# 732A99/TDDE01 Machine Learning
## Lecture 3d Block 1: Deep Learning

Jose M. Peña
IDA, Linköping University, Sweden

# Contents

# Literature

- Main sources
  - Bengio, Y. Learning Deep Architectures for AI. *Foundations and Trends in Machine Learning*, 2:1-127, 2009. Chapters 1-3, 4.2, 4.5-4.6, 6.2-6.3.
  - Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006. Section 5.5.6.
  - LeCun, Y., Bengio, Y. and Hinton, G. Deep Learning. *Nature*, 521:436-44, 2015.
- Additional source
  - Goodfellow, I., Bengio, Y. and Courville, A. *Deep Learning*. Book in preparation for MIT Press, 2016. Available at www.deeplearningbook.org. Chapters 1, 6.

# Limitations of Neural Networks

## Theorem (Universal approximation theorem)

*For every continuous function $f : [a, b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

# Limitations of Neural Networks

### Theorem (Universal approximation theorem)

*For every continuous function $f : [a, b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

### Theorem (Universal classification theorem)

*Let $\mathcal{C}^{(k)}$ contain all classifiers defined by NNs of one hidden layer with $k$ hidden units and the sigmoid activation function. Then, for any distribution $p(\mathbf{x}, t)$,*

$$\lim_{k \to \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

*where $L()$ is the 0/1 loss function.*

## Limitations of Neural Networks

### Theorem (Universal approximation theorem)

*For every continuous function $f : [a, b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

### Theorem (Universal classification theorem)

*Let $\mathcal{C}^{(k)}$ contain all classifiers defined by NNs of one hidden layer with $k$ hidden units and the sigmoid activation function. Then, for any distribution $p(\mathbf{x}, t)$,*

$$\lim_{k \to \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

*where $L()$ is the 0/1 loss function.*

- How many hidden units has such a NN ?

# Limitations of Neural Networks

### Theorem (Universal approximation theorem)

*For every continuous function $f : [a, b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

### Theorem (Universal classification theorem)

*Let $\mathcal{C}^{(k)}$ contain all classifiers defined by NNs of one hidden layer with $k$ hidden units and the sigmoid activation function. Then, for any distribution $p(\mathbf{x}, t)$,*

$$\lim_{k \to \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

*where $L()$ is the 0/1 loss function.*

- How many hidden units has such a NN ?
- How much data do we need to learn such a NN (and avoid overfitting) via the backpropagation algorithm ?

# Limitations of Neural Networks

## Theorem (Universal approximation theorem)

*For every continuous function $f : [a, b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

## Theorem (Universal classification theorem)

*Let $\mathcal{C}^{(k)}$ contain all classifiers defined by NNs of one hidden layer with $k$ hidden units and the sigmoid activation function. Then, for any distribution $p(\mathbf{x}, t)$,*

$$\lim_{k \to \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

*where $L()$ is the 0/1 loss function.*

- ▸ How many hidden units has such a NN ?
- ▸ How much data do we need to learn such a NN (and avoid overfitting) via the backpropagation algorithm ?
- ▸ How fast does the backpropagation algorithm converge to such a NN ? Assuming that it does not get trapped in a local minimum...

# Limitations of Neural Networks

### Theorem (Universal approximation theorem)

*For every continuous function $f : [a,b]^D \to \mathbb{R}$ and for every $\epsilon > 0$, there exists a NN with one hidden layer such that*

$$\sup_{\mathbf{x} \in [a,b]^D} |f(\mathbf{x}) - y(\mathbf{x})| < \epsilon$$

### Theorem (Universal classification theorem)

*Let $\mathcal{C}^{(k)}$ contain all classifiers defined by NNs of one hidden layer with $k$ hidden units and the sigmoid activation function. Then, for any distribution $p(\mathbf{x}, t)$,*

$$\lim_{k \to \infty} \inf_{y \in \mathcal{C}^{(k)}} L(y(\mathbf{x})) - L(p(t|\mathbf{x})) = 0$$

*where $L()$ is the 0/1 loss function.*

- How many hidden units has such a NN ?
- How much data do we need to learn such a NN (and avoid overfitting) via the backpropagation algorithm ?
- How fast does the backpropagation algorithm converge to such a NN ? Assuming that it does not get trapped in a local minimum...
- The answer to the last two questions depends on the first: More hidden units implies more training time and higher generalization error.

# Limitations of Neural Networks

- How many hidden units does the NN need ?

# Limitations of Neural Networks

- ▶ How many hidden units does the NN need ?
- ▶ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.

## Limitations of Neural Networks

- ▸ How many hidden units does the NN need ?
- ▸ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.
- ▸ For most Boolean functions, the size of the circuit is exponential in the size of the input.

# Limitations of Neural Networks

- ▸ How many hidden units does the NN need ?
- ▸ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.
- ▸ For most Boolean functions, the size of the circuit is exponential in the size of the input.
- ▸ However, there are Boolean functions that have a polynomial-size circuit of depth $k$ and an exponential-size circuit of depth $k - 1$.
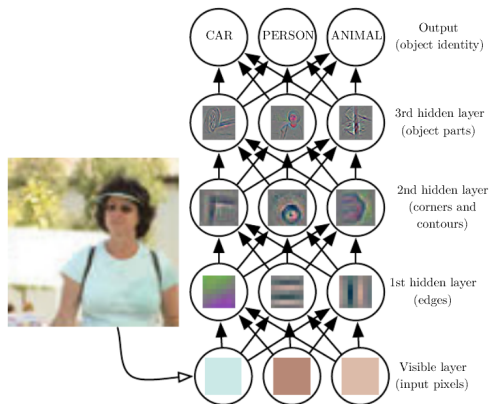
# Limitations of Neural Networks

- ▸ How many hidden units does the NN need ?
- ▸ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.
- ▸ For most Boolean functions, the size of the circuit is exponential in the size of the input.
- ▸ However, there are Boolean functions that have a polynomial-size circuit of depth $k$ and an exponential-size circuit of depth $k - 1$.
- ▸ Then, there is no universally right depth. Ideally, we should let the data determine the right depth.

# Limitations of Neural Networks

- ▸ How many hidden units does the NN need ?
- ▸ Any Boolean function can be written in disjunctive normal form (OR of ANDs) or conjunctive normal form (AND of ORs). This is a depth-two logical circuit.
- ▸ For most Boolean functions, the size of the circuit is exponential in the size of the input.
- ▸ However, there are Boolean functions that have a polynomial-size circuit of depth $k$ and an exponential-size circuit of depth $k - 1$.
- ▸ Then, there is no universally right depth. Ideally, we should let the data determine the right depth.

## Theorem (No free lunch theorem)

*For any algorithm, good performance on some problems comes at the expense of bad performance on some others.*

# Deep Neural Networks



- A deep NN is a function that maps input to output.
- The mapping is formed by composing many simpler functions.
- Each layer provides a new representation of the input, i.e. complex concepts are built from simpler ones.
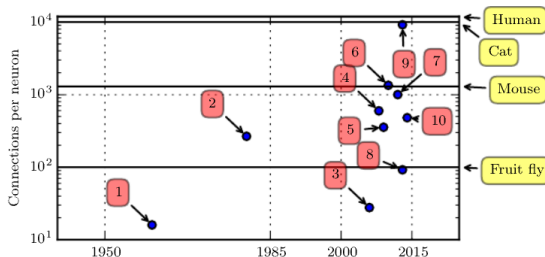- The representation is learned automatically from data.

# Deep Neural Networks



Figure 1.10: Initially, the number of connections between neurons in artificial neural networks was limited by hardware capabilities. Today, the number of connections between neurons is mostly a design consideration. Some artificial neural networks have nearly as many connections per neuron as a cat, and it is quite common for other neural networks to have as many connections per neuron as smaller mammals like mice. Even the human brain does not have an exorbitant amount of connections per neuron. Biological neural network sizes from Wikipedia (2015).

1. Adaptive linear element (Widrow and Hoff, 1960)
2. Neocognitron (Fukushima, 1980)
3. GPU-accelerated convolutional network (Chellapilla et al., 2006)
4. Deep Boltzmann machine (Salakhutdinov and Hinton, 2009a)
5. Unsupervised convolutional network (Jarrett et al., 2009)
6. GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
7. Distributed autoencoder (Le et al., 2012)
8. Multi-GPU convolutional network (Krizhevsky et al., 2012)
9. COTS HPC unsupervised convolutional network (Coates et al., 2013)
10. GoogLeNet (Szegedy et al., 2014a)
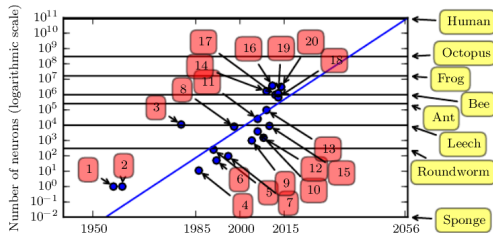
# Deep Neural Networks



Figure 1.11: Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from Wikipedia (2015).

1. Perceptron (Rosenblatt, 1958, 1962)
2. Adaptive linear element (Widrow and Hoff, 1960)
3. Neocognitron (Fukushima, 1980)
4. Early back-propagation network (Rumelhart et al., 1986b)
5. Recurrent neural network for speech recognition (Robinson and Fallside, 1991)
6. Multilayer perceptron for speech recognition (Bengio et al., 1991)
7. Mean field sigmoid belief network (Saul et al., 1996)
8. LeNet-5 (LeCun et al., 1998b)
9. Echo state network (Jaeger and Haas, 2004)
10. Deep belief network (Hinton et al., 2006)
11. GPU-accelerated convolutional network (Chellapilla et al., 2006)
12. Deep Boltzmann machine (Salakhutdinov and Hinton, 2009a)
13. GPU-accelerated deep belief network (Raina et al., 2009)
14. Unsupervised convolutional network (Jarrett et al., 2009)
15. GPU-accelerated multilayer perceptron (Ciresan et al., 2010)
16. OMP-1 network (Coates and Ng, 2011)
17. Distributed autoencoder (Le et al., 2012)
18. Multi-GPU convolutional network (Krizhevsky et al., 2012)
19. COTS HPC unsupervised convolutional network (Coates et al., 2013)
20. GoogLeNet (Szegedy et al., 2014a)

22 layers DNN, but 12 times fewer weights than DNN 19

# Deep Neural Networks

- ▸ Training DNNs is difficult:
  - ▸ Typically, poorer generalization than (shallow) NNs.

# Deep Neural Networks

▸ Training DNNs is difficult:
  ▸ Typically, poorer generalization than (shallow) NNs.
  ▸ The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1]$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.

# Deep Neural Networks

- ▸ Training DNNs is difficult:
    - ▸ Typically, poorer generalization than (shallow) NNs.
    - ▸ The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1)$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
    - ▸ There may be larger plateaus and many more local minima than with NNs.

# Deep Neural Networks
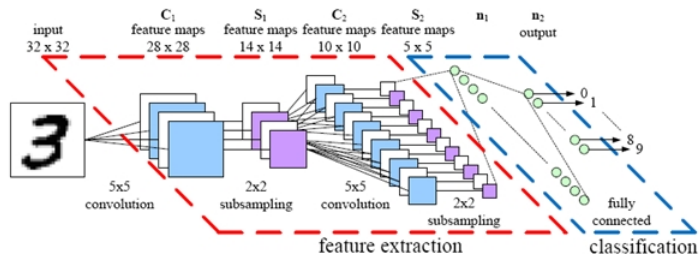
- Training DNNs is difficult:
    - Typically, poorer generalization than (shallow) NNs.
    - The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1]$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
    - There may be larger plateaus and many more local minima than with NNs.

- Training DNNs is doable:
    - Convolutional networks, particularly suitable for image processing.

# Deep Neural Networks

- Training DNNs is difficult:
  - Typically, poorer generalization than (shallow) NNs.
  - The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1)$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
  - There may be larger plateaus and many more local minima than with NNs.

- Training DNNs is doable:
  - Convolutional networks, particularly suitable for image processing.
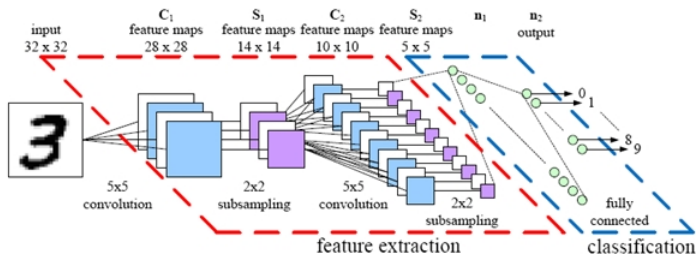  - Rectifier activation function, a new activation function.

# Deep Neural Networks

- Training DNNs is difficult:
  - Typically, poorer generalization than (shallow) NNs.
  - The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1)$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
  - There may be larger plateaus and many more local minima than with NNs.

- Training DNNs is doable:
  - Convolutional networks, particularly suitable for image processing.
  - Rectifier activation function, a new activation function.
  - Layer-wise pre-training, to find a good starting point for training.

# Deep Neural Networks

- Training DNNs is difficult:
  - Typically, poorer generalization than (shallow) NNs.
  - The gradient may vanish/explode as we move away from the output layer, due to multiplying small/big quantities. E.g. the gradient of $\sigma$ and tanh is in $[0, 1)$. So, they may only suffer the gradient vanishing problem. Other activations functions may suffer the gradient exploding problem.
  - There may be larger plateaus and many more local minima than with NNs.

- Training DNNs is doable:
  - Convolutional networks, particularly suitable for image processing.
  - Rectifier activation function, a new activation function.
  - Layer-wise pre-training, to find a good starting point for training.

- In addition to performance, the computational demands of the training must be considered, e.g. CPU, GPU, memory, parallelism, etc.
  - The authors state that GoogLeNet was trained "using modest amount of model and data-parallelism. Although we used a CPU based implementation only, a rough estimate suggests that the GoogLeNet network could be trained to convergence using few high-end GPUs within a week, the main limitation being the memory usage".
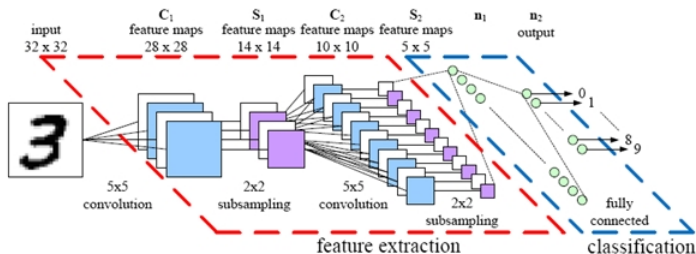
# Convolutional Networks



‣ DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
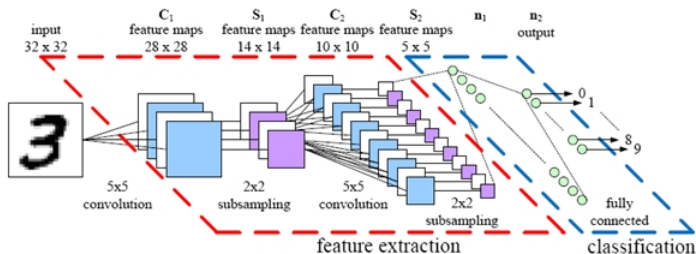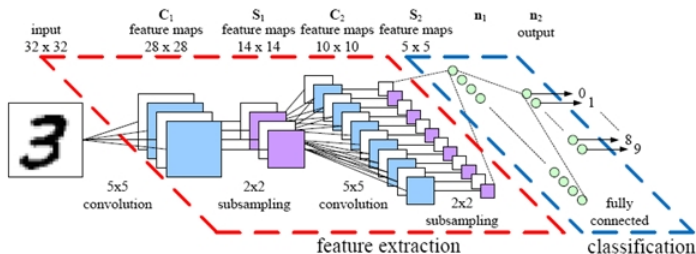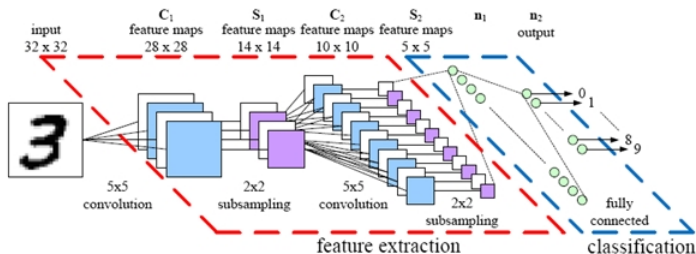
# Convolutional Networks



input 32 x 32 — $C_1$ feature maps 28 x 28 — $S_1$ feature maps 14 x 14 — $C_2$ feature maps 10 x 10 — $S_2$ feature maps 5 x 5 — $n_1$ — $n_2$ output

5x5 convolution — 2x2 subsampling — 5x5 convolution — 2x2 subsampling — fully connected

feature extraction — classification

▸ DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
▸ Convolution: Detection of local features, e.g. $a_j$ is computed from a 5x5 pixel patch of the image.
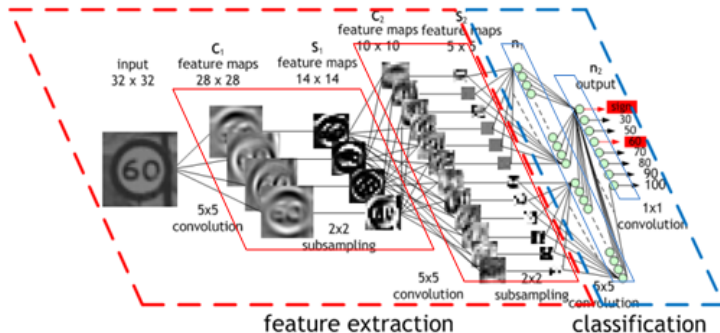
# Convolutional Networks



- DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
- Convolution: Detection of local features, e.g. $a_j$ is computed from a 5x5 pixel patch of the image.
- To achieve invariance, the units in the convolution layer share the same activation function and weights.
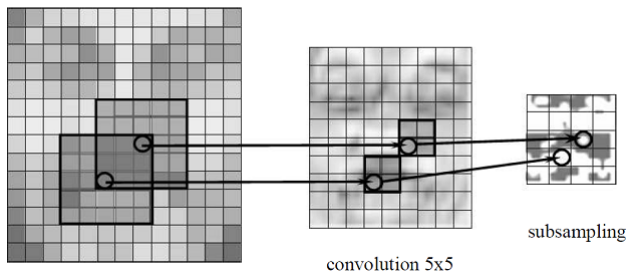
# Convolutional Networks



- DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
- Convolution: Detection of local features, e.g. $a_j$ is computed from a 5x5 pixel patch of the image.
- To achieve invariance, the units in the convolution layer share the same activation function and weights.
- Subsampling: Combination of local features into higher-order features, e.g. $a_k$ is compute from a 2x2 pixel patch of the convoluted image.

# Convolutional Networks



- DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
- Convolution: Detection of local features, e.g. $a_j$ is computed from a 5x5 pixel patch of the image.
- To achieve invariance, the units in the convolution layer share the same activation function and weights.
- Subsampling: Combination of local features into higher-order features, e.g. $a_k$ is compute from a 2x2 pixel patch of the convoluted image.
- There are several feature maps in each layer, to compensate the reduction in resolution by increasing in the number of features being detected.

# Convolutional Networks



- ▸ DNNs suitable for image recognition, since they exhibit invariance to translation, scaling, rotations, and warping.
- ▸ Convolution: Detection of local features, e.g. $a_j$ is computed from a 5x5 pixel patch of the image.
- ▸ To achieve invariance, the units in the convolution layer share the same activation function and weights.
- ▸ Subsampling: Combination of local features into higher-order features, e.g. $a_k$ is compute from a 2x2 pixel patch of the convoluted image.
- ▸ There are several feature maps in each layer, to compensate the reduction in resolution by increasing in the number of features being detected.
- ▸ The final layer is a regular NN for classification.

# Convolutional Networks



subsampling

convolution 5x5



feature extraction    classification

## Convolutional Networks

- DNNs allow increased depth because
  - they are sparse, which allows the gradient to propagate further, and

## Convolutional Networks

- DNNs allow increased depth because
  - they are sparse, which allows the gradient to propagate further, and
  - they have relatively few weights to fit due to feature locality and weight sharing.

## Convolutional Networks

- DNNs allow increased depth because
  - they are sparse, which allows the gradient to propagate further, and
  - they have relatively few weights to fit due to feature locality and weight sharing.
- The backpropagation algorithm needs to be adapted, by modifying the derivatives with respect to the weights in each convolution layer $m$.

## Convolutional Networks

- DNNs allow increased depth because
  - they are sparse, which allows the gradient to propagate further, and
  - they have relatively few weights to fit due to feature locality and weight sharing.
- The backpropagation algorithm needs to be adapted, by modifying the derivatives with respect to the weights in each convolution layer $m$.
- Since $E_n$ depends on $w_i^{(m)}$ only via $a_j^{(m)}$, and $a_j^{(m)} = \sum_{i \in L_j^{(m)}} w_i^{(m)} z_i^{(m-1)}$ where $L_j^{(m)}$ is the set of indexes of the input units, then

$$\frac{\partial E_n}{\partial w_i^{(m)}} = \sum_j \frac{\partial E_n}{\partial a_j^{(m)}} \frac{\partial a_j^{(m)}}{\partial w_i^{(m)}} = \sum_j \delta_j^{(m)} z_i^{(m-1)}$$

## Convolutional Networks

- DNNs allow increased depth because
  - they are sparse, which allows the gradient to propagate further, and
  - they have relatively few weights to fit due to feature locality and weight sharing.
- The backpropagation algorithm needs to be adapted, by modifying the derivatives with respect to the weights in each convolution layer $m$.
- Since $E_n$ depends on $w_i^{(m)}$ only via $a_j^{(m)}$, and $a_j^{(m)} = \sum_{i \in L_j^{(m)}} w_i^{(m)} z_i^{(m-1)}$

  where $L_j^{(m)}$ is the set of indexes of the input units, then

$$\frac{\partial E_n}{\partial w_i^{(m)}} = \sum_j \frac{\partial E_n}{\partial a_j^{(m)}} \frac{\partial a_j^{(m)}}{\partial w_i^{(m)}} = \sum_j \delta_j^{(m)} z_i^{(m-1)}$$

- Note that $w_i^{(m)}$ does not depend on $j$ by weight sharing, whereas $i \in L_j^{(m)}$ by feature locality.
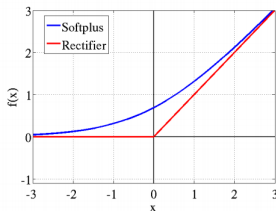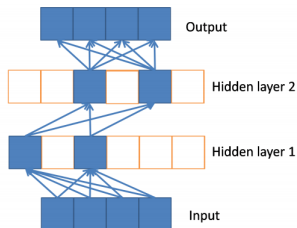
# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

▸ *rectifier*$(x)$ = max$\{0, x\}$, i.e. hidden units are off or operating in a linear regime.
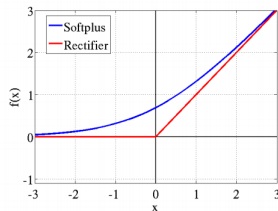
# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- *rectifier*$(x)$ = max$\{0, x\}$, i.e. hidden units are off or operating in a linear regime.
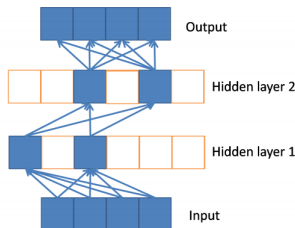- The most popular choice nowadays.

# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- *rectifier*$(x)$ = max$\{0, x\}$, i.e. hidden units are off or operating in a linear regime.
- The most popular choice nowadays.
- Sparsity promoting: Uniform initialization of the weights implies that around 50 % of the hidden units are off.
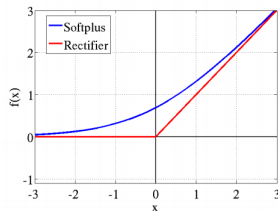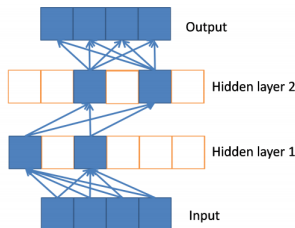
# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- *rectifier*$(x)$ = max$\{0, x\}$, i.e. hidden units are off or operating in a linear regime.
- The most popular choice nowadays.
- Sparsity promoting: Uniform initialization of the weights implies that around 50 % of the hidden units are off.
- Piece-wise linear mapping: The input selects which hidden units are active, and the output is a liner function of the input in the selected hidden units.
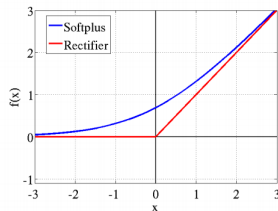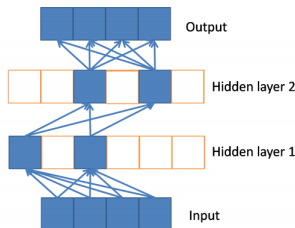
# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- It simplifies the backpropagation algorithm as $h'(a_j) = 1$ for the selected units. So, there is no gradient vanishing on the paths of selected units. Compare with the sigmoid or hyperbolic tangent, for which
  - the gradient is smaller than one, or
  - even zero due to saturation.

# Rectifier Activation Function
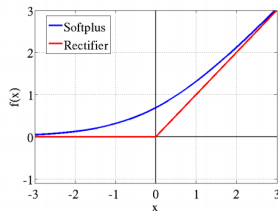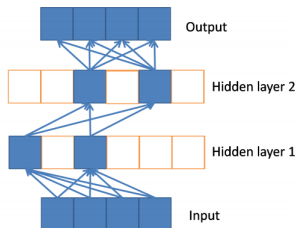


Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- It simplifies the backpropagation algorithm as $h'(a_j) = 1$ for the selected units. So, there is no gradient vanishing on the paths of selected units. Compare with the sigmoid or hyperbolic tangent, for which
  - the gradient is smaller than one, or
  - even zero due to saturation.
- Note that $h'(0)$ does not exist since $h'_+(0) \neq h'_-(0)$. We can get around this problem by simply returning one of two one-sided derivatives. Or using a generalization of the rectifier function.
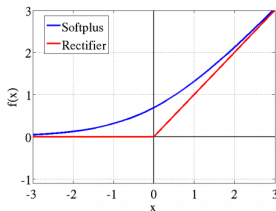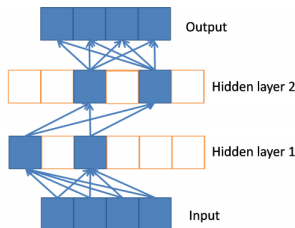
# Rectifier Activation Function



Figure 2: *Left:* **Sparse propagation of activations and gradients in a network of rectifier units.** The input selects a subset of active neurons and computation is linear in this subset. *Right:* **Rectifier and softplus activation functions.** The second one is a smooth version of the first.

- It simplifies the backpropagation algorithm as $h'(a_j) = 1$ for the selected units. So, there is no gradient vanishing on the paths of selected units. Compare with the sigmoid or hyperbolic tangent, for which
  - the gradient is smaller than one, or
  - even zero due to saturation.

- Note that $h'(0)$ does not exist since $h'_+(0) \neq h'_-(0)$. We can get around this problem by simply returning one of two one-sided derivatives. Or using a generalization of the rectifier function.

- Regularization is typically added to prevent numerical problems due to the activation being unbounded, e.g. when forward propagating.

# Layer-Wise Pre-Training

- The pre-training aims to find a good starting point for the subsequent run of the backpropagation algorithm.

# Layer-Wise Pre-Training

- The pre-training aims to find a good starting point for the subsequent run of the backpropagation algorithm.

- Supervised version:
  1. Train each layer of the DNN as if it was the hidden layer in a depth-two NN. As input, use the output of the last of the previously trained layers. As output, use the original classification or regression function.
  2. Run the backpropagation algorithm to fine-tune the weights.

# Layer-Wise Pre-Training

- The pre-training aims to find a good starting point for the subsequent run of the backpropagation algorithm.

- Supervised version:
  1. Train each layer of the DNN as if it was the hidden layer in a depth-two NN. As input, use the output of the last of the previously trained layers. As output, use the original classification or regression function.
  2. Run the backpropagation algorithm to fine-tune the weights.

- Unsupervised version: Similar to the supervised one but the hidden layers (except the last one) are trained to learn an encoding of the output of the previous layer, instead of the original classification or regression function.

# Summary

- Direct application of the backpropagation algorithm to DNNs produces poor results.
- Convolutional networks: It makes the backpropagation algorithm more efficient by using local features and weight sharing. This also achieves invariance, which is particularly important for image processing.
- Rectifier activation function: Free of gradient vanishing problem and it simplifies the backpropagation algorithm.
- Layer-wise pre-training: Heuristic weight initialization to alleviate the local optima problem.