

Perceptron (Recap)
oooooo

Multi-Layer Perceptron
oooooooo
ooo

Training
oooooooooooo

Practical considerations
oooo

Different tastes of NNs
oooooo

Summary
ooo

Cognitive Algorithms Lecture 6

Multilayer Perceptrons

Klaus-Robert Müller, Jan Gerken, **Lorenz Vaitl**,
Augustin Krause, Ken Schreiber

Berlin Institute of Technology
Dept. Machine Learning

Perceptron (Recap)
oooooo

Multi-Layer Perceptron
oooooooo
ooo

Training
oooooooooooo

Practical considerations
oooo

Different tastes of NNs
oooooo

Summary
ooo

Cognitive Algorithms Lecture 6

Multilayer Perceptrons

Klaus-Robert Müller, Jan Gerken, **Lorenz Vaitl**,
Augustin Krause, Ken Schreiber

Berlin Institute of Technology
Dept. Machine Learning

Perceptron (Recap)
oooooo

Multi-Layer Perceptron
oooooooo
ooo

Training
oooooooooooo

Practical considerations
oooo

Different tastes of NNs
oooooo

Summary
ooo

Survey evaluation



Last time

Unsupervised Learning:

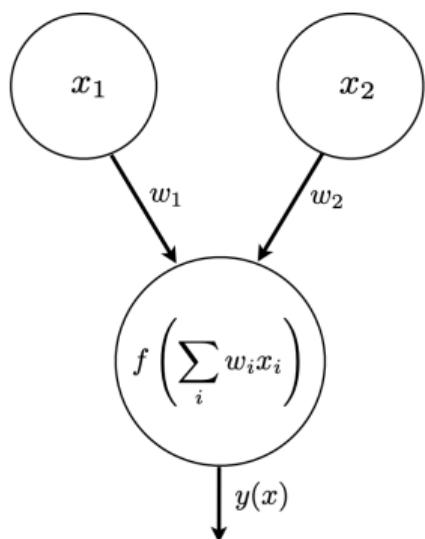
Dimensionality Reduction

- PCA
- Nonnegative Matrix Factorization

Clustering

- K-Means

Single(-Layer) Perceptron



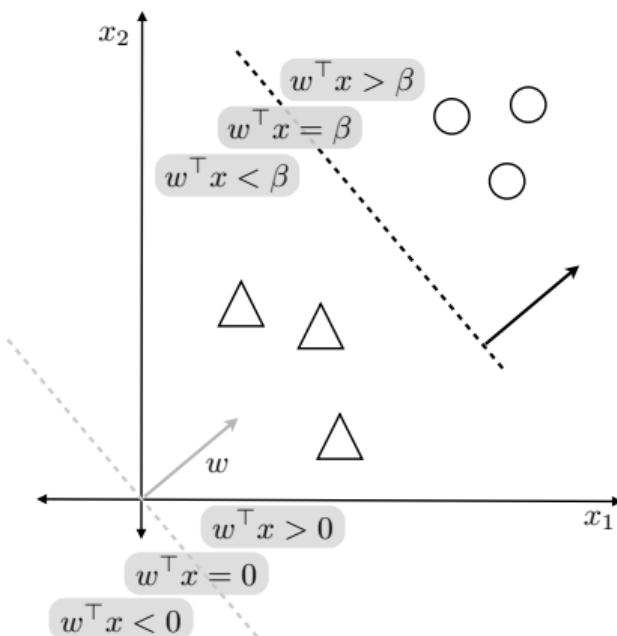
Input nodes x_i receive information

Inputs are multiplied with a weighting factor w_i and summed up

Integrated input is mapped through some (non-linear) function $f(\cdot)$

$$\text{e.g. } f(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{x} \text{ else} \end{cases}$$

Linear Classification and the Perceptron



$$\mathbf{w}^\top \mathbf{x} - \beta = \begin{cases} > 0 & \text{if } \mathbf{x} \text{ belongs to } o \\ < 0 & \text{if } \mathbf{x} \text{ belongs to } \Delta \end{cases}$$

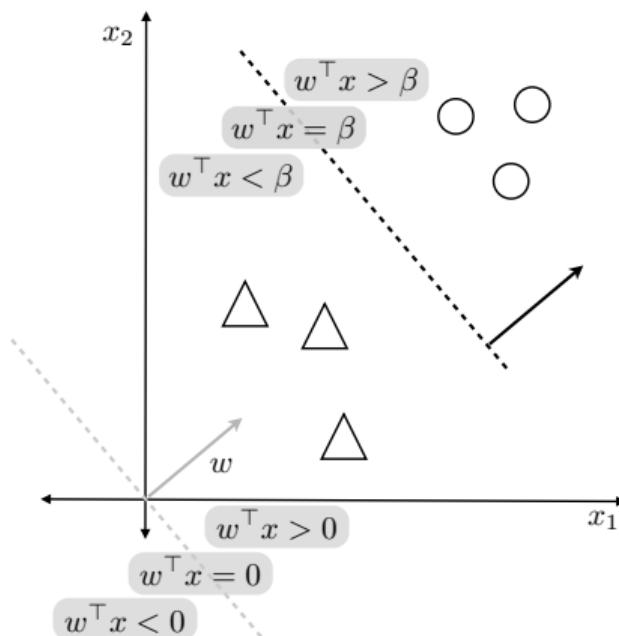
The *offset* β can be included in \mathbf{w}

$$\tilde{\mathbf{x}} \leftarrow \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \quad \tilde{\mathbf{w}} \leftarrow \begin{bmatrix} -\beta \\ \mathbf{w} \end{bmatrix}$$

such that

$$\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}} = \mathbf{w}^\top \mathbf{x} - \beta.$$

Linear Classification and the Perceptron



What is a good w ?

→ We need an **error function** that tells us how good w is.

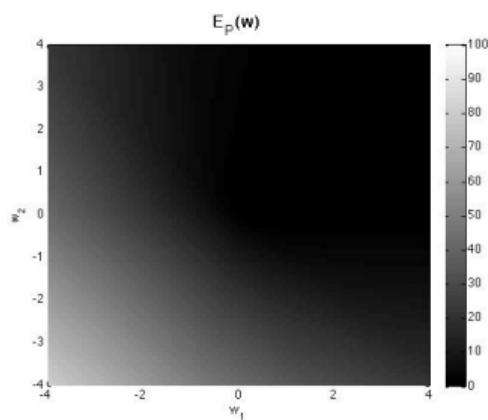
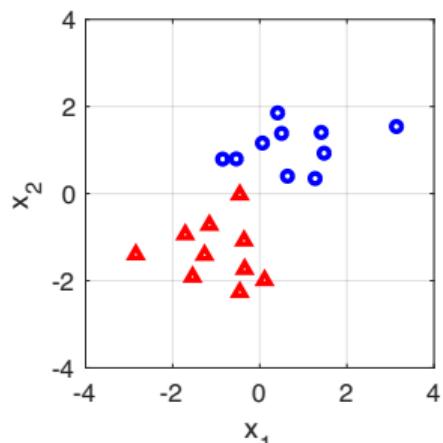
Then we chose w such that the error function is minimized.

The Perceptron Error Function

Perceptron error \mathcal{E}_P is a function of the weights w

$$\operatorname{argmin}_{\mathbf{w}} \left(\mathcal{E}_P(\mathbf{w}) = - \sum_{m \in \mathcal{M}} \mathbf{w}^\top \mathbf{x}_m y_m \right) \quad (1)$$

where \mathcal{M} denotes the index set of all *misclassified* data \mathbf{x}_m
Data $\mathbf{x} \in \mathbb{R}^2$



The Perceptron Learning Algorithm

$$\text{Perceptron error } \mathcal{E}_P(\mathbf{w}) = -\sum_{m \in \mathcal{M}} \mathbf{w}^\top \mathbf{x}_m y_m$$

can be minimized *iteratively* using **stochastic gradient descent**
[Bottou, 2010; Robbins and Monro, 1951]

1. Initialize \mathbf{w}^{old} (randomly, $1/n$, ...)
2. While there are misclassified data points

Pick a random misclassified data point \mathbf{x}_m

Descent in direction of the gradient at single data point \mathbf{x}_m

$$\mathcal{E}_m(\mathbf{w}) = -\mathbf{w}^\top \mathbf{x}_m y_m$$

$$\nabla \mathcal{E}_m(\mathbf{w}) = -\mathbf{x}_m y_m$$

$$\mathbf{w}^{\text{new}} \leftarrow \mathbf{w}^{\text{old}} - \eta \nabla \mathcal{E}_m(\mathbf{w}^{\text{old}}) = \mathbf{w}^{\text{old}} + \eta \mathbf{x}_m y_m$$

We get

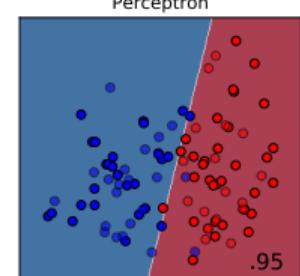
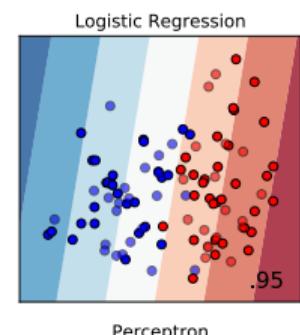
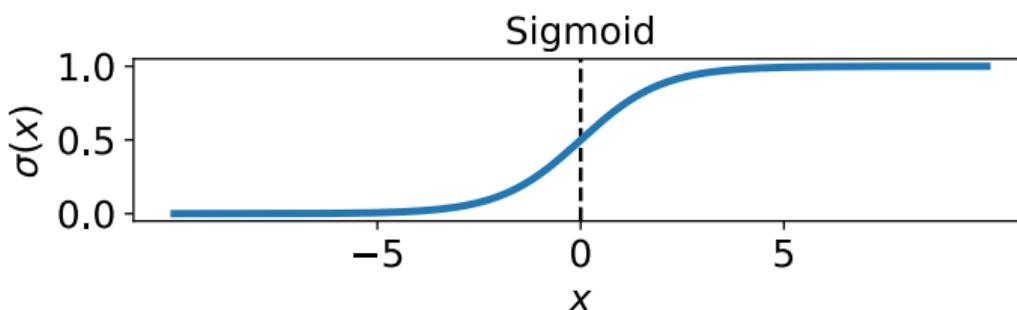
Logistic Regression

If we replace $f(\cdot)$ with the logistic (a.k.a. sigmoid) function

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

Interpret $\sigma(x)$ as probability that x in class +1

Trained with Gradient Descent



Perceptron (Recap)
oooooo

Multi-Layer Perceptron
●oooooo
ooo

Training
oooooooo

Practical considerations
oooo

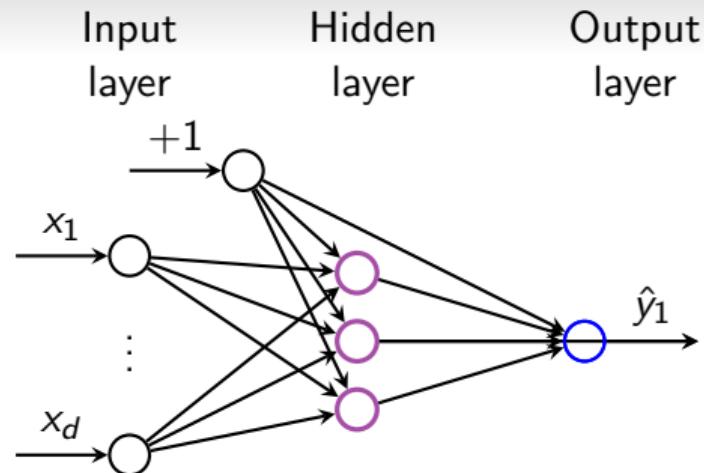
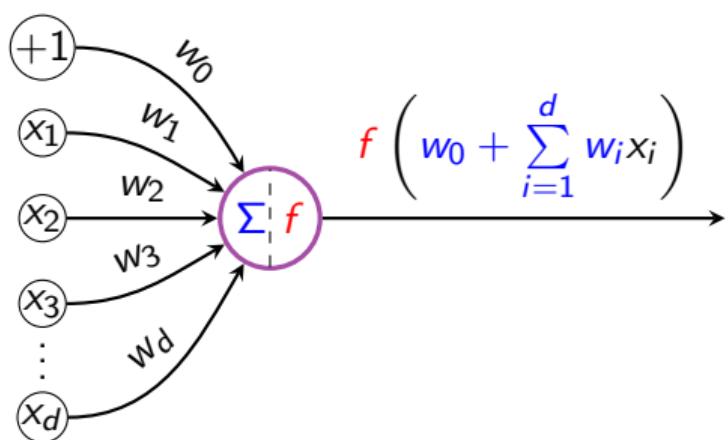
Different tastes of NNs
ooooo

Summary
ooo

What is a Multi-Layer Perceptron

And what is it capable of?

Perceptron



Adding more perceptrons in parallel yields
a multi-layer perceptron (MLP) with a single hidden layer

$$\hat{y}_1(x) = \sum_j w_j^o f(w_j^h{}^T x + w_{j,0}^h) + w_0^o$$

Multi-Layer Perceptron

Output can be (depending on):

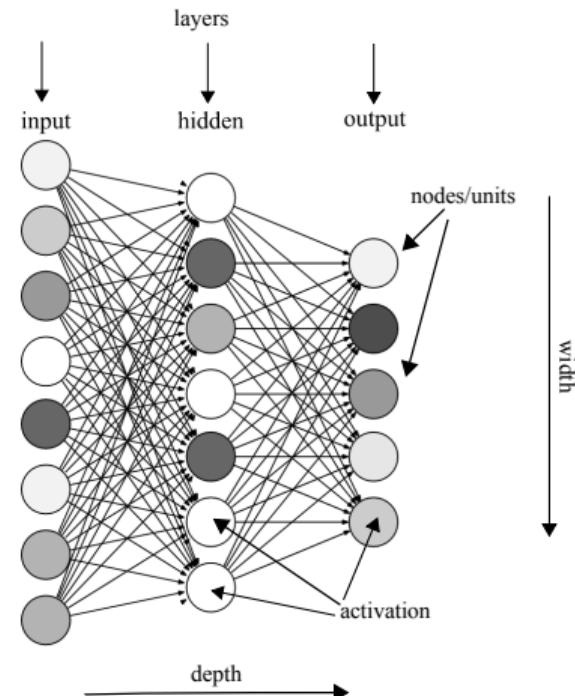
- Class or Regression (activation function of output node)
- Multidimensional (number of output nodes)

Let's see what that does:

<https://playground.tensorflow.org/>

(Homework) How does the network perform if we:

- add more layers (deeper net)
- add more nodes to layers (wider net)
- change the activation ($f(\cdot)$)



Universal approximation theorem

One can show that a single hidden layer **MLP** is a universal function approximator, meaning it **can model any** suitably smooth **function, given enough hidden units**, to any desired level of accuracy [Hornik, 1991].

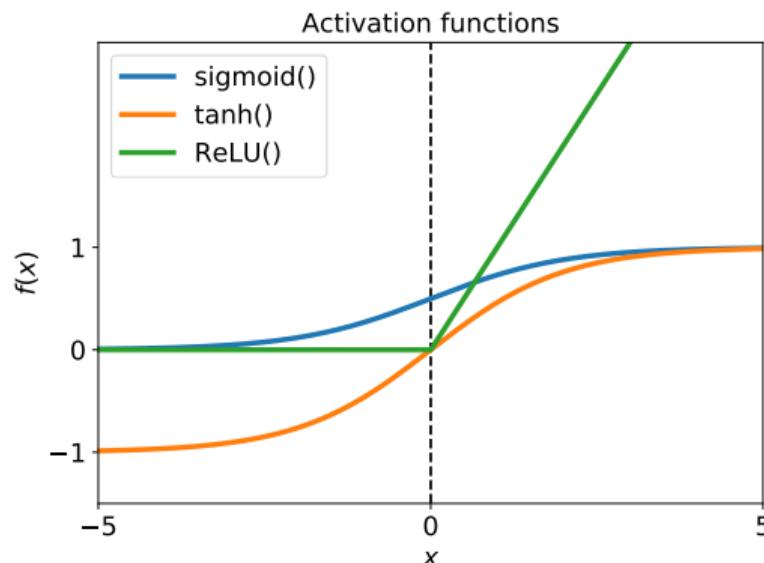
We only need a nonpolynomial activation function [Leshno et al., 1993]

Activation Functions

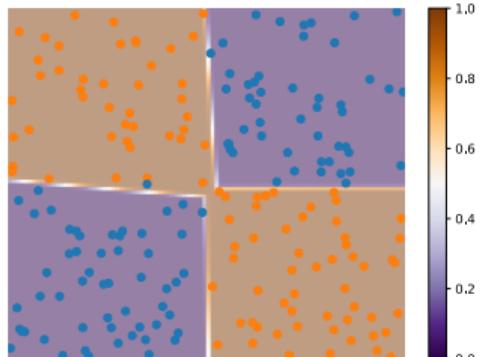
We can use different activation functions. In order for the MLP to be powerful we need **non-linear functions**

Popular activation functions are:

- Sigmoid $\sigma(x) = (1 + \exp(-x))^{-1}$
- $\tanh(x)$ ($= 2 \cdot \sigma(x) - 1$)
- **Rectified Linear Unit**
 $f(x) = \max(0, x)$

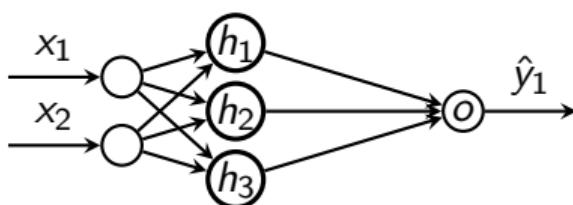
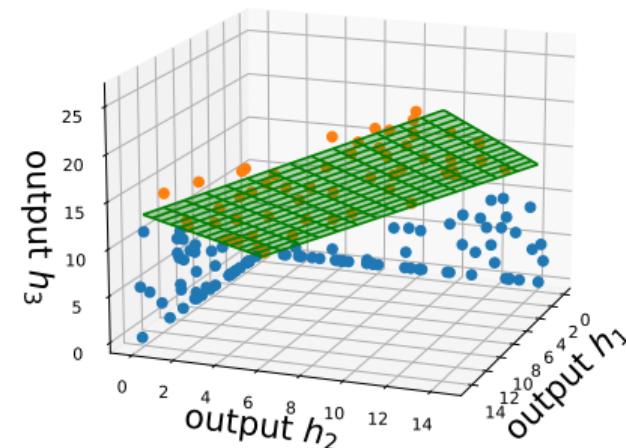


Single hidden layer MLP for XOR



Representation of data points at hidden layer:

- Decision boundary output

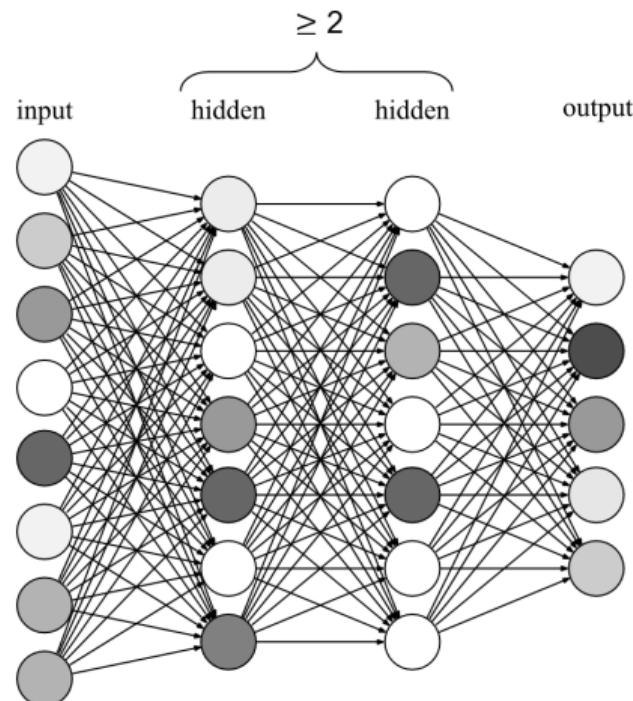


Input layer ReLU activation Sigmoid activation

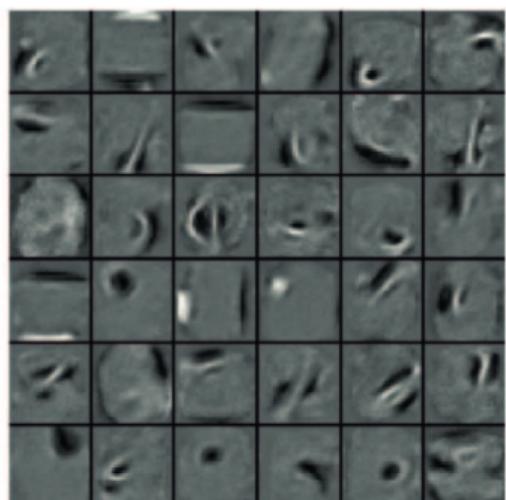
→ hidden layer learns representation in feature space

Deep Learning

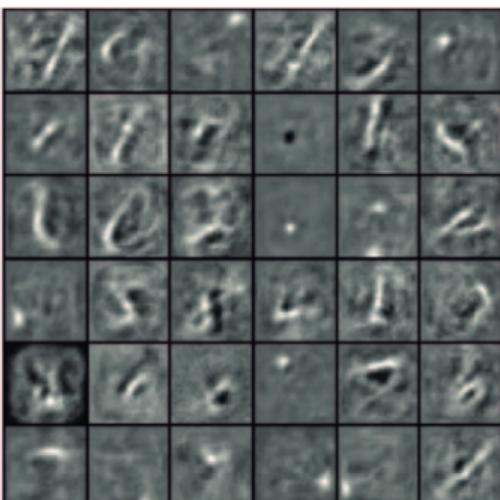
We talk about Deep Neural Networks (DNN) when there are 2 or more hidden layers



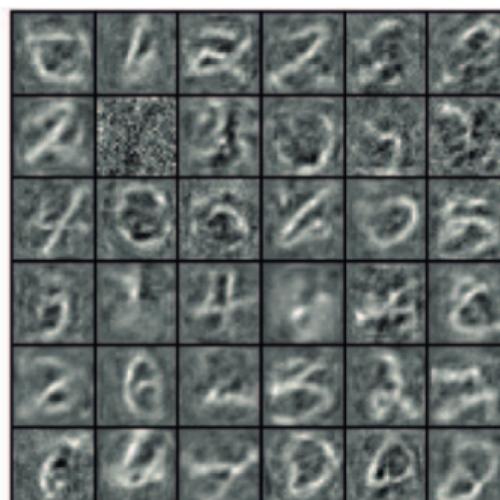
What do the hidden layers learn?



1st layer



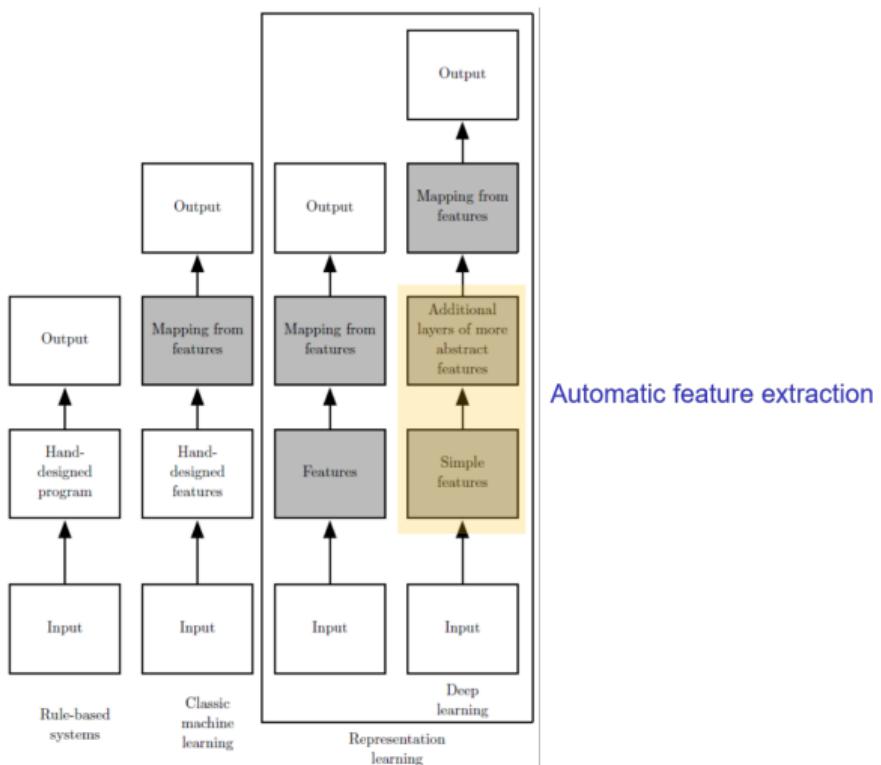
2nd layer



3rd layer

Erhan et al. [2010]

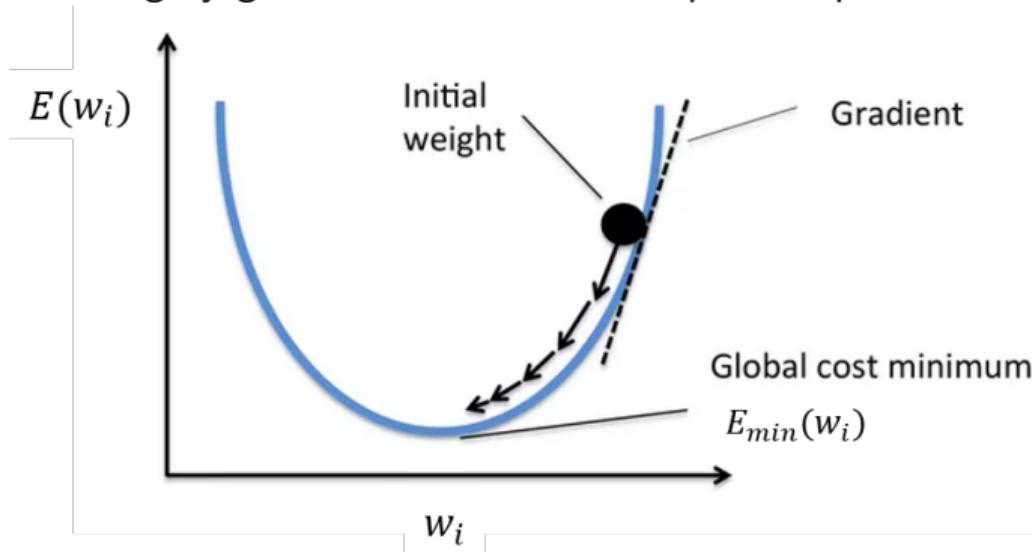
Deep Learning



How do we train the weights of a neural network?

Gradient Descent

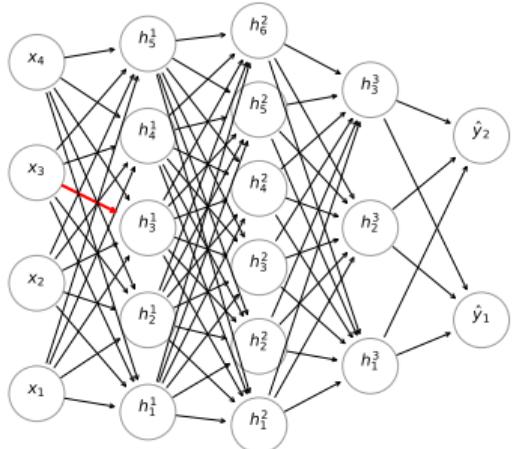
Learning by gradient descent in the space of parameters:



$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial E}{\partial w_{ij}}$$

where E is some error function and η is the learning rate.

How to get gradient for every **weight**

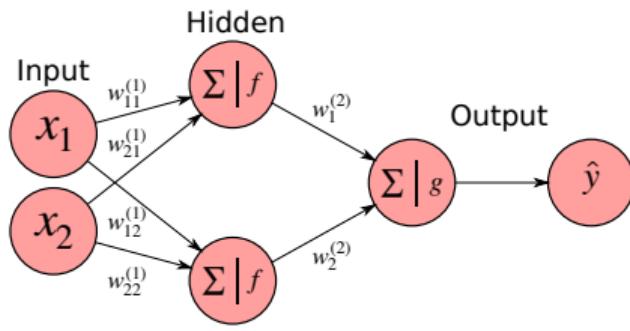


we use the
Chain Rule:

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

$(f \circ g)(x)$ is equivalent to writing $f(g(x))$

Applying the chain rule



$$f(x) = \sigma(x) = \frac{\exp(x)}{1 + \exp(x)}$$

$$g(x) = x$$

$$E = \frac{1}{2}(y - \hat{y})^2$$

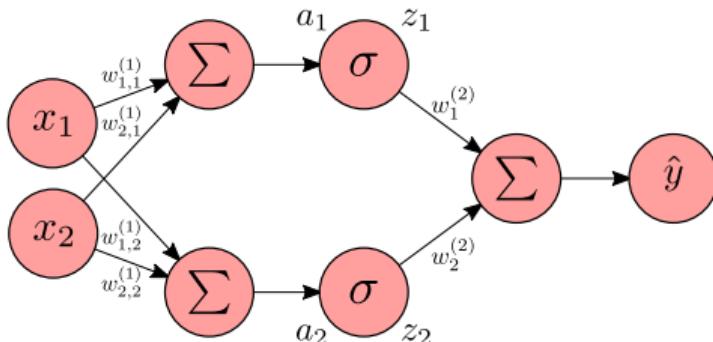
The MLP calculates:

$$\hat{y} = g \left(\sum_i w_i^{(2)} f \left(\sum_j w_{j,i}^{(1)} x_j \right) \right) = g \left(\mathbf{w}^{(2)\top} f \left(\mathbf{W}^{(1)\top} \mathbf{X} \right) \right)$$

in matrix notation:

Applying the chain rule

Input Hidden



$$f(x) = \sigma(x) = \frac{\exp(x)}{1 + \exp(x)}$$

$$\sigma(x)' = \sigma(x)(1 - \sigma(x))$$

$$E = \frac{1}{2}(\hat{y} - y)^2$$

$$\hat{y} = w_1^{(2)} \underbrace{\sigma(w_{11}^{(1)}x_1 + w_{21}^{(1)}x_2)}_{z_1} + w_2^{(2)} \underbrace{\sigma(w_{12}^{(1)}x_1 + w_{22}^{(1)}x_2)}_{z_2}$$

$$\frac{\partial E}{\partial w_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1} \frac{\partial z_1}{\partial w_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1} \frac{\partial z_1}{\partial a_1} \frac{\partial a_1}{\partial w_{11}^{(1)}}$$

Applying the chain rule

Remember:

$$\sigma(x)' = \sigma(x)(1 - \sigma(x)), \quad E = \frac{1}{2}(\hat{y} - y)^2$$

$$\frac{\partial E}{\partial \hat{y}} = (\hat{y} - y)$$

$$\hat{y} = w_1^{(2)} \sigma(\underbrace{w_{11}^{(1)} x_1 + w_{21}^{(1)} x_2}_{a_1}) + w_2^{(2)} \sigma(\underbrace{w_{12}^{(1)} x_1 + w_{22}^{(1)} x_2}_{z_2})$$

$$\frac{\partial \hat{y}}{\partial z_1} = w_1^{(2)}$$

$$\frac{\partial z_1}{\partial a_1} = \sigma(a_1)(1 - \sigma(a_1))$$

$$\frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1} \frac{\partial z_1}{\partial a_1} \frac{\partial a_1}{\partial w_{11}^{(1)}}$$

$$\frac{\partial a_1}{\partial w_{11}^{(1)}} = x_1$$

$$\Rightarrow \frac{\partial E}{\partial w_{11}^{(1)}} = (\hat{y} - y) \cdot w_1^{(2)} \cdot \sigma(a_1)(1 - \sigma(a_1)) \cdot x_1$$

Backpropagation

Algorithm that applies chain rule with the help of dynamic programming.

1. Take datapoints X

2. Predict label \hat{y} (forward)

Save activation value a for every node

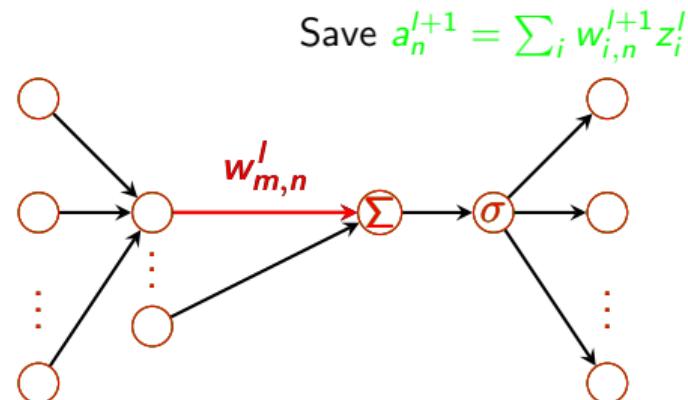
3. Backpropagate error $E(\hat{y}, y)$

Calculate $\frac{\partial E}{\partial z^L}$ ($= \frac{\partial E}{\partial \hat{y}}$)

go backwards to calculate gradients

$$\frac{\partial E}{\partial w_{m,n}^l} = \frac{\partial E}{\partial a_n^{l+1}} \frac{\partial a_n^{l+1}}{\partial w_{m,n}^l}$$

4. Update weights with $w \leftarrow w - \eta \frac{\partial E}{\partial w}$



$$\frac{\partial E}{\partial z_n^l} = \sum_p \frac{\partial E}{\partial a_p^{l+1}} \frac{\partial a_p^{l+1}}{\partial z_n^l}$$

$$\frac{\partial E}{\partial a_n^l} = \frac{\partial E}{\partial z_n^l} \frac{\partial z_n^l}{\partial a_n^l}$$

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}$$

$$(f \circ g): \mathbb{R}^n \rightarrow \mathbb{R}$$

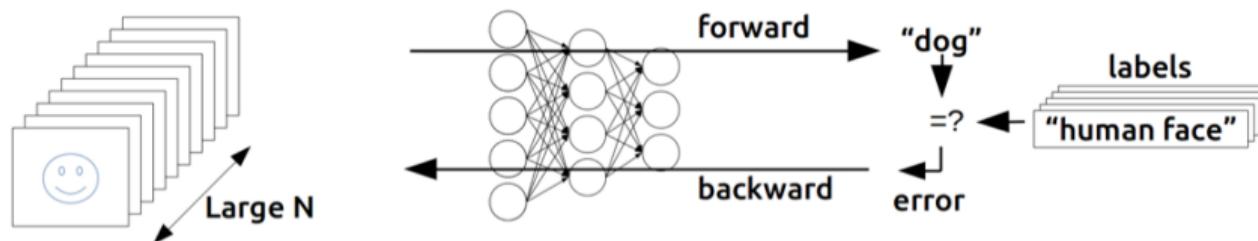
$$\frac{\partial}{\partial x_i} (f \circ g) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right) \begin{pmatrix} \frac{\partial g_1}{\partial x_i} \\ \vdots \\ \frac{\partial g_m}{\partial x_i} \end{pmatrix} = \sum_{j=1}^m \frac{\partial f}{\partial x_j} \cdot \frac{\partial g_j}{\partial x_i}$$

$$= (\nabla f) \cdot \frac{\partial x}{\partial x_i} \begin{pmatrix} g_1 \\ \vdots \\ g_m \end{pmatrix}$$

Summary: Backpropagation

In a neural network with hidden layers there are many more parameters as in the Perceptron. So an efficient way to compute the gradients is needed.

Training



<https://devblogs.nvidia.com>

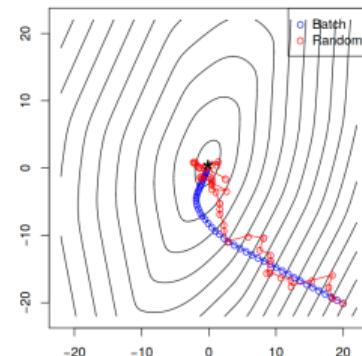
Repeat until convergence/stopping criteria:

1. Forward pass: Apply the current weights to your data
2. Backward pass: Error backpropagation for every individual w_{ij}
3. Update weights: $w_{ij} \leftarrow w_{ij} - \eta \frac{\partial E}{\partial w_{ij}}$

Gradient Descent (GD) vs. Stochastic Gradient Descent (SGD)

Updates using the entire data set are computationally expensive

- full-batch mode (GD): use all training samples before updating the weights ($\mathcal{O}(pn)$)
- mini-batch mode (SGD): use a subset of the training data ($\mathcal{O}(pm)$)

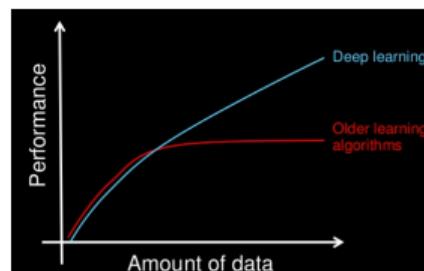


from Tibshirani's lecture

Algorithms for adaptive learning rate η can yield speed-ups and can guarantee convergence (e.g. Adam)

Practical considerations

- MLP is a standard feed forward (artificial) Neural Network
- Need a lot of data to work
- Best for high-dimensional datasets
- MLP's tend to overfit really easily
 - Dropout, Regularisation, early stopping, ...



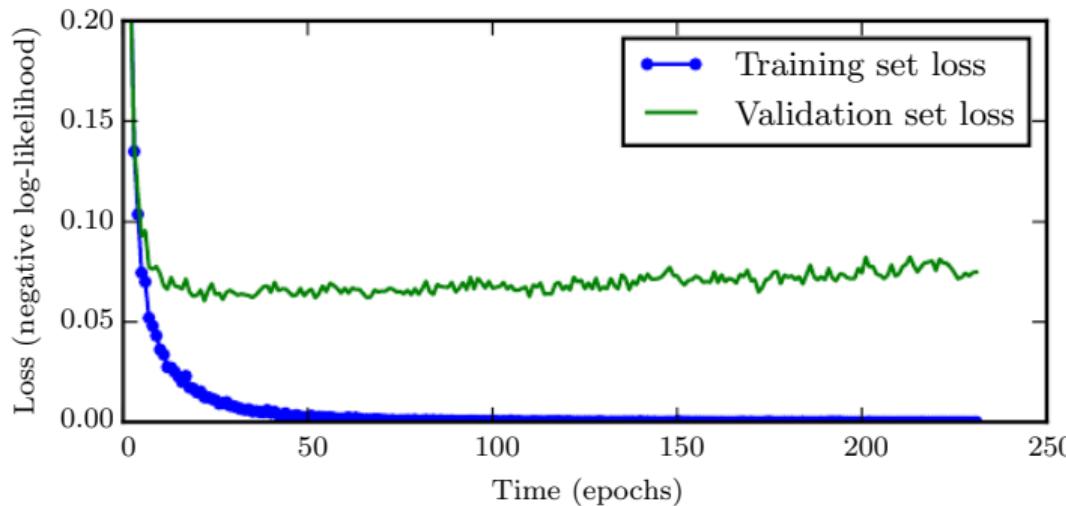
from Andrew Ng



Output of MLP with 1,2 and 4 hidden layers

from <https://www.kdnuggets.com>

How to prevent overfitting: Early Stopping

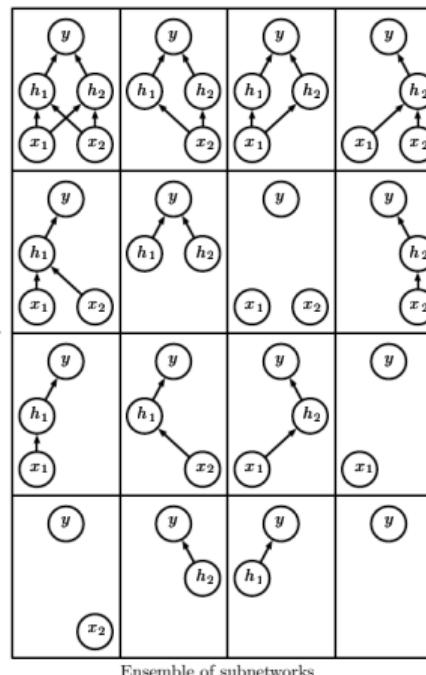
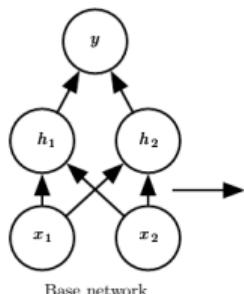


“The algorithm terminates when no parameters have improved over the best recorded validation error for some pre-specified number of iterations” [Goodfellow et al., 2016]

当没有参数比最佳记录值有所改善时，该算法终止 图 7.3.. 学习曲线显示了在一些预先指定的迭代次数下，负对数似然损失如何随验证误差而变化

How to prevent overfitting: Dropout

- Randomly set weights to 0 temporarily during training
- Makes the Network more robust



[Goodfellow et al., 2016]

When training a neural network from scratch there are many parameters you need to choose!

This is mostly done by **trial and error**.

- Which regularisation
- Which architecture
- Which activation function
- Which learning rate
- ...

You can also use grid search for some parameters.

Perceptron (Recap)
oooooo

Multi-Layer Perceptron
oooooooo
ooo

Training
oooooooooooo

Practical considerations
oooo

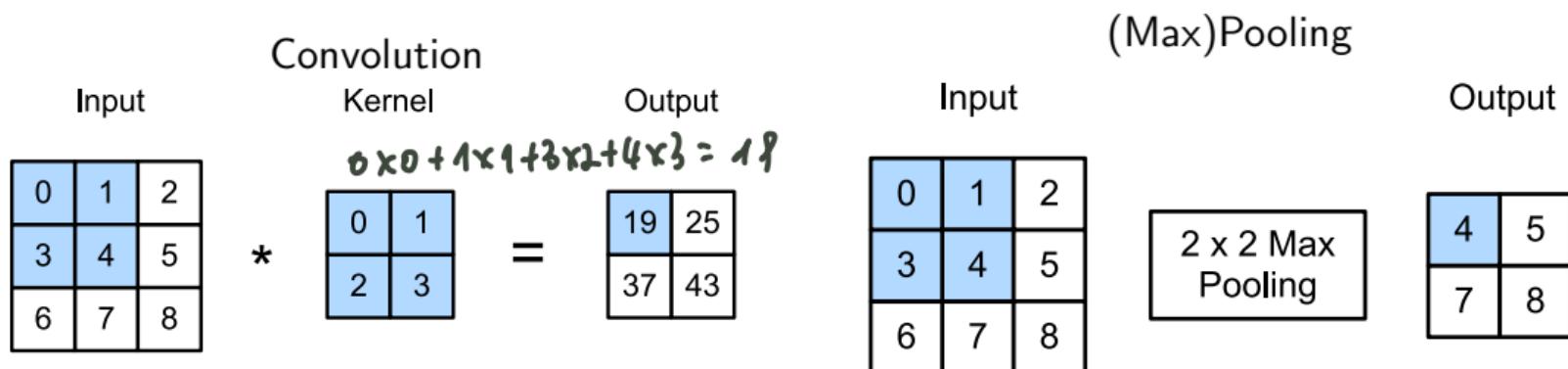
Different tastes of NNs
●oooo

Summary
ooo

Let's see some more advanced examples for Neural Networks

Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are widely used in image processing such as object recognition. They usually consist of various convolutional and pooling layers:



[Zhang et al., 2019]

Other architectures: Convolutional Neural Networks

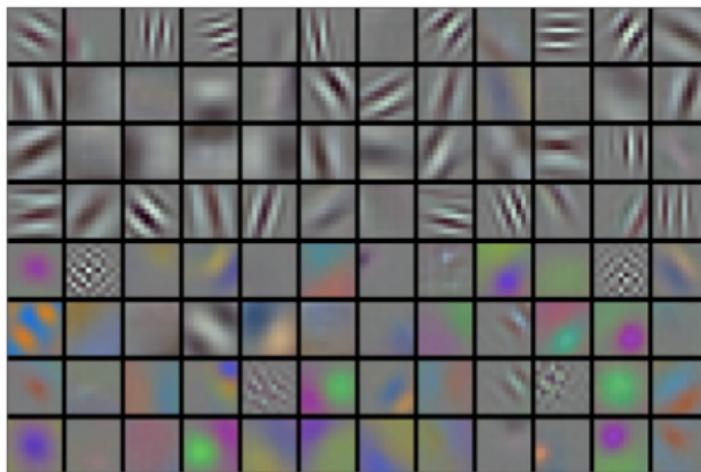


Fig. 9.1.1: Image filters learned by the first layer of AlexNet

Zhang et al. [2019]

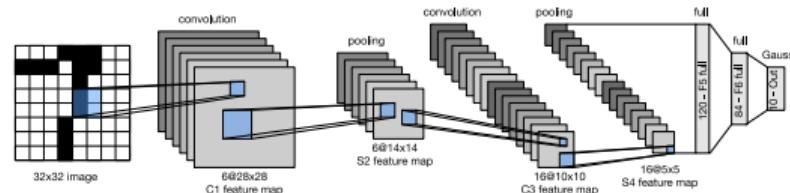


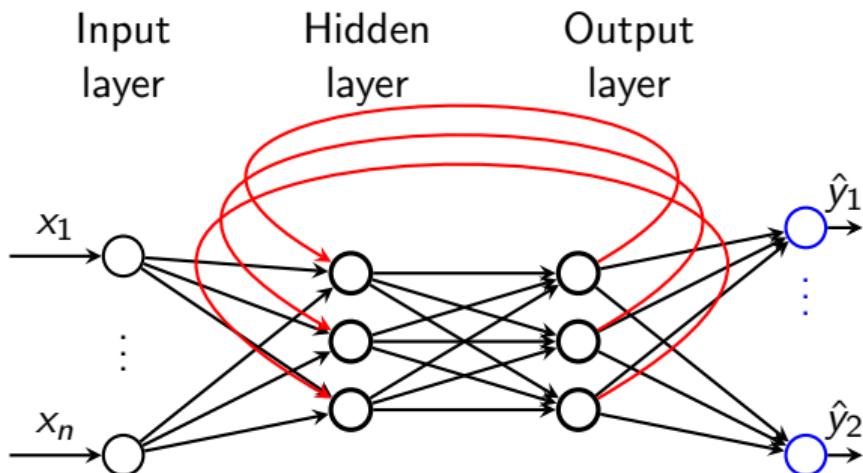
Fig. 8.6.1: Data flow in LeNet 5. The input is a handwritten digit, the output a probability over 10 possible outcomes.

Zhang et al. [2019]

Recurrent Neural Network

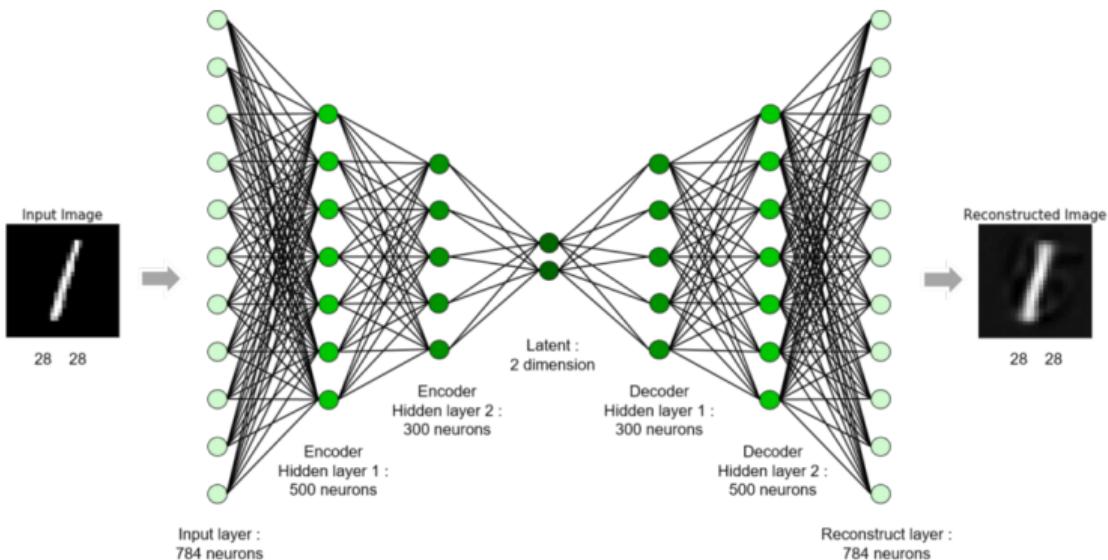
循环神经网络

- Used for time-series data
- Can use recurrence for representing states



Auto-encoders

- Used for unsupervised learning
- Finds low-dimensional representation of data



from mc.ai

Summary I

- Multilayer Perceptrons/Neural Networks are very powerful
- more complicated to train than linear models
 - Non-convex problem
 - Stochastic Gradient Descent
- many choices need to be made
- training can be costly
 - Since we have many parameters and big datasets
- in times of Deep Learning things are changing fast
- trial and error!

Summary II

Open source books and lectures for further reading (and watching):

- Deep Learning by Goodfellow, Bengio & Courville (2016)
- Dive into Deep Learning by Zhang, Lipton, Li & Smola (2019)
with many coding examples
- 3Blue1Brown videos on MLPs and Backpropagation Different Notation!
- Video Lectures by Geoff Hinton

References

- L. Bottou. Large-scale machine learning with stochastic gradient descent. In Y. Lechevallier and G. Saporta, editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010)*, pages 177–187, Paris, France, 2010. Springer.
- D. Erhan, A. Courville, and Y. Bengio. Understanding representations learned in deep architectures. *Department dInformatique et Recherche Operationnelle, University of Montreal, QC, Canada, Tech. Rep.* 1355:1, 2010.
- I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- M. Leshno, V. Y. Lin, A. Pinkus, and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6):861–867, 1993.
- H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3):400—407, 1951.
- A. Zhang, Z. C. Lipton, M. Li, and A. J. Smola. *Dive into Deep Learning*. 2019. <http://www.d2l.ai>.