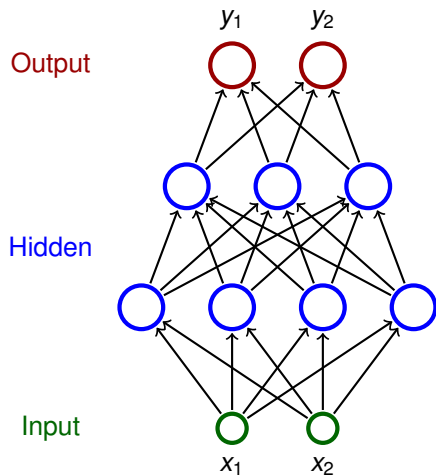


## MLP training – practical issues

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the  $i$ -th layer are connected with all neurons in the  $i + 1$ -st layer
- ▶ Architecture of a MLP is typically described by numbers of neurons in individual layers (e.g. 2-4-3-2)

## Notation:

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  *after the computation stops*
  - ▶  $y_j$  is the output of the neuron  $j$  *after the computation stops*

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the weight of the connection **from  $i$  to  $j$**   
(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )
- ▶  $j_{\leftarrow}$  is a set of all  $i$  such that  $j$  is adjacent from  $i$   
(i.e. there is an arc **to**  $j$  from  $i$ )
- ▶  $j_{\rightarrow}$  is a set of all  $i$  such that  $j$  is adjacent to  $i$   
(i.e. there is an arc **from**  $j$  to  $i$ )

## Learning:

- ▶ Given a **training set**  $\mathcal{T}$  of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \mid k = 1, \dots, p \right\}$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $j \in Y$ , denote by  $d_{kj}$  the desired output of the neuron  $j$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{kj})_{j \in Y}$ ).

- ▶ **Error function:**  $E(\vec{w}) = \sum_{k=1}^p E_k(\vec{w})$



- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:
  - ▶ Choose (randomly) a set of training examples  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

where

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)})$$

- ▶  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$
- ▶  $\nabla E_k(\vec{w}^{(t)})$  is the gradient of the error of the example  $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially.

## MLP – mse gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \frac{1}{p} \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

# MLP – mse gradient

For every  $w_{ji}$  we have

$$\frac{\partial E}{\partial w_{ji}} = \frac{1}{p} \sum_{k=1}^p \frac{\partial E_k}{\partial w_{ji}}$$

where for every  $k = 1, \dots, p$  holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

and for every  $j \in Z \setminus X$  we get (for squared error)

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

(Here all  $y_j$  are in fact  $y_j(\vec{w}, \vec{x}_k)$ ).

## (Some) error functions

- ▶ **squared error:**

$$E(\vec{w}) = \sum_{k=1}^p E_k(\vec{w})$$

where  $E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} (y_j(\vec{w}, \vec{x}_k) - d_{kj})^2$

- ▶ **mean squared error (mse):**

$$E(\vec{w}) = \frac{1}{p} \sum_{k=1}^p E_k(\vec{w})$$

- ▶ **(categorical) cross entropy:**

$$E(\vec{w}) = -\frac{1}{p} \sum_{k=1}^p \sum_{j \in Y} d_{kj} \ln(y_j)$$

# Practical issues of gradient descent

- ▶ Training efficiency:
  - ▶ What size of a minibatch?
  - ▶ How to choose the learning rate  $\varepsilon(t)$  and control SGD ?
  - ▶ How to pre-process the inputs?
  - ▶ How to initialize weights?
  - ▶ How to choose desired output values of the network?

# Practical issues of gradient descent

- ▶ Training efficiency:
  - ▶ What size of a minibatch?
  - ▶ How to choose the learning rate  $\varepsilon(t)$  and control SGD ?
  - ▶ How to pre-process the inputs?
  - ▶ How to initialize weights?
  - ▶ How to choose desired output values of the network?
- ▶ Quality of the resulting model:
  - ▶ When to stop training?
  - ▶ Regularization techniques.
  - ▶ How large network?

For simplicity, I will illustrate the reasoning on MLP + mse. Later we will see other topologies and error functions with different but always somewhat related issues.

## Issues in gradient descent

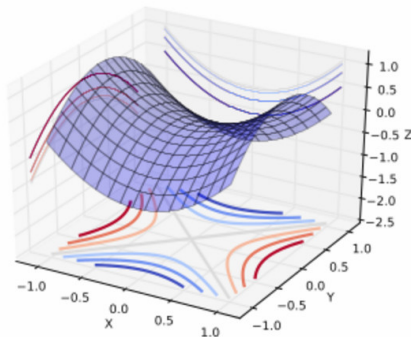
- ▶ Small networks: Lots of local minima where the descent gets stuck.
- ▶ The model identifiability problem: Swapping incoming weights of neurons  $i$  and  $j$  leaves the same network topology – **weight space symmetry**.
- ▶ Recent studies show that for sufficiently large networks all local minima have low values of the error function.

# Issues in gradient descent

- ▶ Small networks: Lots of local minima where the descent gets stuck.
- ▶ The model identifiability problem: Swapping incoming weights of neurons  $i$  and  $j$  leaves the same network topology – **weight space symmetry**.
- ▶ Recent studies show that for sufficiently large networks all local minima have low values of the error function.

## Saddle points

One can show (by a combinatorial argument) that larger networks have exponentially more saddle points than local minima.

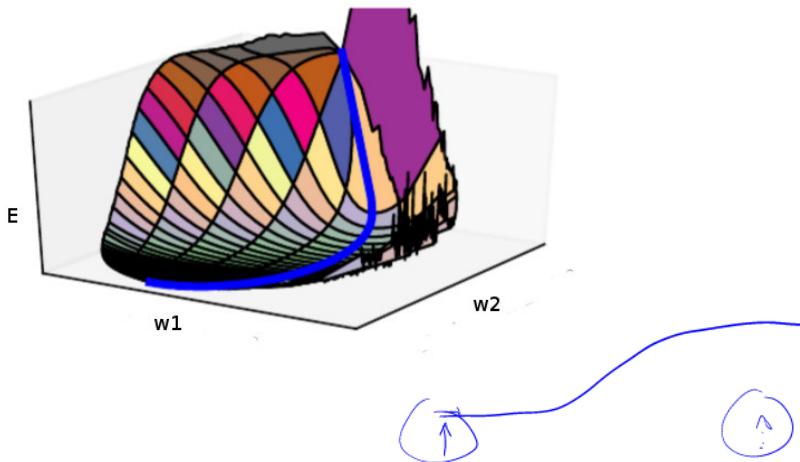




# Issues in gradient descent – too slow descent

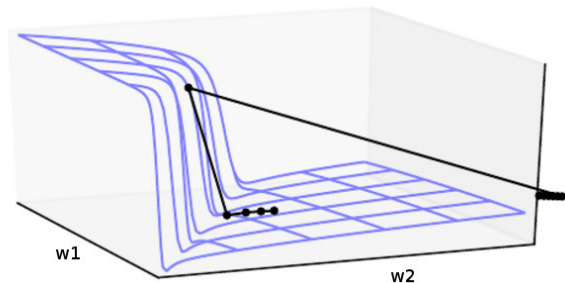
- ▶ flat regions

E.g. if the inner potentials are too large (in abs. value), then their derivative is extremely small.

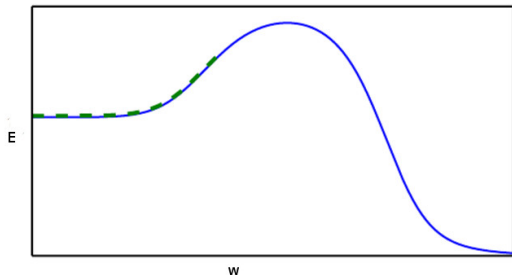


# Issues in gradient descent – too fast descent

- ▶ steep cliffs: the gradient is extremely large, descent skips important weight vectors



## Issues in gradient descent – local vs global structure



What if we initialize on the left?

# Gradient Descent in Large Networks

## Theorem

Assume (roughly),

- ▶ activation functions: "smooth" ReLU (softplus)

$$\sigma(z) = \log(1 + \exp(z))$$



*In general: Smooth, non-polynomial, analytic, Lipschitz.*

- ▶ inputs  $\vec{x}_k$  of Euclidean norm equal to 1, desired values  $d_k$  satisfying  $|d_k| \in O(1)$ ,
- ▶ the number of hidden neurons per layer sufficiently large (polynomial in certain numerical characteristics of inputs roughly measuring their similarity, and exponential in the depth of the network),
- ▶ the learning rate constant and sufficiently small.

*The gradient descent converges (with high probability) to a global minimum with zero error at linear rate.*

Later we get to a special type of networks called ResNet where the above result demands only polynomially many neurons per layer (w.r.t. depth).

# Issues in computing the gradient

- ▶ vanishing and exploding gradients

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

# Issues in computing the gradient

- ▶ vanishing and exploding gradients

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \quad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \quad \text{for } j \in Z \setminus (Y \cup X)$$

- ▶ inexact gradient computation:

- ▶ Minibatch gradient is only an estimate of the true gradient.
- ▶ Note that the variance of the estimate is (roughly)  $\sigma / \sqrt{m}$  where  $m$  is the size of the minibatch and  $\sigma$  is the variance of the gradient estimate for a single training example.  
(E.g. minibatch size 10 000 means 100 times more computation than the size 100 but gives only 10 times less variance.)

## Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.

## Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.



# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.
- ▶ It is common (especially when using GPUs) for power of 2 batch sizes to offer better runtime. Typical power of 2 batch sizes range from 32 to 256, with 16 sometimes being attempted for large models.

# Minibatch size

- ▶ Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- ▶ Multicore architectures are usually underutilized by extremely small batches.
- ▶ If all examples in the batch are to be processed in parallel (as is the typical case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.
- ▶ It is common (especially when using GPUs) for power of 2 batch sizes to offer better runtime. Typical power of 2 batch sizes range from 32 to 256, with 16 sometimes being attempted for large models.
- ▶ Small batches can offer a regularizing effect, perhaps due to the noise they add to the learning process.

It has been observed in practice that when using a larger batch there is a degradation in the quality of the model, as measured by its ability to generalize.

# Momentum

Issue in the gradient descent:

- ▶  $\nabla E(\vec{w}^{(t)})$  constantly changes direction (but the error steadily decreases).



# Momentum

Issue in the gradient descent:

- ▶  $\nabla E(\vec{w}^{(t)})$  constantly changes direction (but the error steadily decreases).

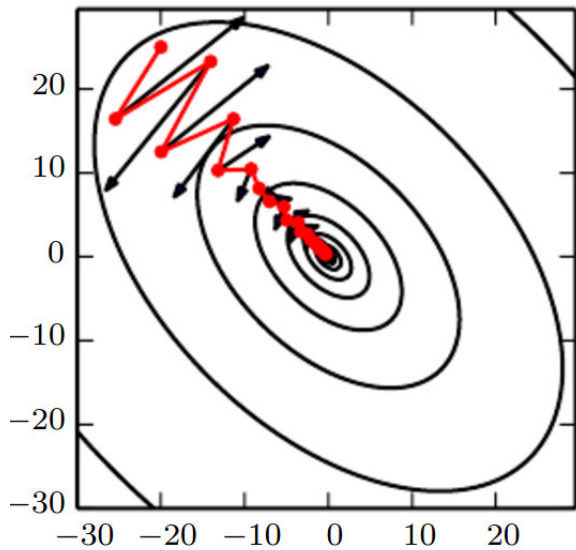


**Solution:** In every step add the change made in the previous step (weighted by a factor  $\alpha$ ):

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)}) + \alpha \cdot \Delta w_{ji}^{(t-1)}$$

where  $0 < \alpha < 1$ .

## Momentum – illustration



# SGD with momentum

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:
  - ▶ Choose (randomly) a set of training examples  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

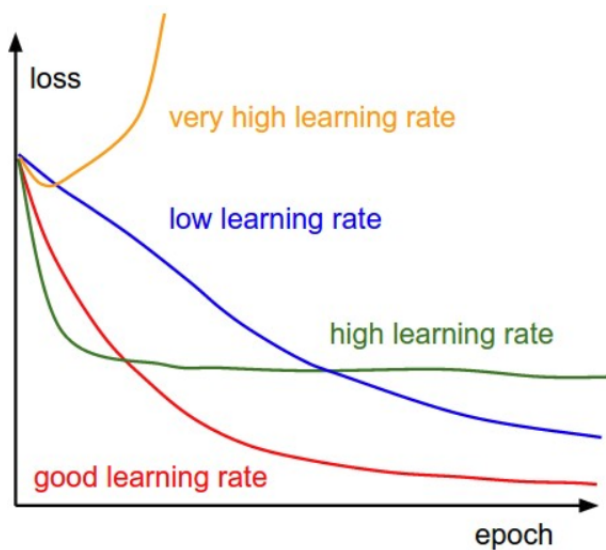
where

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)}) + \alpha \Delta \vec{w}^{(t-1)}$$

- ▶  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$
- ▶  $0 < \alpha < 1$  measures the "influence" of the momentum
- ▶  $\nabla E_k(\vec{w}^{(t)})$  is the gradient of the error of the example  $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially.

# Learning rate





# Search for the learning rate

- ▶ Use settings from a successful solution of a similar problem as a baseline.
- ▶ Search for the learning rate using the learning monitoring:
  - ▶ Search through values from small (e.g. 0.001) to (0.1), possibly multiplying by 2.
  - ▶ Train for several epochs, observe the learning curves (see cross-validation later).

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  a number of steps  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  a number of steps  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  a number of steps  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)
- ▶ Piecewise constant scheduling: A constant learning rate for a number of steps/epochs, then a smaller learning rate, and so on.

# Adaptive learning rate

- ▶ Power scheduling: Set  $\epsilon(t) = \epsilon_0 / (1 + t/s)$  where  $\epsilon_0$  is an initial learning rate and  $s$  a number of steps  
(after  $s$  steps the learning rate is  $\epsilon_0/2$ , after  $2s$  it is  $\epsilon_0/3$  etc.)
- ▶ Exponential scheduling: Set  $\epsilon(t) = \epsilon_0 \cdot 0.1^{t/s}$ .  
(the learning rate decays faster than in the power scheduling)
- ▶ Piecewise constant scheduling: A constant learning rate for a number of steps/epochs, then a smaller learning rate, and so on.
- ▶ 1cycle scheduling: Start by increasing the initial learning rate from  $\epsilon_0$  linearly to  $\epsilon_1$  (approx.  $\epsilon_1 = 10\epsilon_0$ ) halfway through training. Then decrease from  $\epsilon_1$  linearly to  $\epsilon_0$ . Finish by dropping the learning rate by several orders of magnitude (still linearly).  
According to a 2018 paper by Leslie Smith this may converge much faster (100 epochs vs 800 epochs on CIFAR10 dataset).

For comparison of some methods see: AN EMPIRICAL STUDY OF LEARNING RATES IN DEEP NEURAL NETWORKS FOR SPEECH RECOGNITION, Senior et al

So far we have considered fixed schedules for learning rates.

It is better to have

- ▶ larger rates for weights with smaller updates,
- ▶ smaller rates for weights with larger updates.

AdaGrad uses individually adapting learning rate for each weight.

# SGD with AdaGrad

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# SGD with AdaGrad

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2, \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\frac{\eta}{\sqrt{r_{ji}^{(t)} + \delta}} \cdot \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)})$$

and

$$r_{ji}^{(t)} = r_{ji}^{(t-1)} + \left( \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)}) \right)^2$$

- ▶  $\eta$  is a constant expressing the influence of the learning rate, typically 0.01.
- ▶  $\delta > 0$  is a smoothing term (typically 1e-8) avoiding division by 0.



The main disadvantage of AdaGrad is the accumulation of the gradient throughout the whole learning process.

In case the learning needs to get over several "hills" before settling in a deep "valley", the weight updates get far too small before getting to it.

RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl.

# SGD with RMSProp

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# SGD with RMSProp

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), compute  $\vec{w}^{(t+1)}$  :
  - ▶ Choose (randomly) a minibatch  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\frac{\eta}{\sqrt{r_{ji}^{(t)} + \delta}} \cdot \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)})$$

and

$$r_{ji}^{(t)} = \rho r_{ji}^{(t-1)} + (1 - \rho) \left( \sum_{k \in T} \frac{\partial E_k}{\partial w_{ji}}(\vec{w}^{(t)}) \right)^2$$

- ▶  $\eta$  is a constant expressing the influence of the learning rate (Hinton suggests  $\rho = 0.9$  and  $\eta = 0.001$ ).
- ▶  $\delta > 0$  is a smoothing term (typically  $1e-8$ ) avoiding division by 0.

## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

Unfortunately, there is currently no consensus on this point.

According to a recent study, the family of algorithms with adaptive learning rates (represented by RMSProp and AdaDelta) performed fairly robustly, no single best algorithm has emerged.

## Other optimization methods

There are more methods such as AdaDelta, Adam (roughly RMSProp combined with momentum), etc.

A natural question: Which algorithm should one choose?

Unfortunately, there is currently no consensus on this point.

According to a recent study, the family of algorithms with adaptive learning rates (represented by RMSProp and AdaDelta) performed fairly robustly, no single best algorithm has emerged.

Currently, the most popular optimization algorithms actively in use include SGD, SGD with momentum, RMSProp, RMSProp with momentum, AdaDelta and Adam.

The choice of which algorithm to use, at this point, seems to depend largely on the user's familiarity with the algorithm.

# Choice of (hidden) activations

Generic requirements imposed on activation functions:

1. differentiability

(to do gradient descent)

2. non-linearity

(linear multi-layer networks are equivalent to single-layer)

3. monotonicity

(local extrema of activation functions induce local extrema of the error function)

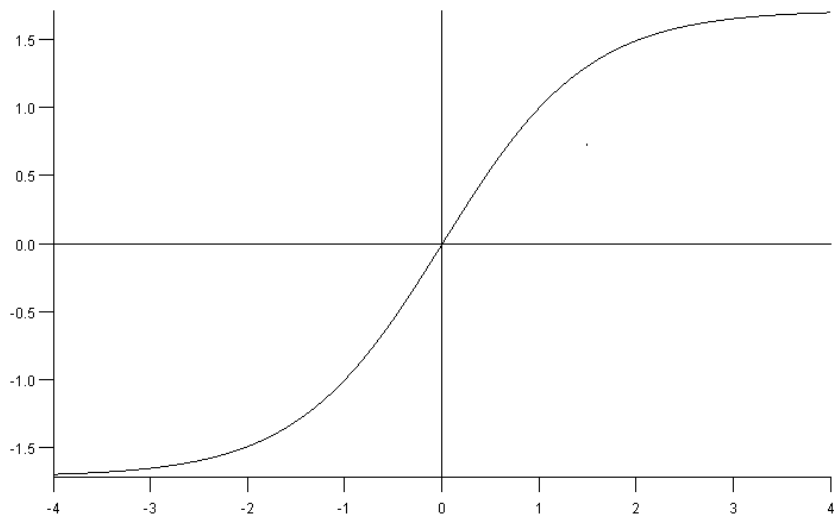
4. "linearity"

(i.e. preserve as much linearity as possible; linear models are easiest to fit; find the "minimum" non-linearity needed to solve a given task)



The choice of activation functions is closely related to input preprocessing and the initial choice of weights. I will illustrate the reasoning on sigmoidal functions; say few words about other activation functions later.

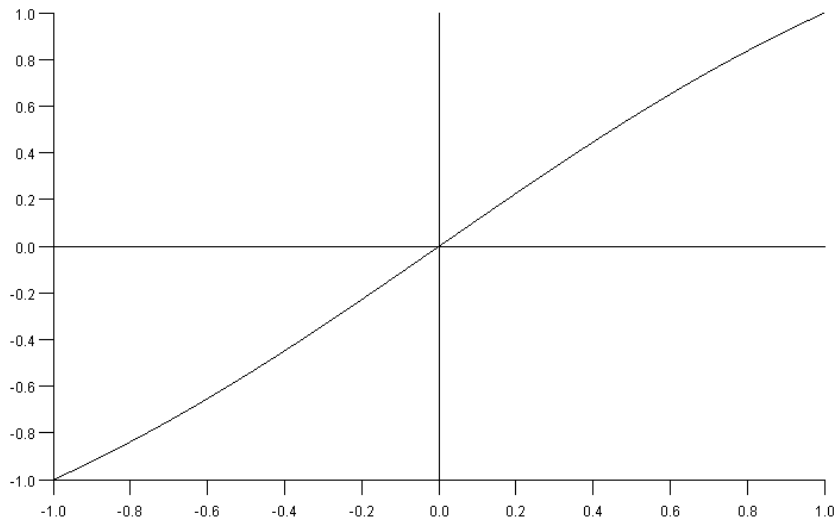
## Activation functions – tanh



$\sigma(\xi) = 1.7159 \cdot \tanh\left(\frac{2}{3} \cdot \xi\right)$ , we have  $\lim_{\xi \rightarrow \infty} \sigma(\xi) = 1.7159$  and  $\lim_{\xi \rightarrow -\infty} \sigma(\xi) = -1.7159$

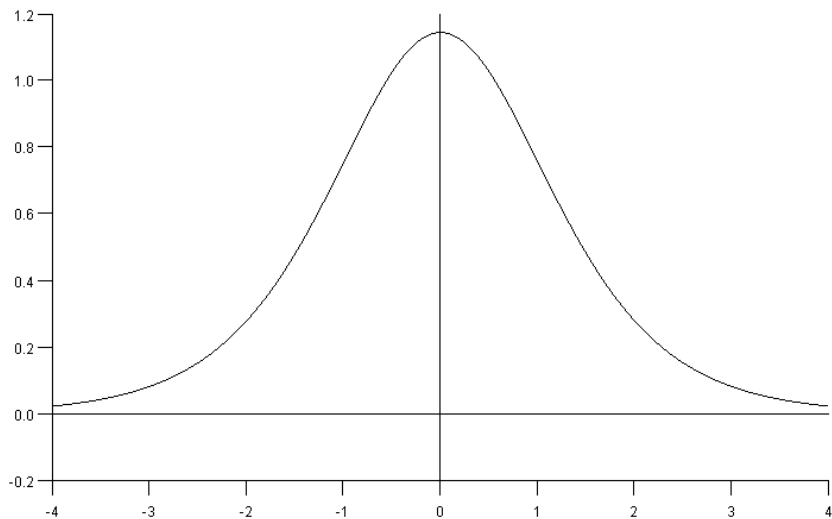


## Activation functions – tanh



$\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  is almost linear on  $[-1, 1]$

## Activation functions – tanh



first derivative:  $\sigma(\xi) = 1.7159 \cdot \tanh\left(\frac{2}{3} \cdot \xi\right)$

# Input preprocessing

- ▶ Some inputs may be much larger than others.

E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.

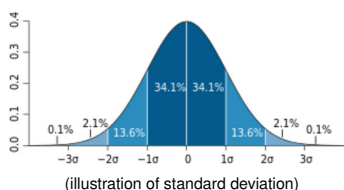
# Input preprocessing

- ▶ Some inputs may be much larger than others.  
E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.
- ▶ Large inputs have greater influence on the training than the small ones. In addition, too large inputs may slow down learning (saturation of activation functions).

# Input preprocessing

- ▶ Some inputs may be much larger than others.  
E.g.: Height vs weight of a person, maximum speed of a car (in km/h) vs its price (in CZK), etc.
- ▶ Large inputs have greater influence on the training than the small ones. In addition, too large inputs may slow down learning (saturation of activation functions).
- ▶ Typical standardization:
  - ▶ average = 0 (subtract the mean)
  - ▶ variance = 1 (divide by the standard deviation)

Here the mean and standard deviation may be estimated from data (the training set).



## Initial weights (for tanh)

- ▶ Assume weights chosen uniformly in random from an interval  $[-w, w]$  where  $w$  depends on the number of inputs of a given neuron.

## Initial weights (for tanh)

- ▶ Assume weights chosen uniformly in random from an interval  $[-w, w]$  where  $w$  depends on the number of inputs of a given neuron.
- ▶ Consider the activation function  $\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  for all neurons.
  - ▶  $\sigma$  is almost linear on  $[-1, 1]$ ,
  - ▶  $\sigma$  saturates out of the interval  $[-4, 4]$  (i.e. it is close to its limit values and its derivative is close to 0).

## Initial weights (for tanh)

- ▶ Assume weights chosen uniformly in random from an interval  $[-w, w]$  where  $w$  depends on the number of inputs of a given neuron.
- ▶ Consider the activation function  $\sigma(\xi) = 1.7159 \cdot \tanh(\frac{2}{3} \cdot \xi)$  for all neurons.
  - ▶  $\sigma$  is almost linear on  $[-1, 1]$ ,
  - ▶  $\sigma$  saturates out of the interval  $[-4, 4]$  (i.e. it is close to its limit values and its derivative is close to 0).

Thus

- ▶ for too small  $w$  we may get (almost) linear model.
- ▶ for too large  $w$  (i.e. much larger than 1) the activations may get saturated and the learning will be very slow.

Hence, we want to choose  $w$  so that the inner potentials of neurons will be roughly in the interval  $[-1, 1]$ .



## Initial weights (for tanh)

- ▶ Standardization gives mean = 0 and variance = 1 of the input data.

## Initial weights (for tanh)

- ▶ Standardization gives mean = 0 and variance = 1 of the input data.
- ▶ Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume that its weights are chosen uniformly from  $[-w, w]$ .

## Initial weights (for tanh)

- ▶ Standardization gives mean = 0 and variance = 1 of the input data.
- ▶ Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume that its weights are chosen uniformly from  $[-w, w]$ .
- ▶ **The rule:** choose  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $\sigma_j$ ) is close to the border of the interval on which  $\sigma_j$  is linear. In our case:  $\sigma_j \approx 1$ .

## Initial weights (for tanh)

- ▶ Standardization gives mean = 0 and variance = 1 of the input data.
- ▶ Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume that its weights are chosen uniformly from  $[-w, w]$ .
- ▶ **The rule:** choose  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $\sigma_j$ ) is close to the border of the interval on which  $\sigma_j$  is linear. In our case:  $\sigma_j \approx 1$ .
- ▶ Our assumptions imply:  $\sigma_j = \sqrt{\frac{n}{3}} \cdot w$ .

Thus we put  $w = \frac{\sqrt{3}}{\sqrt{n}}$ .

## Initial weights (for tanh)

- ▶ Standardization gives mean = 0 and variance = 1 of the input data.
- ▶ Consider a neuron  $j$  from the first layer with  $n$  inputs. Assume that its weights are chosen uniformly from  $[-w, w]$ .
- ▶ **The rule:** choose  $w$  so that the *standard deviation* of  $\xi_j$  (denote by  $\sigma_j$ ) is close to the border of the interval on which  $\sigma_j$  is linear. In our case:  $\sigma_j \approx 1$ .
- ▶ Our assumptions imply:  $\sigma_j = \sqrt{\frac{n}{3}} \cdot w$ .  
Thus we put  $w = \frac{\sqrt{3}}{\sqrt{n}}$ .
- ▶ The same works for higher layers,  $n$  corresponds to the number of neurons in the layer one level lower.

## Glorot & Bengio initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

# Glorot & Bengio initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

Glorot & Bengio (2010) presented a **normalized initialization** by choosing  $w$  uniformly from the interval:

$$\left( -\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}} \right) = \left( -\sqrt{\frac{3}{(m+n)/2}}, \sqrt{\frac{3}{(m+n)/2}} \right)$$

Here  $n$  is the number of inputs to the layer,  $m$  is the number of outputs of the layer (i.e. the number of neurons in the layer).

# Glorot & Bengio initialization

The previous heuristics for weight initialization ignores variance of the gradient (i.e. it is concerned only with the "size" of activations in the forward pass).

Glorot & Bengio (2010) presented a **normalized initialization** by choosing  $w$  uniformly from the interval:

$$\left( -\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}} \right) = \left( -\sqrt{\frac{3}{(m+n)/2}}, \sqrt{\frac{3}{(m+n)/2}} \right)$$

Here  $n$  is the number of inputs to the layer,  $m$  is the number of outputs of the layer (i.e. the number of neurons in the layer).

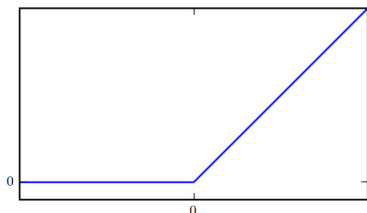
This is designed to compromise between the goal of initializing all layers to have the same activation variance and the goal of initializing all layers to have the same gradient variance.

The formula is derived using the assumption that the network consists only of a chain of matrix multiplications, with no non-linearities. Real neural networks obviously violate this assumption, but many strategies designed for the linear model perform reasonably well on its non-linear counterparts.



# Modern activation functions

For hidden neurons sigmoidal functions are often substituted with piece-wise linear activations functions. Most prominent is ReLU:

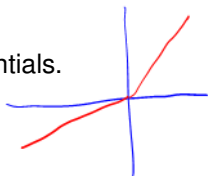


$$\sigma(\xi) = \max\{0, \xi\}$$

- ▶ THE default activation function recommended for use with most feedforward neural networks.
- ▶ As close to linear function as possible; very simple; does not saturate for large potentials.
- ▶ Dead for negative potentials.

## More modern activation functions

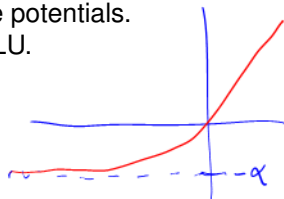
- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.



# More modern activation functions

- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.
- ▶ ELU: "Smoothed" ReLU:

$$\sigma(\xi) = \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$



Here  $\alpha$  is a parameter, ELU converges to  $-\alpha$  as  $\xi \rightarrow -\infty$ . As opposed to ReLU: Smooth, always non-zero gradient (but saturates), slower to compute.

## More modern activation functions

- ▶ Leaky ReLU (greenboard):
  - ▶ Generalizes ReLU, not dead for negative potentials.
  - ▶ Experimentally not much better than ReLU.
- ▶ ELU: "Smoothed" ReLU:

$$\sigma(\xi) = \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

Here  $\alpha$  is a parameter, ELU converges to  $-\alpha$  as  $\xi \rightarrow -\infty$ . As opposed to ReLU: Smooth, always non-zero gradient (but saturates), slower to compute.

- ▶ SELU: Scaled variant of ELU: :

$$\sigma(\xi) = \lambda \begin{cases} \alpha(\exp(\xi) - 1) & \text{for } \xi < 0 \\ \xi & \text{for } \xi \geq 0 \end{cases}$$

*Self-normalizing*, i.e. output of each layer will tend to preserve a mean (close to) 0 and a standard deviation (close to) 1 for  $\lambda \approx 1.050$  and  $\alpha \approx 1.673$ , properly initialized weights (see below) and normalized inputs (zero mean, standard deviation 1).

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ Glorot & Bengio (2010): Choose weights randomly from the normal distribution with mean 0 and variance  $2/(n + m)$

Suitable for activation functions: None, tanh, logistic, softmax

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ Glorot & Bengio (2010): Choose weights randomly from the normal distribution with mean 0 and variance  $2/(n + m)$

Suitable for activation functions: None, tanh, logistic, softmax

- ▶ He (2015): Choose weights randomly from the normal distribution with mean 0 and variance  $2/n$   
Designed for ReLU, leaky ReLU

# Initializing with Normal Distribution

Denote by  $n$  the number of inputs to the initialized layer, and  $m$  the number of neurons in the layer.

- ▶ Glorot & Bengio (2010): Choose weights randomly from the normal distribution with mean 0 and variance  $2/(n + m)$

Suitable for activation functions: None, tanh, logistic, softmax

- ▶ He (2015): Choose weights randomly from the normal distribution with mean 0 and variance  $2/n$   
Designed for ReLU, leaky ReLU

- ▶ LeCun (1990): Choose weights randomly from the normal distribution with mean 0 and variance  $1/n$   
Suitable for SELU

# How to choose activation of hidden neurons

- ▶ Default is ReLU.
- ▶ According to Aurélien Géron:

*SELU > ELU > leakyReLU > ReLU > tanh > logistic*

For discussion see: Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, Aurélien Géron



## Output neurons

The choice of activation functions for output units depends on the concrete applications.

For regression (function approximation) the output is typically linear.

# Output neurons

The choice of activation functions for output units depends on the concrete applications.

For regression (function approximation) the output is typically linear.

For classification, the current activation functions of choice are

- ▶ logistic sigmoid – binary classification
- ▶ softmax: Given an output neuron  $j \in Y$

$$y_j = \sigma_j(\xi_j) = \frac{e^{\xi_j}}{\sum_{i \in Y} e^{\xi_i}}$$

for multi-class classification.

# Output neurons

The choice of activation functions for output units depends on the concrete applications.

For regression (function approximation) the output is typically linear.

For classification, the current activation functions of choice are

- ▶ logistic sigmoid – binary classification
- ▶ softmax: Given an output neuron  $j \in Y$

$$y_j = \sigma_j(\xi_j) = \frac{e^{\xi_j}}{\sum_{i \in Y} e^{\xi_i}}$$

for multi-class classification.

The error function used with softmax (assuming that the target values  $d_{kj}$  are from  $\{0, 1\}$ ) is typically **cross-entropy**:

$$-\frac{1}{p} \sum_{k=1}^p \sum_{j \in Y} d_{kj} \ln(y_j)$$

... which somewhat corresponds to the maximum likelihood principle.

## Sigmoidal outputs with cross-entropy – in detail

Consider

- ▶ Binary classification, two classes  $\{0, 1\}$
- ▶ One output neuron  $j$ , its activation logistic sigmoid

$$\sigma_j(\xi_j) = \frac{1}{1 + e^{-\xi_j}}$$

The output of the network is  $y = \sigma_j(\xi_j)$ .

# Sigmoidal outputs with cross-entropy – in detail

Consider

- ▶ Binary classification, two classes  $\{0, 1\}$
- ▶ One output neuron  $j$ , its activation logistic sigmoid

$$\sigma_j(\xi_j) = \frac{1}{1 + e^{-\xi_j}}$$

The output of the network is  $y = \sigma_j(\xi_j)$ .

- ▶ For a training set

$$\mathcal{T} = \left\{ \left( \vec{x}_k, d_k \right) \mid k = 1, \dots, p \right\}$$

(here  $\vec{x}_k \in \mathbb{R}^{|\mathcal{X}|}$  and  $d_k \in \mathbb{R}$ ), the cross-entropy looks like this:

$$E^{\text{cross}} = -\frac{1}{p} \sum_{k=1}^p [d_k \ln(y_k) + (1 - d_k) \ln(1 - y_k)]$$

where  $y_k$  is the output of the network for the  $k$ -th training input  $\vec{x}_k$ , and  $d_k$  is the  $k$ -th desired output.

# Generalization

**Intuition:** Generalization = ability to cope with new unseen instances.

Data are mostly noisy, so it is not good idea to fit exactly.

In case of function approximation, the network should not return exact results as in the training set.

# Generalization

**Intuition:** Generalization = ability to cope with new unseen instances.

Data are mostly noisy, so it is not good idea to fit exactly.

In case of function approximation, the network should not return exact results as in the training set.

More formally: It is typically assumed that the training set has been generated as follows:

$$d_{kj} = g_j(\vec{x}_k) + \Theta_{kj}$$

where  $g_j$  is the "underlying" function corresponding to the output neuron  $j \in Y$  and  $\Theta_{kj}$  is random noise.

The network should fit  $g_j$  not the noise.

Methods improving generalization are called **regularization methods**.

# Regularization

Regularization is a big issue in neural networks, as they typically use a huge amount of parameters and thus are very susceptible to overfitting.



Regularization is a big issue in neural networks, as they typically use a huge amount of parameters and thus are very susceptible to overfitting.

von Neumann: **"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk."**

... and I ask you prof. Neumann:

What can you fit with 40GB of parameters??

## Early stopping

Early stopping means that we stop learning before it reaches a minimum of the error  $E$ .

When to stop?

# Early stopping

Early stopping means that we stop learning before it reaches a minimum of the error  $E$ .

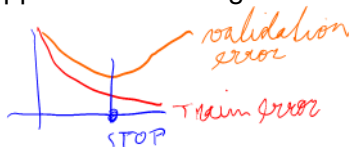
When to stop?

In many applications the error function is not the main thing we want to optimize.

E.g. in the case of a trading system, we typically want to maximize our profit not to minimize (strange) error functions designed to be easily differentiable.

Also, as noted before, minimizing  $E$  completely is not good for generalization.

For start: We may employ standard approach of training on one set and stopping on another one.



## Early stopping

Divide your dataset into several subsets:

- ▶ **training set** (e.g. 60%) – train the network here
- ▶ **validation set** (e.g. 20%) – use to stop the training
- ▶ (possibly) **test set** (e.g. 20%) – use to compare trained models

What to use as a stopping rule?

# Early stopping

Divide your dataset into several subsets:

- ▶ **training set** (e.g. 60%) – train the network here
- ▶ **validation set** (e.g. 20%) – use to stop the training
- ▶ (possibly) **test set** (e.g. 20%) – use to compare trained models

What to use as a stopping rule?

You may observe  $E$  (or any other function of interest) on the validation set, if it does not improve for last  $k$  steps, stop.

Alternatively, you may observe the gradient, if it is small for some time, stop.

(recent studies shown that this traditional rule is not too good: it may happen that the gradient is larger close to minimum values; on the other hand,  $E$  does not have to be evaluated which saves time.

To compare models you may use ML techniques such as various types of cross-validation etc.

## Size of the network

Similar problem as in the case of the training duration:

- ▶ Too small network is not able to capture intrinsic properties of the training set.
- ▶ Large networks overfit faster.

**Solution:** Optimal number of neurons :-)

# Size of the network

Similar problem as in the case of the training duration:

- ▶ Too small network is not able to capture intrinsic properties of the training set.
- ▶ Large networks overfit faster.

**Solution:** Optimal number of neurons :-)

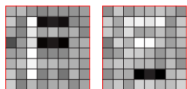
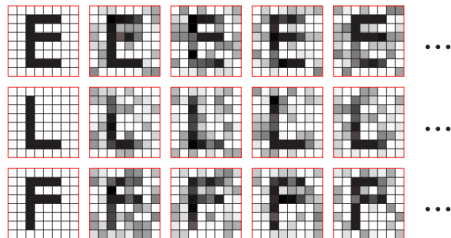
- ▶ there are some (useless) theoretical bounds
- ▶ there are algorithms dynamically adding/removing neurons (not much use nowadays)
- ▶ In practice:
  - ▶ start using a rule of thumb: the number of neurons  $\approx$  ten times less than the number of training instances.
  - ▶ experiment, experiment, experiment.

# Feature extraction

Consider a two layer network. Hidden neurons are supposed to represent "patterns" in the inputs.

Example: Network 64-2-3 for letter classification:

*sample training patterns*



*learned input-to-hidden weights*



# Ensemble methods

Techniques for reducing generalization error by combining several models.

The reason that ensemble methods work is that different models will usually not make all the same errors on the test set.

**Idea:** Train several different models separately, then have all of the models vote on the output for test examples.

# Ensemble methods

Techniques for reducing generalization error by combining several models.

The reason that ensemble methods work is that different models will usually not make all the same errors on the test set.

**Idea:** Train several different models separately, then have all of the models vote on the output for test examples.

## Bagging:

- ▶ Generate  $k$  training sets  $T_1, \dots, T_k$  by *sampling from  $\mathcal{T}$  uniformly with replacement*.

If the number of samples is  $|\mathcal{T}|$ , then on average  $|T_i| = (1 - 1/e)|\mathcal{T}|$ .

- ▶ For each  $i$ , train a model  $M_i$  on  $T_i$ .
- ▶ Combine outputs of the models: for regression by averaging, for classification by (majority) voting.

# Dropout

**The algorithm:** In every step of the gradient descent

- ▶ choose randomly a set  $N$  of neurons, each neuron is included in  $N$  independently with probability  $1/2$ ,  
(in practice, different probabilities are used as well).
- ▶ do forward and backward propagations only using the selected neurons  
(i.e. leave weights of the other neurons unchanged)

# Dropout

**The algorithm:** In every step of the gradient descent

- ▶ choose randomly a set  $N$  of neurons, each neuron is included in  $N$  independently with probability  $1/2$ ,  
(in practice, different probabilities are used as well).
- ▶ do forward and backward propagations only using the selected neurons  
(i.e. leave weights of the other neurons unchanged)

Dropout resembles bagging: Large ensemble of neural networks is trained "at once" on parts of the data.

Dropout is not exactly the same as bagging: The models share parameters, with each model inheriting a different subset of parameters from the parent neural network. This parameter sharing makes it possible to represent an exponential number of models with a tractable amount of memory.

In the case of bagging, each model is trained to convergence on its respective training set. This would be infeasible for large networks/training sets.

## Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

## Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

In every step we decrease weights (multiplicatively) as follows:

$$w_{ji}^{(t+1)} = (1 - \zeta)(w_{ji}^{(t)} + \Delta w_{ji}^{(t)})$$

Intuition: Unimportant weights will be pushed to 0, important weights will survive the decay.

## Weight decay and L2 regularization

Generalization can be improved by removing "unimportant" weights.

Penalising large weights gives stronger indication about their importance.

In every step we decrease weights (multiplicatively) as follows:

$$w_{ji}^{(t+1)} = (1 - \zeta)(w_{ji}^{(t)} + \Delta w_{ji}^{(t)})$$

Intuition: Unimportant weights will be pushed to 0, important weights will survive the decay.

Weight decay is equivalent to the gradient descent with a constant learning rate  $\varepsilon$  and the following error function:

$$E'(\vec{w}) = E(\vec{w}) + \frac{2\zeta}{\varepsilon}(\vec{w} \cdot \vec{w})$$

Here  $\frac{2\zeta}{\varepsilon}(\vec{w} \cdot \vec{w})$  is the L2 regularization that penalizes large weights.

There are many more practical tips, optimization methods, regularization methods, etc.

For a very nice survey see

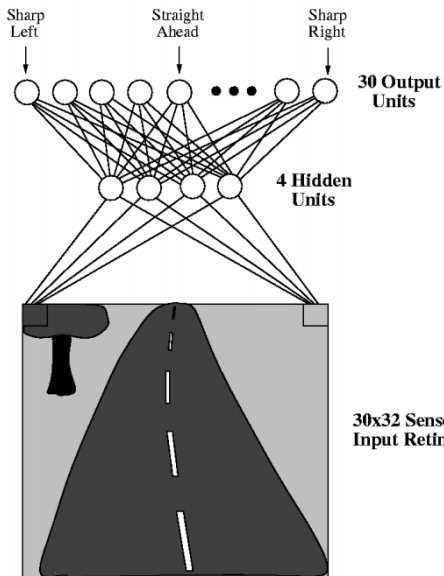
<http://www.deeplearningbook.org/>

... and also all other infinitely many urls concerned with deep learning.



## Some applications

# ALVINN (history)



## Architecture:

- ▶ MLP, 960 – 4 – 30 (also 960 – 5 – 30)
- ▶ inputs correspond to pixels

## Architecture:

- ▶ MLP, 960 – 4 – 30 (also 960 – 5 – 30)
- ▶ inputs correspond to pixels

## Activity:

- ▶ activation functions: logistic sigmoid
- ▶ Steering wheel position determined by "center of mass" of neuron values.



**Learning:** Trained during (live) drive.

- ▶ Front window view captured by a camera, 25 images per second.
- ▶ Training samples of the form  $(\vec{x}_k, \vec{d}_k)$  where
  - ▶  $\vec{x}_k$  = image of the road
  - ▶  $\vec{d}_k$  = corresponding position of the steering wheel
- ▶ position of the steering wheel "blurred" by Gaussian distribution:

$$d_{ki} = e^{-D_i^2/10}$$



where  $D_i$  is the distance of the  $i$ -th output from the one which corresponds to the correct position of the wheel.

(The authors claim that this was better than the binary output.)

## ALVINN – Selection of training samples

Naive approach: take images directly from the camera and adapt accordingly.

# ALVINN – Selection of training samples

Naive approach: take images directly from the camera and adapt accordingly.

Problems:

- ▶ If the driver is gentle enough, the car never learns how to get out of dangerous situations. A solution may be
  - ▶ turn off learning for a moment, then suddenly switch on, and let the net catch on,
  - ▶ let the driver drive as if being insane (dangerous, possibly expensive).
- ▶ The real view out of the front window is repetitive and boring, the net would overfit on few examples.

## ALVINN – Selection of training examples

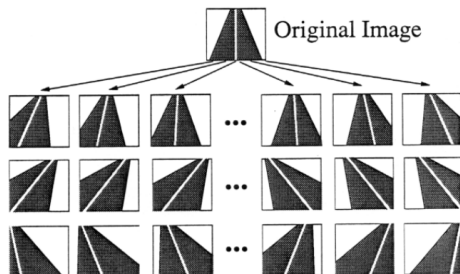
Problem with a "good" driver is solved as follows:



# ALVINN – Selection of training examples

Problem with a "good" driver is solved as follows:

- ▶ 15 distorted copies of each image:

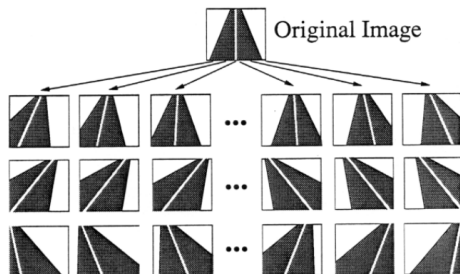


- ▶ desired output generated for each copy

# ALVINN – Selection of training examples

Problem with a "good" driver is solved as follows:

- ▶ 15 distorted copies of each image:



- ▶ desired output generated for each copy

"Boring" images solved as follows:

- ▶ a buffer of 200 images (including 15 copies of the original), in every step the system trains on the buffer
- ▶ after several updates a new image is captured, 15 copies are made and they will substitute 15 images in the buffer (5 chosen randomly, 10 with the **smallest** error).

- ▶ pure backpropagation
- ▶ constant learning rate
- ▶ momentum, slowly increasing.

## Results:

- ▶ Trained for 5 minutes, speed 4 miles per hour.
- ▶ ALVINN was able to drive well on a new road it has never seen (in different weather conditions).

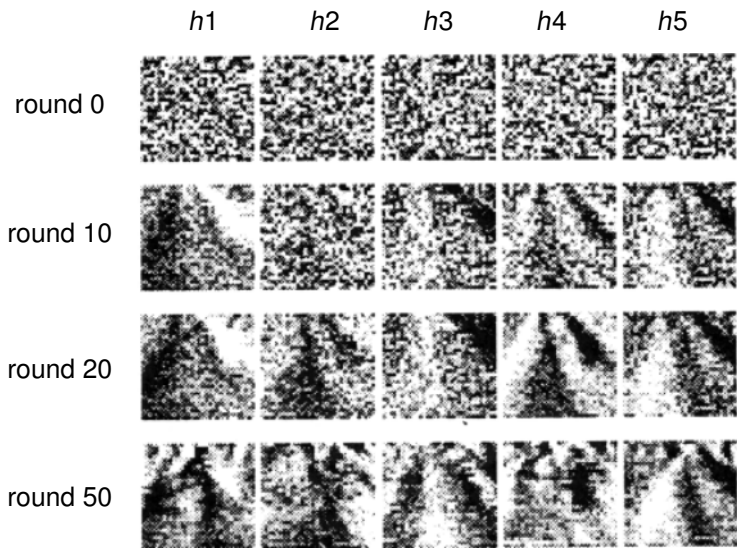
# ALVINN - learning

- ▶ pure backpropagation
- ▶ constant learning rate
- ▶ momentum, slowly increasing.

## Results:

- ▶ Trained for 5 minutes, speed 4 miles per hour.
- ▶ ALVINN was able to drive well on a new road it has never seen (in different weather conditions).
- ▶ The maximum speed was limited by the hydraulic controller of the steering wheel, not the learning algorithm.

# ALVINN - weight development



Here  $h_1, \dots, h_5$  are hidden neurons.

# MNIST – handwritten digits recognition

- ▶ Database of labelled images of handwritten digits: 60 000 training examples, 10 000 testing.
- ▶ Dimensions: 28 x 28, digits are centered to the "center of gravity" of pixel values and normalized to fixed size.
- ▶ More at <http://yann.lecun.com/exdb/mnist/>

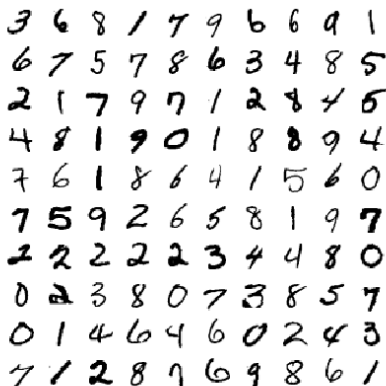


Fig. 4. Size-normalized examples from the MNIST database.

The database is used as a standard benchmark in lots of publications.

# MNIST – handwritten digits recognition

- ▶ Database of labelled images of handwritten digits: 60 000 training examples, 10 000 testing.
- ▶ Dimensions: 28 x 28, digits are centered to the "center of gravity" of pixel values and normalized to fixed size.
- ▶ More at <http://yann.lecun.com/exdb/mnist/>

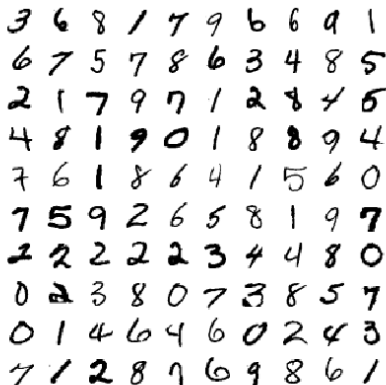


Fig. 4. Size-normalized examples from the MNIST database.

The database is used as a standard benchmark in lots of publications.

Allows comparison of various methods.

One of the best "old" results is the following:

6-layer NN 784-2500-2000-1500-1000-500-10 (on GPU)  
(Ciresan et al. 2010)

**Abstract:** Good old on-line back-propagation for plain multi-layer perceptrons yields a very low 0.35 error rate on the famous MNIST handwritten digits benchmark. All we need to achieve this best result so far are many hidden layers, many neurons per layer, numerous deformed training images, and graphics cards to greatly speed up learning.

A famous application of a learning convolutional network LeNet-1 in 1998.



# MNIST – LeNet1

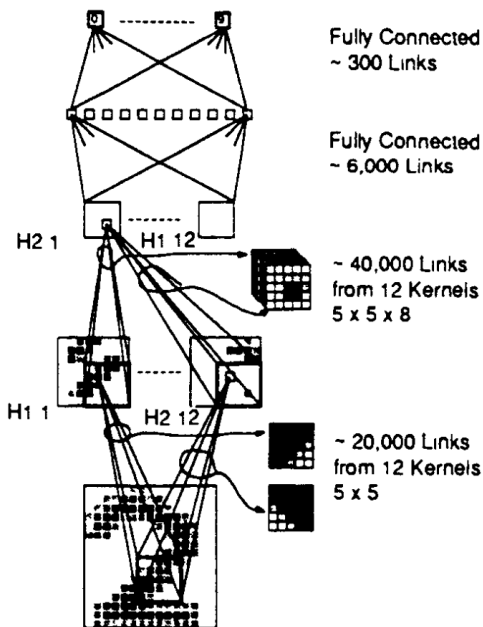
10 Output Units

Layer H3  
30 Hidden Units

Layer H2  
 $12 \times 16 = 192$   
Hidden Units

Layer H1  
 $12 \times 64 = 768$   
Hidden Units

256 Input Units



Interpretation of output:

- ▶ the output neuron with the highest value identifies the digit.
- ▶ the same, but if the two largest neuron values are too close together, the input is rejected (i.e. no answer).

## **Learning:**

Inputs:

- ▶ training on 7291 samples, tested on 2007 samples

## **Results:**

- ▶ error on test set without rejection: 5%
- ▶ error on test set with rejection: 1% (12% rejected)
  
- ▶ compare with dense MLP with 40 hidden neurons: error 1% (19.4% rejected)

# Modern convolutional networks

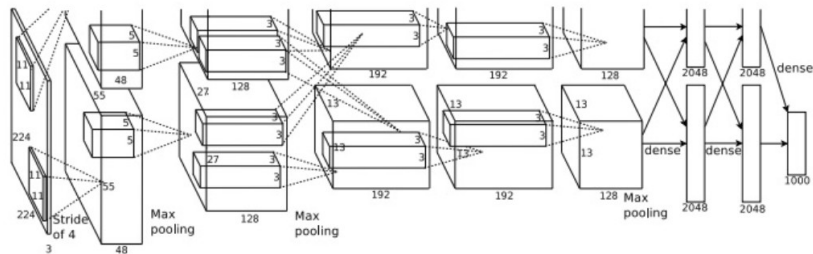
The rest of the lecture is based on the online book Neural Networks and Deep Learning by Michael Nielsen.

<http://neuralnetworksanddeeplearning.com/index.html>

- ▶ Convolutional networks are currently the best networks for image classification.
- ▶ Their common ancestor is LeNet-5 (and other LeNets) from nineties.

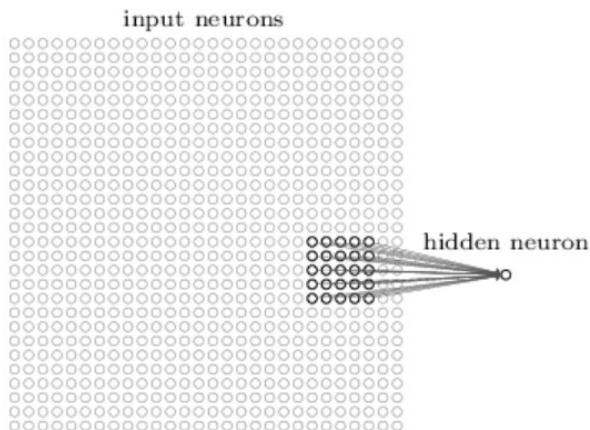
Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998

In 2012 this network made a breakthrough in ILVSCR competition, taking the classification error from around 28% to 16%:



A convolutional network, trained on two GPUs.

# Convolutional networks - local receptive fields

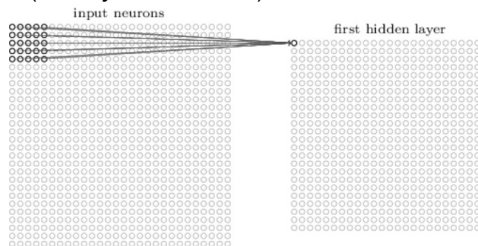


Every neuron is connected with a field of  $k \times k$  (in this case  $5 \times 5$ ) neurons in the lower layer (this field is *receptive field*).

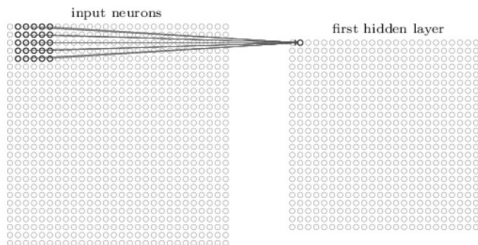
Neuron is "standard": Computes a weighted sum of its inputs, applies an activation function.

# Convolutional networks - stride length

Then we slide the local receptive field over by one pixel to the right (i.e., by one neuron), to connect to a second hidden neuron:

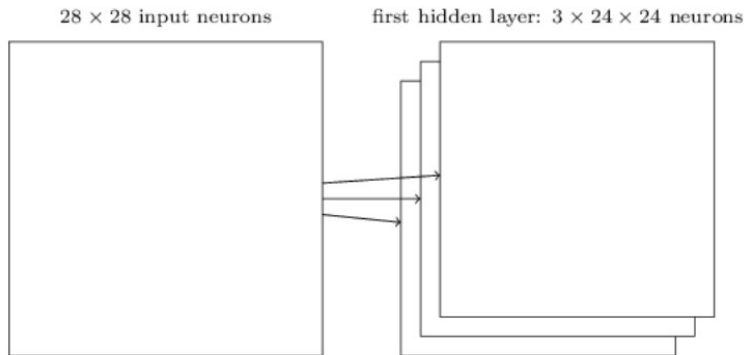


The "size" of the slide is called *stride length*.



The group of all such neurons is *feature map*. all these neurons *share weights and biases*!

# Feature maps

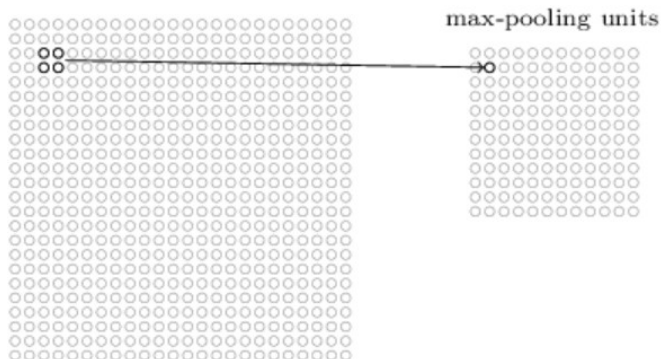


Each feature map represents a property of the input that is supposed to be spatially invariant.

Typically, we consider several feature maps in a single layer.

# Pooling

hidden neurons (output from feature map)

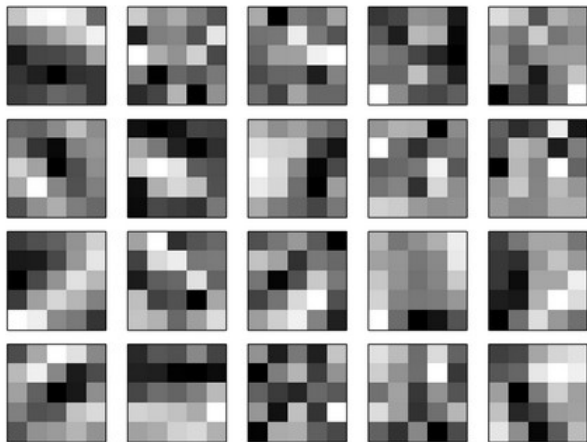


Neurons in the pooling layer compute functions of their receptive fields:

- ▶ **Max-pooling** : maximum of inputs
- ▶ **L2-pooling** : square root of the sum of squares
- ▶ **Average-pooling** : mean
- ▶ ...

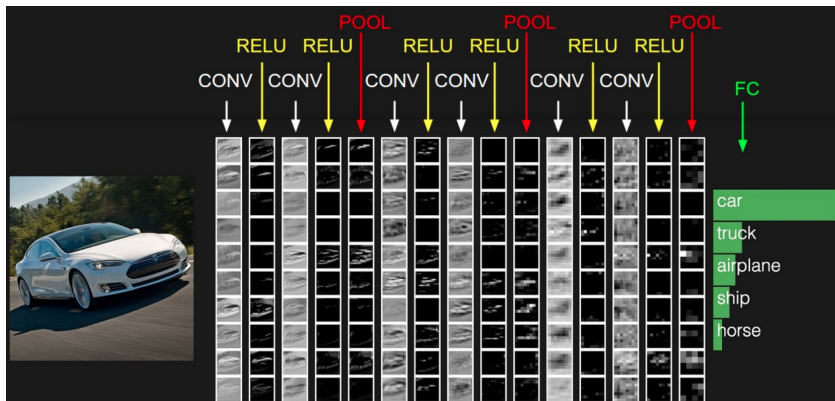


## Trained feature maps

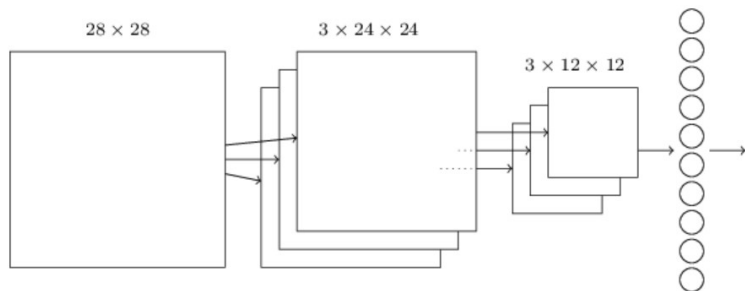


(20 feature maps, receptive fields  $5 \times 5$ )

# Trained feature maps



# Simple convolutional network

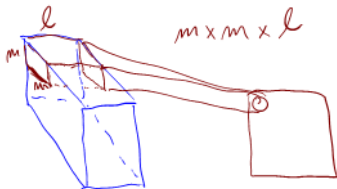
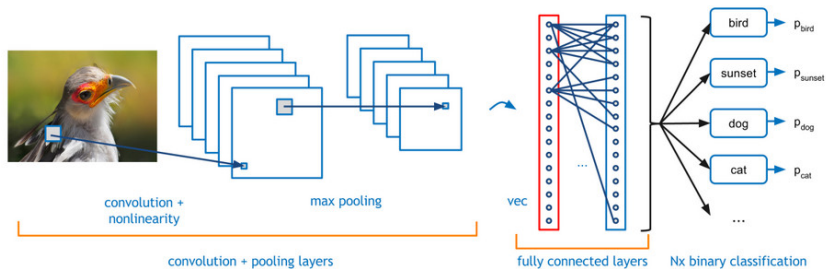


$28 \times 28$  input image, 3 feature maps, each feature map has its own max-pooling (field  $5 \times 5$ , stride = 1), 10 output neurons.

Each neuron in the output layer gets input from each neuron in the pooling layer.

Trained using backprop, which can be easily adapted to convolutional networks.

# Convolutional network

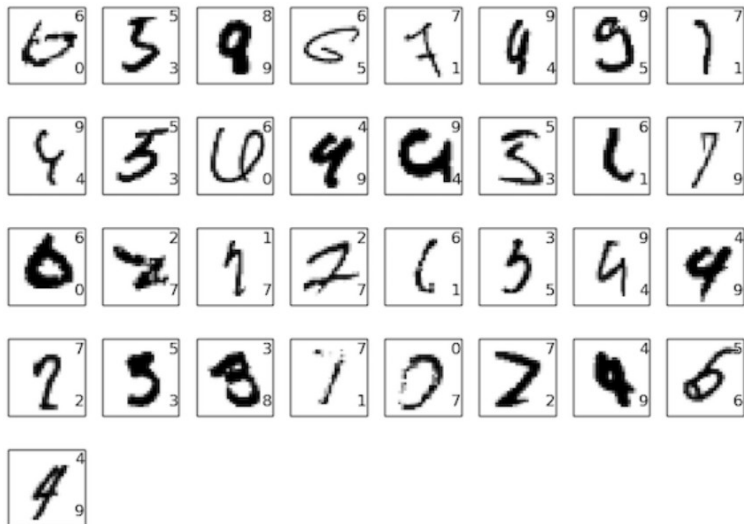


## Simple convolutional network vs MNIST

two convolutional-pooling layers, one 20, second 40 feature maps, two dense (MLP) layers (1000-1000), outputs (10)

- ▶ Activation functions of the feature maps and dense layers: ReLU
- ▶ max-pooling
- ▶ output layer: soft-max
- ▶ Error function: negative log-likelihood (= cross-entropy)
- ▶ Training: SGD, mini-batch size 10
- ▶ learning rate 0.03
- ▶ L2 regularization with "weight"  $\lambda = 0.1$  + dropout with prob. 1/2
- ▶ training for 40 epochs (i.e. every training example is considered 40 times)
- ▶ Expanded dataset: displacement by one pixel to an arbitrary direction.
- ▶ Committee voting of 5 networks.

Out of 10 000 images in the test set, only these 33 have been incorrectly classified:





# ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)

Competition in classification over a subset of images from ImageNet.

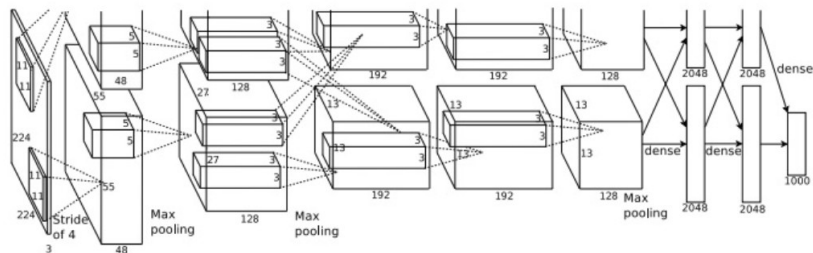
Started in 2010, assisted in breakthrough in image recognition.

Training set 1.2 million images, 1000 classes. Validation set: 50 000, test set: 150 000.

Many images contain more than one object  $\Rightarrow$  model is allowed to choose five classes, the correct label must be among the five. (top-5 criterion).



ImageNet classification with deep convolutional neural networks, by Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton (2012).



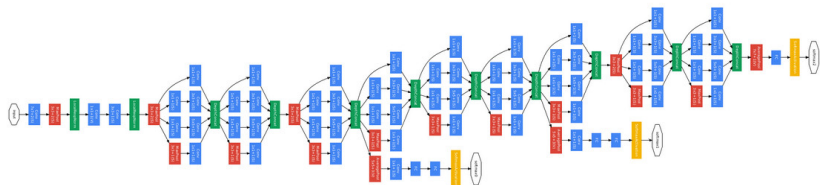
Trained on two GPUs (NVIDIA GeForce GTX 580)

Výsledky:

- ▶ accuracy 84.7% in top-5 (second best algorithm at the time 73.8%)
- ▶ 63.3% "perfect" (top-1) classification

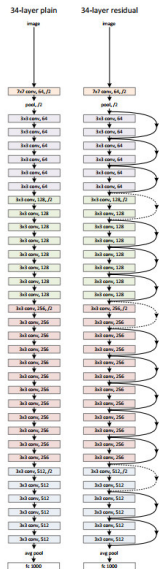
The same set as in 2012, top-5 criterion.

GoogLeNet: deep convolutional network, 22 layers

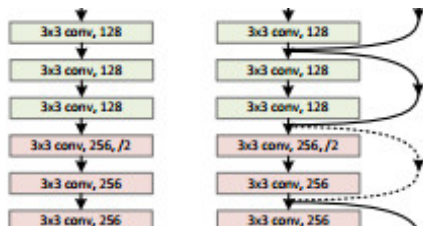


Results:

- ▶ Accuracy 93.33% top-5



- ▶ Deep convolutional network
- ▶ Various numbers of layers, the winner has 152 layers
- ▶ Skip connections implementing residual learning
- ▶ Error **3.57%** in top-5.



Trimps-Soushen (The Third Research Institute of Ministry of Public Security)

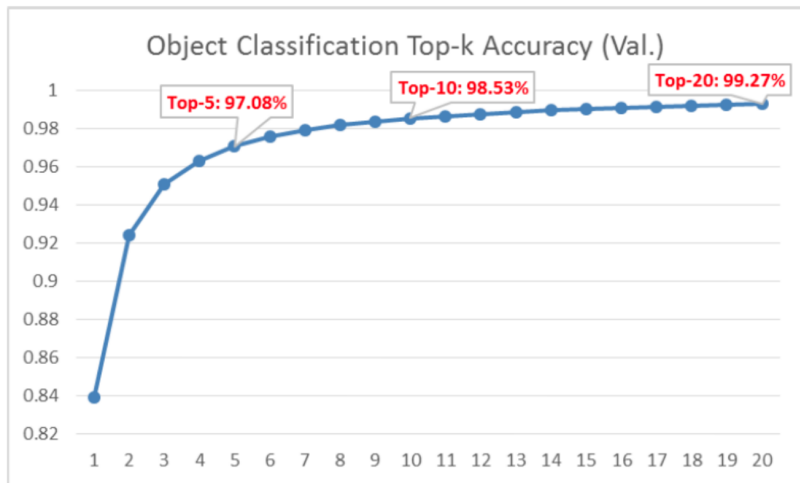
There is no new innovative technology or novelty by Trimps-Soushen.

Ensemble of the pretrained models from Inception-v3, Inception-v4, Inception-ResNet-v2, Pre-Activation ResNet-200, and Wide ResNet (WRN-68-2).

Each of the models are strong at classifying some categories, but also weak at classifying some categories.

Test error: 2.99%

# Top-k accuracy analyzed



<https://towardsdatascience.com/review-trimps-soushen-winner-in-ilsvrc-2016-image-classification-dfbc423111dd>

# Top-20 typical errors

Out of 1458 misclassified images in Top-20:

Error Categories	Numbers	Percentages(%)
Label May Wrong	221	15.16
Multiple Objects (>5)	118	8.09
Non-Obvious Main Object	355	24.35
Confusing Label	206	14.13
Fine-grained Label	258	17.70
Obvious Wrong	234	16.05
Partial Object	66	4.53

<https://towardsdatascience.com/review-trimps-soushen-winner-in-ilsvc-2016-image-classification-dfbc423111dd>

# Top-k accuracy analyzed

Predict:

1 *pencil box*

2 *diaper*

3 *bib*

4 *purse*

5 *running shoe*

Ground Truth:

*sleeping bag*



<https://towardsdatascience.com/review-trimps-soushen-winner-in-ilsvrc-2016-image-classification-dfbc423111dd>

# Top-k accuracy analyzed

Predict:

1 *dock*

2 *submarine*

3 *boathouse*

4 *breakwater*

5 *lifeboat*

Ground Truth:

*paper towel*





# Top-k accuracy analyzed

Predict:

1 *bolete*

2 *earthstar*

3 *gyromitra*

4 *hen of the woods*

5 *mushroom*

Ground Truth:

*stinkhorn*



# Top-k accuracy analyzed

Predict:

1 *apron*

2 *plastic bag*

3 *sleeping bag*

4 *umbrella*

5 *bulletproof vest*

Ground Truth:

*poncho*



# Superhuman convolutional nets?!

Andrej Karpathy: ...the task of labeling images with 5 out of 1000 categories quickly turned out to be extremely challenging, even for some friends in the lab who have been working on ILSVRC and its classes for a while. First we thought we would put it up on [Amazon Mechanical Turk]. Then we thought we could recruit paid undergrads. Then I organized a labeling party of intense labeling effort only among the (expert labelers) in our lab. Then I developed a modified interface that used GoogLeNet predictions to prune the number of categories from 1000 to only about 100. It was still too hard - people kept missing categories and getting up to ranges of 13-15% error rates. In the end I realized that to get anywhere competitively close to GoogLeNet, it was most efficient if I sat down and went through the painfully long training process and the subsequent careful annotation process myself... The labeling happened at a rate of about 1 per minute, but this decreased over time... Some images are easily recognized, while some images (such as those of fine-grained breeds of dogs, birds, or monkeys) can require multiple minutes of concentrated effort. I became very good at identifying breeds of dogs... Based on the sample of images I worked on, the GoogLeNet classification error turned out to be 6.8%... My own error in the end turned out to be 5.1%, approximately 1.7% better.

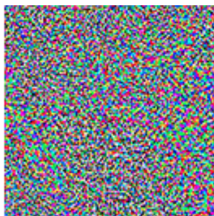
# Does it really work?



**"panda"**

57.7% confidence

+  $\epsilon$



=

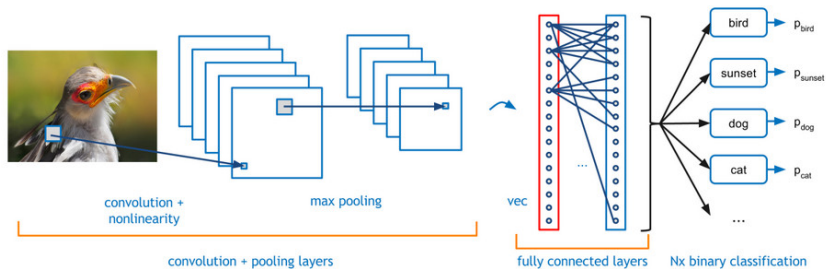


**"gibbon"**

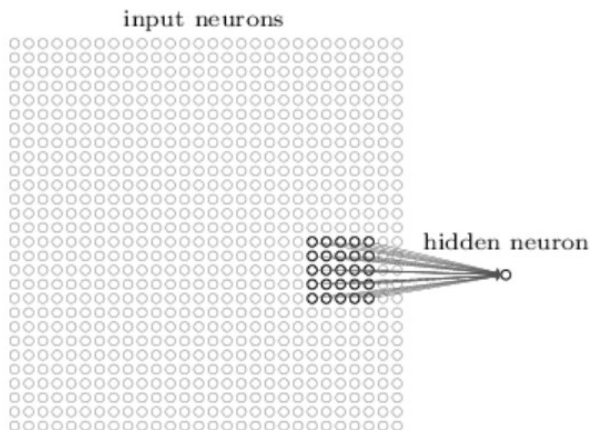
99.3% confidence

# Convolutional networks – theory

# Convolutional network



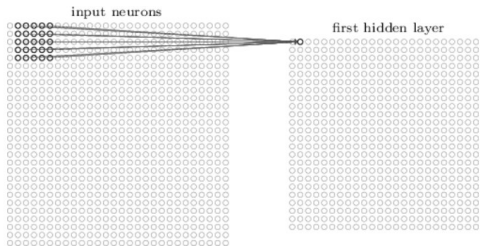
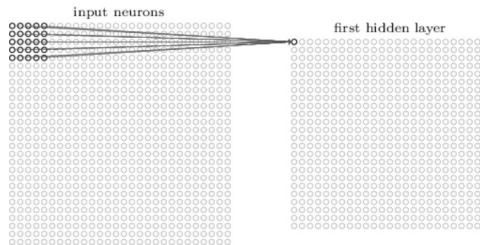
# Convolutional layers



Every neuron is connected with a (typically small) *receptive field* of neurons in the lower layer.

Neuron is "standard": Computes a weighted sum of its inputs, applies an activation function.

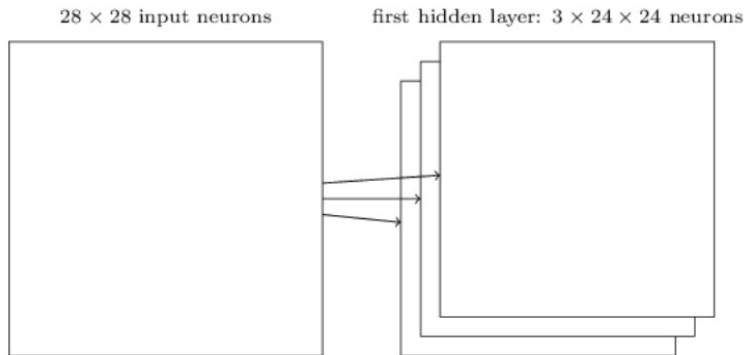
# Convolutional layers



Neurons grouped into *feature maps* sharing weights.



# Convolutional layers

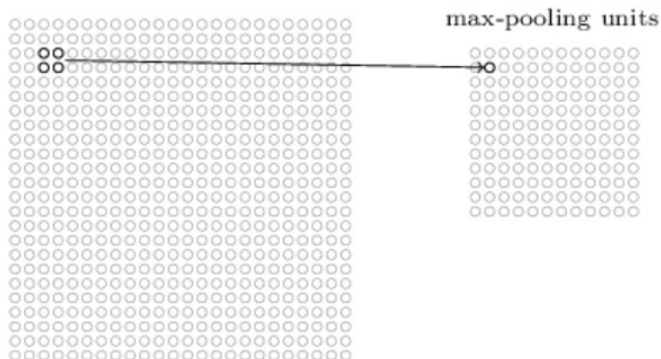


Each feature map represents a property of the input that is supposed to be spatially invariant.

Typically, we consider several feature maps in a single layer.

# Pooling layers

hidden neurons (output from feature map)



Neurons in the pooling layer compute simple functions of their receptive fields (the fields are typically disjoint):

- ▶ **Max-pooling** : maximum of inputs
- ▶ **L2-pooling** : square root of the sum of squares
- ▶ **Average-pooling** : mean
- ▶ ...

## Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$

# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

Several types of layers:

- ▶ **input** layer  $L_0$
- ▶ **dense** layer  $L_m$ : Each neuron of  $L_m$  connected with each neuron of  $L_{m-1}$ .



# Convolutional networks – architecture

Neurons organized in layers,  $L_0, L_1, \dots, L_n$ , connections (typically) only from  $L_m$  to  $L_{m+1}$ .

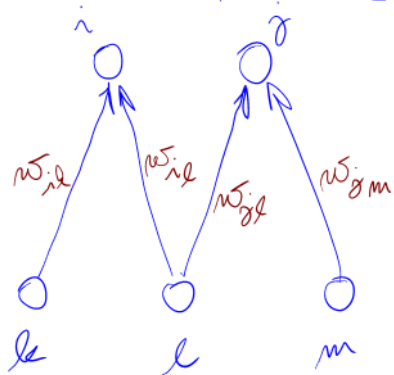
Several types of layers:

- ▶ **input** layer  $L_0$
- ▶ **dense** layer  $L_m$ : Each neuron of  $L_m$  connected with each neuron of  $L_{m-1}$ .
- ▶ **convolutional & pooling** layer  $L_m$ : Contains two sub-layers:
  - ▶ **convolutional layer**: Neurons organized into disjoint **feature maps**, all neurons of a given feature map *share weights* (but have different inputs)
  - ▶ **pooling layer**: Each (convolutional) feature map  $F$  has a corresponding **pooling map**  $P$ . Neurons of  $P$ 
    - ▶ have inputs only from  $F$  (typically few of them),
    - ▶ compute a simple aggregate function (such as max),
    - ▶ have *disjoint inputs*.



## Convolutional networks – architecture

$$[ik] = \{ik, j\ell\} = [j\ell]$$
$$[il] = \{il, jm\} = [jm]$$



$$w_{ik} = w_{j\ell}$$
$$w_{il} = w_{jm}$$

- ▶ Denote
  - ▶  $X$  a set of *input* neurons
  - ▶  $Y$  a set of *output* neurons
  - ▶  $Z$  a set of *all* neurons ( $X, Y \subseteq Z$ )
- ▶ individual neurons denoted by indices  $i, j$  etc.
  - ▶  $\xi_j$  is the inner potential of the neuron  $j$  after the computation stops
  - ▶  $y_j$  is the output of the neuron  $j$  after the computation stops

(define  $y_0 = 1$  is the value of the formal unit input)

- ▶  $w_{ji}$  is the weight of the connection **from**  $i$  **to**  $j$   
(in particular,  $w_{j0}$  is the weight of the connection from the formal unit input, i.e.  $w_{j0} = -b_j$  where  $b_j$  is the bias of the neuron  $j$ )
- ▶  $j^{\leftarrow}$  is a set of all  $i$  such that  $j$  is adjacent from  $i$   
(i.e. there is an arc **to**  $j$  from  $i$ )
- ▶  $j^{\rightarrow}$  is a set of all  $i$  such that  $j$  is adjacent to  $i$   
(i.e. there is an arc **from**  $j$  to  $i$ )
- ▶  $[ji]$  is a set of all connections (i.e. pairs of neurons) sharing the weight  $w_{ji}$ .

# Convolutional networks – activity

- ▶ neurons of dense and convolutional layers:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable):

$$y_j = \sigma_j(\xi_j)$$



# Convolutional networks – activity

- ▶ neurons of dense and convolutional layers:

- ▶ inner potential of neuron  $j$ :

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$

- ▶ activation function  $\sigma_j$  for neuron  $j$  (arbitrary differentiable):

$$y_j = \sigma_j(\xi_j)$$

- ▶ Neurons of pooling layers: Apply the "pooling" function:

- ▶ max-pooling:

$$y_j = \max_{i \in j_{\leftarrow}} y_i$$

- ▶ avg-pooling:

$$y_j = \frac{\sum_{i \in j_{\leftarrow}} y_i}{|j_{\leftarrow}|}$$

A convolutional network is evaluated layer-wise (as MLP), for each  $j \in Y$  we have that  $y_j(\vec{w}, \vec{x})$  is the value of the output neuron  $j$  after evaluating the network with weights  $\vec{w}$  and input  $\vec{x}$ .

# Convolutional networks – learning

## Learning:

- ▶ Given a **training set**  $\mathcal{T}$  of the form

$$\{(\vec{x}_k, \vec{d}_k) \mid k = 1, \dots, p\}$$

$$\vec{d}_k = (0 \dots 0, \overbrace{1, 0 \dots 0}^{1000})$$

Here, every  $\vec{x}_k \in \mathbb{R}^{|X|}$  is an *input vector* and every  $\vec{d}_k \in \mathbb{R}^{|Y|}$  is the desired network output. For every  $j \in Y$ , denote by  $d_{kj}$  the desired output of the neuron  $j$  for a given network input  $\vec{x}_k$  (the vector  $\vec{d}_k$  can be written as  $(d_{kj})_{j \in Y}$ ).

- ▶ **Error function – mean squared error (for example):**

$$E(\vec{w}) = \frac{1}{p} \sum_{k=1}^p E_k(\vec{w})$$

where

$$E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} (y_j(\vec{w}, \vec{x}_k) - d_{kj})^2$$

# Convolutional networks – SGD

The algorithm computes a sequence of weight vectors

$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \dots$

- ▶ weights in  $\vec{w}^{(0)}$  are randomly initialized to values close to 0
- ▶ in the step  $t + 1$  (here  $t = 0, 1, 2 \dots$ ), weights  $\vec{w}^{(t+1)}$  are computed as follows:
  - ▶ Choose (randomly) a set of training examples  $T \subseteq \{1, \dots, p\}$
  - ▶ Compute

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

where

$$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \frac{1}{|T|} \sum_{k \in T} \nabla E_k(\vec{w}^{(t)})$$

Here  $T$  is a *minibatch* (of a fixed size),

- ▶  $0 < \varepsilon(t) \leq 1$  is a *learning rate* in step  $t + 1$
- ▶  $\nabla E_k(\vec{w}^{(t)})$  is the gradient of the error of the example  $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially. **Epoch** consists of one round through all data.

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

First, switch from derivatives w.r.t.  $w_{ji}$  to derivatives w.r.t.  $y_j$ :

- ▶ Recall that for every  $w_{ji}$  where  $j$  is in a dense layer, i.e. does not share weights:

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

# Backprop

Recall that  $\nabla E_k(\vec{w}^{(t)})$  is a vector of all partial derivatives of the form  $\frac{\partial E_k}{\partial w_{ji}}$ .

How to compute  $\frac{\partial E_k}{\partial w_{ji}}$  ?

First, switch from derivatives w.r.t.  $w_{ji}$  to derivatives w.r.t.  $y_j$ :

- ▶ Recall that for every  $w_{ji}$  where  $j$  is in a dense layer, i.e. does not share weights:

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

- ▶ Now for every  $w_{ji}$  where  $j$  is in a convolutional layer:

$$\frac{\partial E_k}{\partial w_{ji}} = \sum_{r \in [ji]} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot y_\ell$$

- ▶ Neurons of pooling layers do not have weights.

# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

- ▶ for every  $j \in Z \setminus Y$  such that  $j^\rightarrow$  is either a dense layer, or a convolutional layer:

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^\rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj}$$



# Backprop

Now compute derivatives w.r.t.  $y_j$ :

- ▶ for every  $j \in Y$ :

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$$

This holds for the squared error, for other error functions the derivative w.r.t. outputs will be different.

- ▶ for every  $j \in Z \setminus Y$  such that  $j \rightarrow$  is either a dense layer, or a convolutional layer:

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j \rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj}$$

*arg max<sub>r ∈ i ←</sub> y<sub>r</sub> = i'*

- ▶ for every  $j \in Z \setminus Y$  such that  $j \rightarrow$  is max-pooling: Then  $j \rightarrow = \{i\}$  for a single "max" neuron and we have

$$\begin{pmatrix} x^1 & x^{i'} \\ x^n & x^n \\ -2 & -3 \end{pmatrix} \rightarrow \frac{\partial E_k}{\partial y_j} = \begin{cases} \frac{\partial E_k}{\partial y_i} & \text{if } j = \arg \max_{r \in i \leftarrow} y_r \\ 0 & \text{otherwise} \end{cases}$$

*max({1, 3, 2}) = 3*  
*arg max({1, 3, 2}) = 2*

i.e. gradient can be propagated from the output layer downwards as in MLP.

## Convolutional networks – summary

- ▶ Conv. nets. are nowadays the most used networks in image processing (and also in other areas where input has some local, "spatially" invariant properties)
- ▶ Typically trained using backpropagation.
- ▶ Due to the weight sharing allow (very) deep architectures.
- ▶ Typically extended with more adjustments and tricks in their topologies.