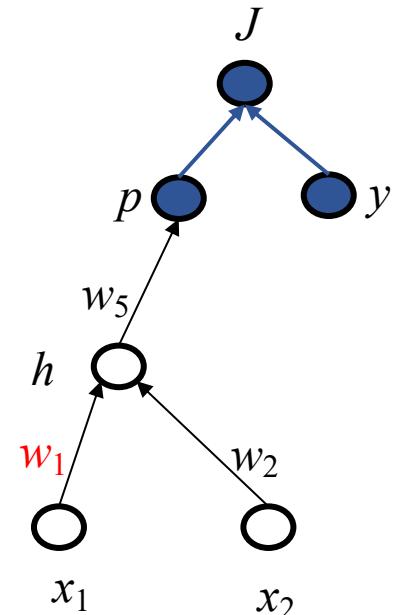
$$h=f_1(h), p = f_2(p)$$

$$\frac{\partial J}{\partial h} = \frac{\partial J}{\partial p} \frac{\partial p}{\partial h}$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(\textcolor{red}{w_1}x_{1,i} + w_2x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$

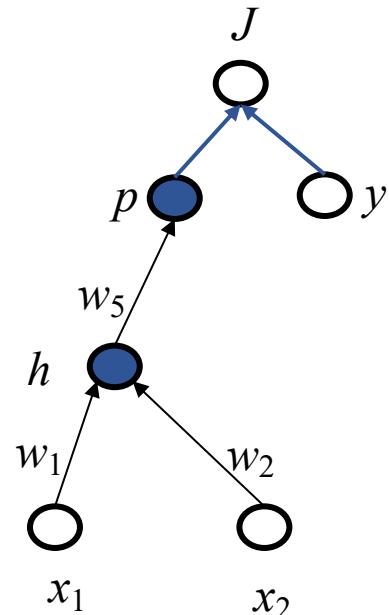
$$\begin{aligned}\frac{\partial J}{\partial w_1} &= \sum_i \boxed{\frac{\partial J}{\partial p_i}} \frac{\partial p_i}{\partial w_1} \\ &= \sum_i \boxed{2(p_i - y_i)} \frac{\partial p_i}{\partial w_1} \\ &= \sum_i \boxed{\Delta p_i} \frac{\partial p_i}{\partial w_1}\end{aligned}\quad \Delta p_i = \frac{\partial J}{\partial p_i} = 2(p_i - y_i)$$



Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1x_{1,i} + w_2x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5 h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$

$$\begin{aligned}
 \frac{\partial J}{\partial w_1} &= \sum_i \Delta p_i \frac{\partial p_i}{\partial w_1} & \Delta p_i &= \frac{\partial J}{\partial p_i} = 2(p_i - y_i) \\
 &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial h_i} \frac{\partial h_i}{\partial w_1} & \Delta h_i &= \frac{\partial J}{\partial h_i} = \Delta p_i \frac{\partial p_i}{\partial h_i} = 2(p_i - y_i)w_5 \\
 &= \sum_i \Delta h_i \frac{\partial h_i}{\partial w_1}
 \end{aligned}$$

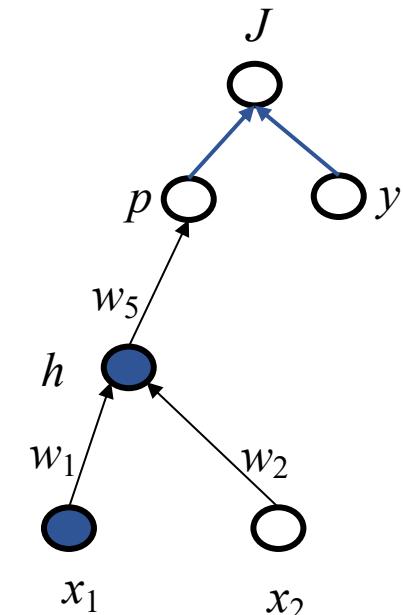


$$u(x) = 1, \text{if } x > 0,$$

$$u(x) = 0, \text{if } x \leq 0.$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5 h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$



$$\begin{aligned}\frac{\partial J}{\partial w_1} &= \sum_i \Delta p_i \frac{\partial p_i}{\partial w_1} \\ &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial h_i} \frac{\partial h_i}{\partial w_1} \\ &= \sum_i \Delta h_i \frac{\partial h_i}{\partial w_1} \\ &= \sum_i \Delta h_i u(w_1 x_{1,i} + w_2 x_{2,i}) x_{1,i}\end{aligned}$$

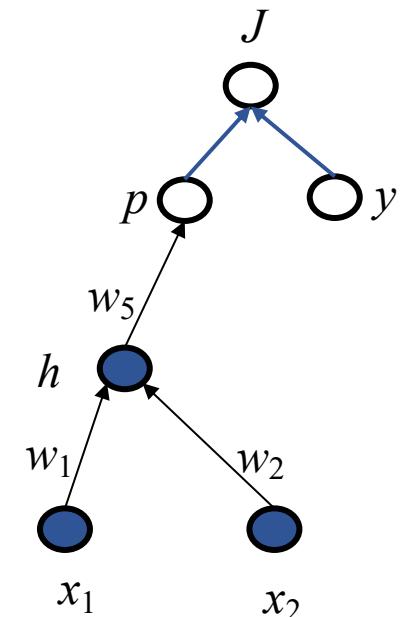
$$\begin{aligned}\Delta p_i &= \frac{\partial J}{\partial p_i} = 2(p_i - y_i) \\ \Delta h_i &= \frac{\partial J}{\partial h_i} = \Delta p_i \frac{\partial p_i}{\partial h_i} = 2(p_i - y_i) w_5 \\ \Delta x_{1,i} &= \frac{\partial J}{\partial x_{1,i}} = \Delta h_i \frac{\partial h_i}{\partial x_{1,i}} = 2(p_i - y_i) w_5 u(w_1 x_{1,i} + w_2 x_{2,i}) w_1\end{aligned}$$

$$u(x) = 1, \text{if } x > 0,$$

$$u(x) = 0, \text{if } x \leq 0.$$

Back-propagation

- A chaining rule for obtaining gradient
- $h_i = f_1(\mathbf{x}_i; w_1, w_2) = \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)$
- $p_i = f_2(h_i; w_5) = w_5 h_i$
- Loss:
- $J(\mathbf{p}, \mathbf{w}) = \sum_i (p_i - y_i)^2$



$$\begin{aligned}\frac{\partial J}{\partial w_1} &= \sum_i \Delta p_i \frac{\partial p_i}{\partial w_1} \\ &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial h_i} \frac{\partial h_i}{\partial w_1} \\ &= \sum_i \Delta h_i \frac{\partial h_i}{\partial w_1} \quad \Delta h_i = 2(p_i - y_i)w_5 \\ &= \sum_i 2(p_i - y_i)w_5 u(w_1 x_{1,i} + w_2 x_{2,i}) x_{1,i}\end{aligned}$$

$$\begin{aligned}\frac{\partial J}{\partial w_2} &= \sum_i \Delta p_i \frac{\partial p_i}{\partial w_2} \\ &= \sum_i \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial h_i} \frac{\partial h_i}{\partial w_2} \\ &= \sum_i \Delta h_i \frac{\partial h_i}{\partial w_2} \\ &= \sum_i 2(p_i - y_i)w_5 u(w_1 x_{1,i} + w_2 x_{2,i}) x_{2,i}\end{aligned}$$

$$\frac{\partial h_i}{\partial w_1} = \frac{\partial \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)}{\partial w_1}$$

$\max(w_1 x_{1,i} + w_2 x_{2,i}, 0) = w_1 x_{1,i} + w_2 x_{2,i}$ if $w_1 x_{1,i} + w_2 x_{2,i} > 0$
or = 0 otherwise

$$\frac{\partial \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)}{\partial w_1} = w_1 x_{1,I} \text{ if } w_1 x_{1,i} + w_2 x_{2,i} > 0 \\ \text{or } = 0 \text{ otherwise}$$

$$\frac{\partial \max(w_1 x_{1,i} + w_2 x_{2,i}, 0)}{\partial w_1} = u(w_1 x_{1,i} + w_2 x_{2,i}) x_{1,i}$$

$$u(x) = 1, \text{if } x > 0,$$

$$u(x) = 0, \text{if } x \leq 0.$$

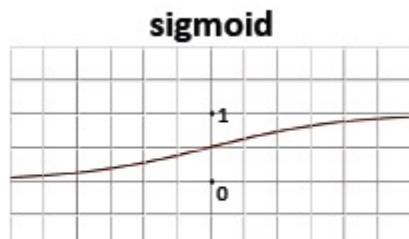
Back-propagation

- A chaining rule for obtaining gradient
- Gradients to non-linear function?

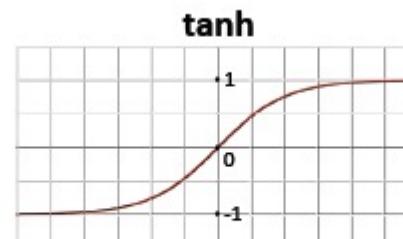
$$\text{Sigmoid}(x)' = \text{Sigmoid}(x)(1 - \text{Sigmoid}(x))$$

$$\text{Tanh}(x)' = 1 - 2 \text{Tanh}^2(x)$$

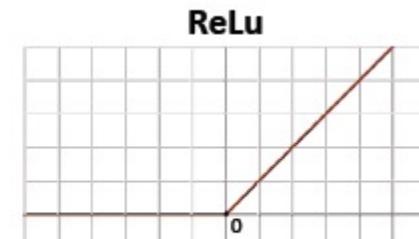
$$\text{Relu}(x)' = u(x)$$



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



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



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