



# Statistical Machine Learning

## *Lecture 9*

### *Convolutional neural networks*

### *How to train neural networks*



UPPSALA  
UNIVERSITET

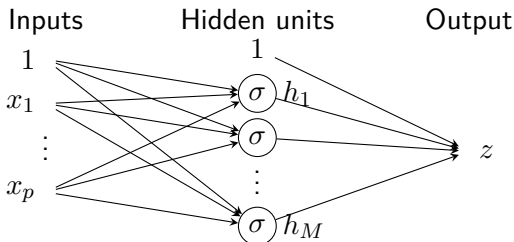
**Niklas Wahlström**

Division of Systems and Control  
Department of Information Technology  
Uppsala University

[niklas.wahlstrom@it.uu.se](mailto:niklas.wahlstrom@it.uu.se)  
[www.it.uu.se/katalog/nikwa778](http://www.it.uu.se/katalog/nikwa778)

# Summary of Lecture 8 (I/III) - Neural network

A neural network is a sequential construction of **several** generalized linear regression models.



$$h_1 = \sigma \left( \beta_{01}^{(1)} + \sum_{j=1}^p \beta_{j1}^{(1)} x_j \right)$$

$$h_2 = \sigma \left( \beta_{02}^{(1)} + \sum_{j=1}^p \beta_{j2}^{(1)} x_j \right)$$

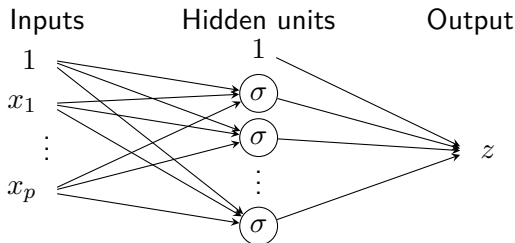
$$\vdots$$

$$h_M = \sigma \left( \beta_{0M}^{(1)} + \sum_{j=1}^p \beta_{jM}^{(1)} x_j \right)$$

$$z = \beta_0^{(2)} + \sum_{m=1}^M \beta_m^{(2)} h_m$$

# Summary of Lecture 8 (I/III) - Neural network

A neural network is a sequential construction of **several** generalized linear regression models.



$$\mathbf{h} = \sigma(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)\top})$$

$$z = \mathbf{W}^{(2)\top} \mathbf{h} + \mathbf{b}^{(2)\top}$$

$$\mathbf{h} = \begin{bmatrix} h_1 \\ \vdots \\ h_M \end{bmatrix}$$

Hidden  
units

$$\mathbf{b}^{(1)} = [\beta_{01}^{(1)} \dots \beta_{0M}^{(1)}]$$

$$\mathbf{W}^{(1)} = \begin{bmatrix} \beta_{11}^{(1)} & \dots & \beta_{1M}^{(1)} \\ \vdots & \dots & \vdots \\ \beta_{p1}^{(1)} & \dots & \beta_{pM}^{(1)} \end{bmatrix}$$

offset vector

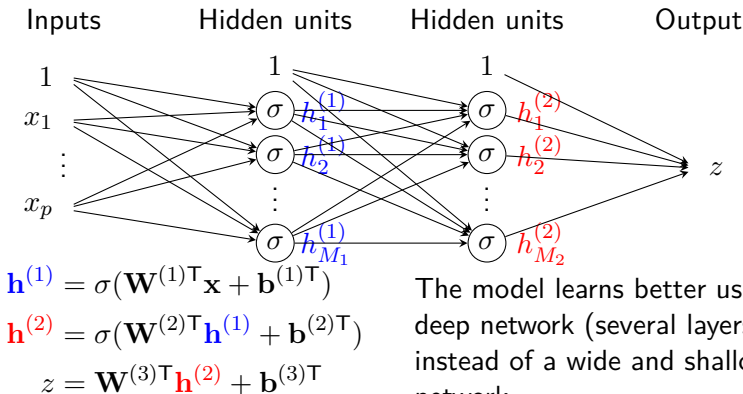
$$\mathbf{b}^{(2)} = [\beta_0^{(2)}]$$

weight matrix

$$\mathbf{W}^{(2)} = \begin{bmatrix} \beta_1^{(2)} \\ \vdots \\ \beta_M^{(2)} \end{bmatrix}$$

# Summary of Lecture 8 (I/III) - Neural network

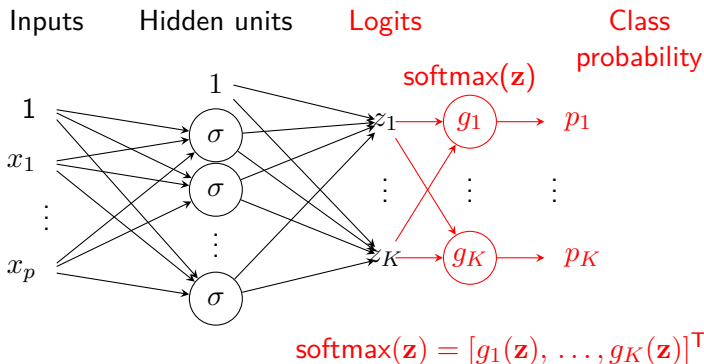
A neural network is a **sequential** construction of several generalized linear regression models.



# Summary of Lecture 8 (II/III) - classification

For  $K > 2$  classes we want to predict the class probability for all  $K$  classes  $p_k = p(y = k|\mathbf{x})$ . We extend the logistic function to the **softmax activation function**

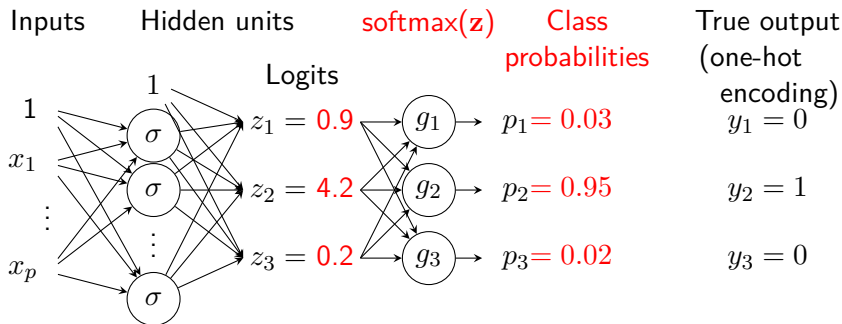
$$p_k = g_k(z) = \frac{e^{z_k}}{\sum_{l=1}^K e^{z_l}}, \quad k = 1, \dots, K.$$



# Summary of Lecture 8 (III/III)

## - Example $K = 3$ classes

Consider an example with three classes  $K = 3$ .



The network is trained by minimizing the **cross-entropy**

$$L(\mathbf{y}, \mathbf{p}) = - \sum_{k=1}^K y_k \log(p_k) = - \log 0.95 = 0.05$$

Cross-entropy is the ML solution. See lecture notes, Sec. 3.2.3.

# Guest lectures and content

---

## Guest lectures

1. **Tuesday March 5 at 15.15-16.00, Siegbahnsalen**  
*Peltarion* - deep learning for audio applications
2. **Monday March 11 at 14.15-15.00, Siegbahnsalen**  
*Spotify* - their use of machine learning

---

## Outline

1. **Previous lecture** The neural network model
  - Neural network for regression
  - Neural network for classification
  - Convolutional neural network (we just started)
2. **This lecture**
  - Convolutional neural network
  - How to train a neural network



# Convolutional neural networks

---

**Convolutional neural networks** (CNN) are a special kind neural networks tailored for problems where the input data has a grid-like structure.

## Examples

- Digital images (2D grid of pixels)
- Audio waveform data (1D grid, times series)
- Volumetric data e.g. CT scans (3D grid)

The description here will focus on images.

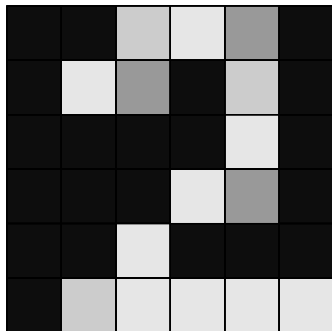


# Data representation of images

Consider a grayscale image of  $6 \times 6$  **pixels**.

- Each pixel value represents the color. The value ranges from 0 (total absence, black) to 1 (total presence, white)
- The pixels are the input variables  $x_{1,1}, x_{1,2}, \dots, x_{6,6}$ .

Image



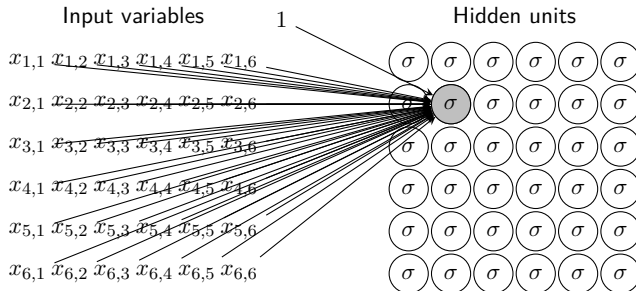
Data representation

0.0	0.0	0.8	0.9	0.6	0.0
0.0	0.9	0.6	0.0	0.8	0.0
0.0	0.0	0.0	0.0	0.9	0.0
0.0	0.0	0.0	0.9	0.6	0.0
0.0	0.0	0.9	0.0	0.0	0.0
0.0	0.8	0.9	0.9	0.9	0.9

# The convolutional layer

Consider a hidden layer with  $6 \times 6$  hidden units.

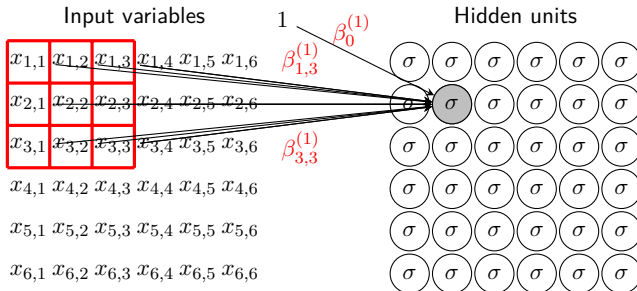
- **Dense layer:** Each hidden unit is connected with **all pixels**.  
Each pixel-hidden-unit-pair has its own **unique parameter**.



# The convolutional layer

Consider a hidden layer with  $6 \times 6$  hidden units.

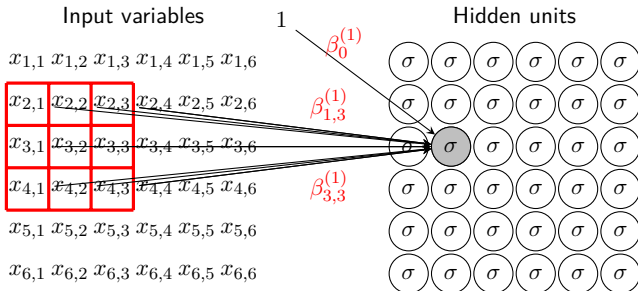
- **Dense layer:** Each hidden unit is connected with **all pixels**. Each pixel-hidden-unit-pair has its own **unique parameter**.
- **Convolutional layer:** Each hidden unit is connected with a **region of pixels** via a set of parameters, so-called **kernel**. Different hidden units have the **same set of parameters**.



# The convolutional layer

Consider a hidden layer with  $6 \times 6$  hidden units.

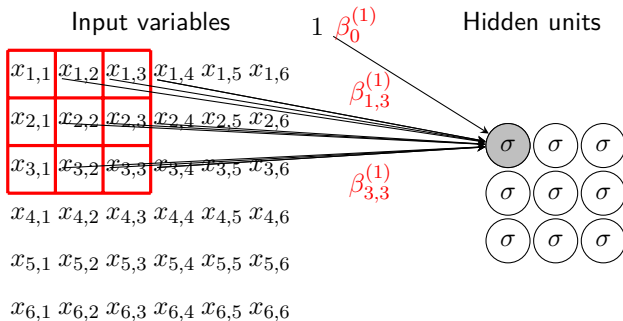
- **Dense layer:** Each hidden unit is connected with **all pixels**. Each pixel-hidden-unit-pair has its own **unique parameter**.
- **Convolutional layer:** Each hidden unit is connected with a **region of pixels** via a set of parameters, so-called **kernel**. Different hidden units have the **same set of parameters**.



Conv. layer uses **sparse interactions** and **parameter sharing**

# Condensing information with strides

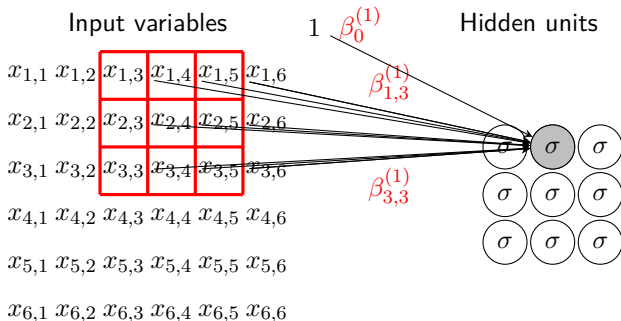
- **Problem:** As we proceed through the network we want to condense the information.
- **Solution:** Apply the kernel to every two pixels. We use a **stride** of 2 (instead of 1).



With stride 2 we get half the number of rows and columns in the hidden layer.

# Condensing information with strides

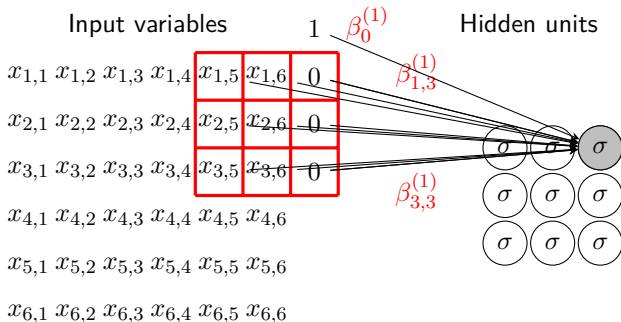
- **Problem:** As we proceed through the network we want to condense the information.
- **Solution:** Apply the kernel to every two pixels. We use a **stride** of 2 (instead of 1).



With stride 2 we get half the number of rows and columns in the hidden layer.

# Condensing information with strides

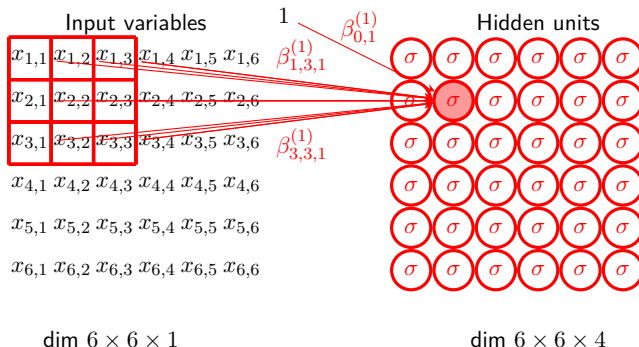
- **Problem:** As we proceed through the network we want to condense the information.
- **Solution:** Apply the kernel to every two pixels. We use a **stride** of 2 (instead of 1).



With stride 2 we get half the number of rows and columns in the hidden layer.

# Multiple channels

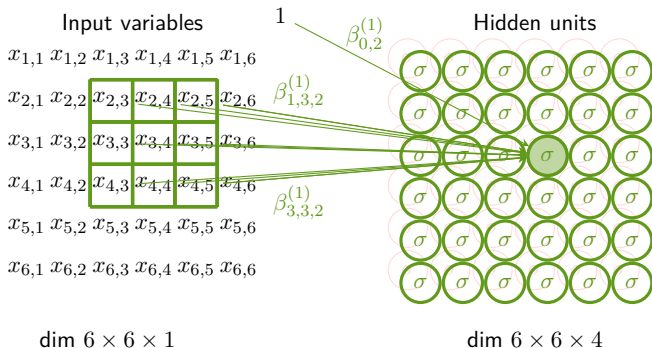
- One kernel per layer does not give enough flexibility.  $\Rightarrow$
- We use **multiple kernels** (visualized with different colors).
- Each kernel produces its own set of hidden units – a **channel**.





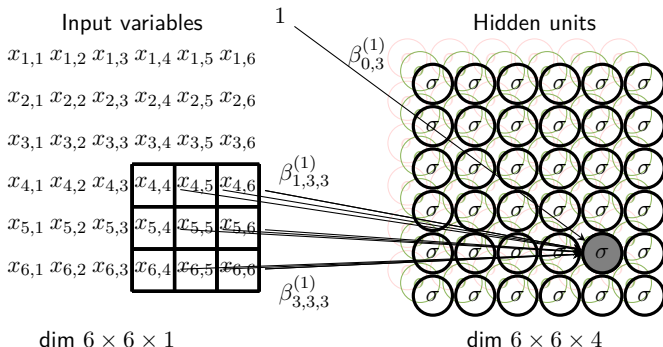
# Multiple channels

- One kernel per layer does not give enough flexibility.  $\Rightarrow$
- We use **multiple kernels** (visualized with different colors).
- Each kernel produces its own set of hidden units – a **channel**.



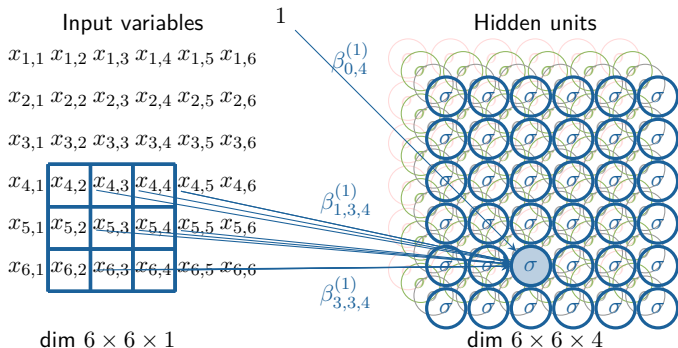
# Multiple channels

- One kernel per layer does not give enough flexibility.  $\Rightarrow$
- We use **multiple kernels** (visualized with different colors).
- Each kernel produces its own set of hidden units – a **channel**.



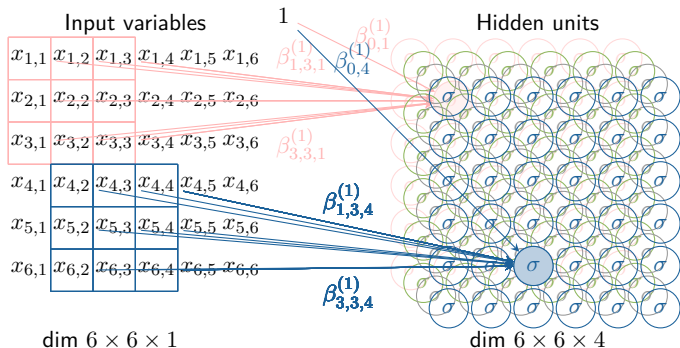
# Multiple channels

- One kernel per layer does not give enough flexibility.  $\Rightarrow$
- We use **multiple kernels** (visualized with different colors).
- Each kernel produces its own set of hidden units – a **channel**.



# Multiple channels

- One kernel per layer does not give enough flexibility.  $\Rightarrow$
- We use **multiple kernels** (visualized with different colors).
- Each kernel produces its own set of hidden units – a **channel**.



Hidden layers are organized in **tensors** of size (rows  $\times$  columns  $\times$  channels).

# What is a tensor?

A **tensor** is a generalization of scalar, vector and matrix to arbitrary **order**.

## Scalar

order 0

$$a = 3$$



## Vector

order 1

$$\mathbf{b} = \begin{bmatrix} 3 \\ -2 \\ -1 \end{bmatrix}$$



## Matrix

order 2

$$\mathbf{W} = \begin{bmatrix} 3 & 2 \\ -2 & 1 \\ -1 & 2 \end{bmatrix}$$



## Tensor

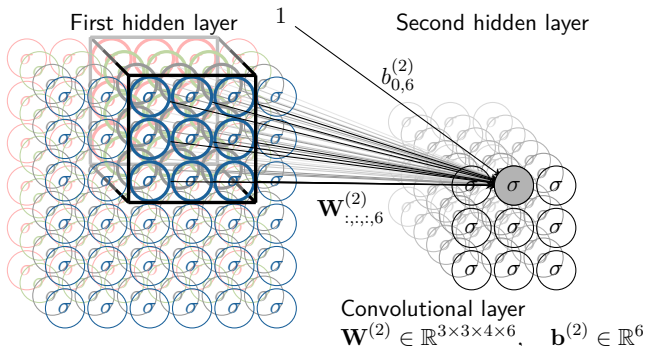
any order  
(here order 3)

$$\mathbf{T}_{::,1} = \begin{bmatrix} 3 & 2 \\ -2 & 1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{T}_{::,2} = \begin{bmatrix} -1 & 4 \\ 1 & 2 \\ -5 & 3 \end{bmatrix}$$



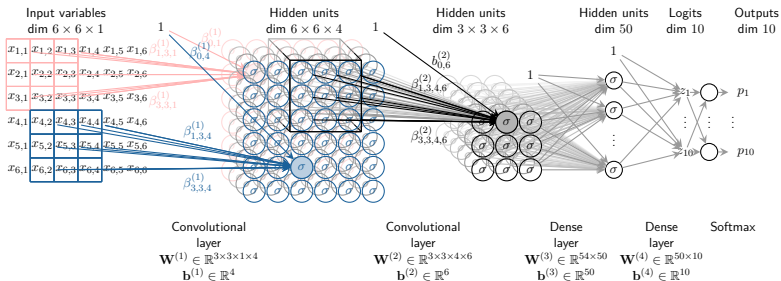
# Multiple channels (cont.)

- A kernel operates on **all channels** in a hidden layer.
- Each kernel has the dimension (kernel rows  $\times$  kernel columns  $\times$  input channels), here  $(3 \times 3 \times 4)$ .
- We stack all kernel parameters in a **weight tensor** with dimensions (kernel rows  $\times$  kernel columns  $\times$  input channels  $\times$  output channels), here  $(3 \times 3 \times 4 \times 6)$



# Full CNN architecture

- A full CNN usually consist of multiple convolutional layers (here two) and a few final dense layers (here two).
- If we have a classification problem at hand, we end with a softmax activation function to produce class probabilities.



Here we use 50 hidden units in the last hidden layer and consider a classification problem with  $K = 10$  classes.

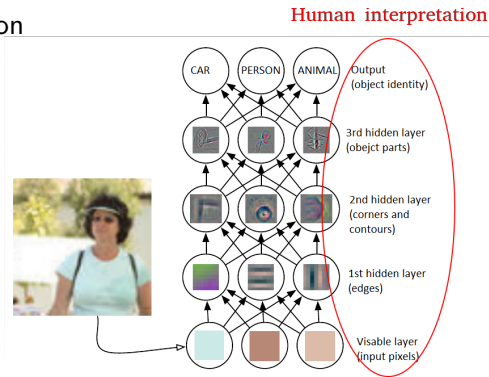
# Why deep?

## Example: Image classification

**Input:** pixels of an **image**

**Output:** **object identity**

Each hidden layer extracts increasingly abstract features.



Zeiler, M. D. and Fergus, R. **Visualizing and understanding convolutional networks**  
*Computer Vision - ECCV (2014)*.



# Skin cancer – background

---

One recent result on the use of deep learning in medicine -  
Detecting skin cancer (February 2017)

Andre Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M. and Thrun, S. **Dermatologist-level classification of skin cancer with deep neural networks.** *Nature*, 542, 115–118, February, 2017.

Some background figures (from the US) on skin cancer:

- Melanomas represents less than 5% of all skin cancers, **but** accounts for 75% of all skin-cancer-related deaths.
- Early detection absolutely critical. Estimated 5-year survival rate for melanoma: Over 99% if detected in its earlier stages and 14% if detected in its later stages.



# Skin cancer – task

---

Image copyright Nature (doi:10.1038/nature21056)

# Skin cancer – solution (ultrabrief)

---

In the paper they used the following network architecture

Image copyright Nature doi:10.1038/nature21056)

- Initialize all parameters from a neural network trained on 1.28 million images (**transfer learning**).
- From this initialization we learn new model parameters using 129 450 clinical images ( $\sim 100$  times more images than any previous study).
- Use the model to predict class based on unseen data.

# Skin cancer – indication of the results

---

$$\text{sensitivity} = \frac{\text{true positive}}{\text{positive}}$$

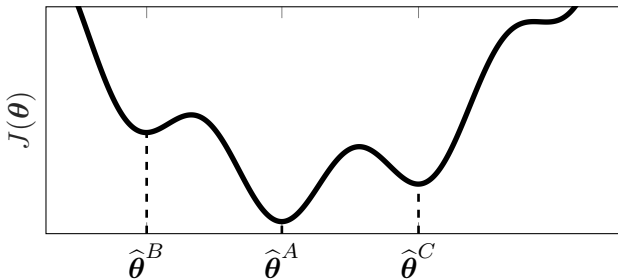
$$\text{specificity} = \frac{\text{true negative}}{\text{negative}}$$

Image copyright Nature (doi:10.1038/nature21056)

# Unconstrained numerical optimization

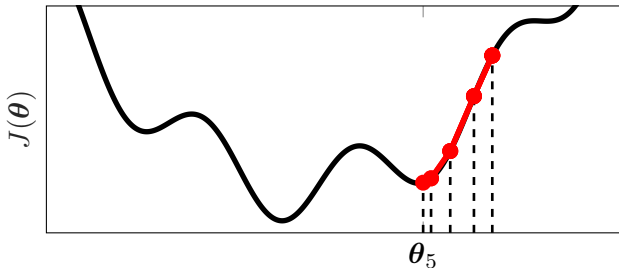
We train a network by considering the optimization problem

$$\hat{\theta} = \arg \min_{\theta} J(\theta), \quad J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i, \theta)$$



- The best possible solution  $\hat{\theta}$  is the **global minimizer** ( $\theta^A$ )
- The global minimizer is typically very hard to find, and we have to settle for a **local minimizer** ( $\theta^A, \theta^B, \theta^C$ )

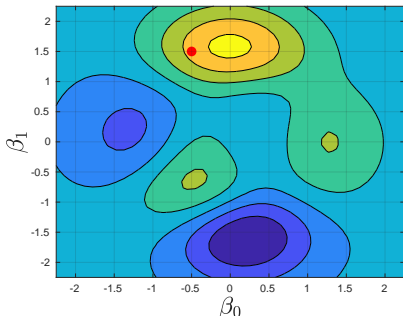
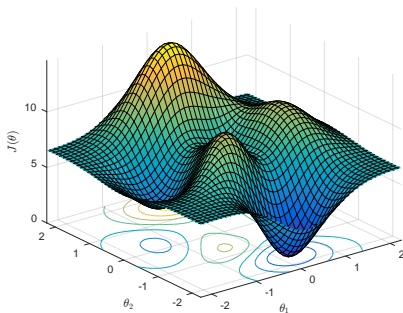
# Iterative solution - Example 1D



In our search for a local optimizer we...

- ... do an initial guess of  $\theta$ ...
- ... and update  $\theta$  iteratively.

# Iterative solution (gradient descent) - Example 2D

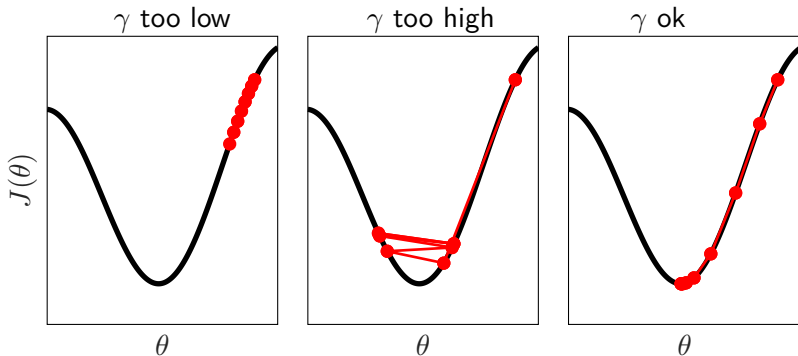


$$[\beta_0, \beta_1]^T \in \mathbb{R}^2$$

1. Pick a  $\theta_0$
2. while(*not converged*)
  - Update  $\theta_{t+1} = \theta_t - \gamma \mathbf{g}_t$ , where  $\mathbf{g}_t = \nabla_{\theta} J(\theta)$
  - Update  $t := t + 1$

We call  $\gamma \in \mathbb{R}$  the **step length** or **learning rate**.

# Learning rate



**Tuning strategy:** If the cost function...

- ...decreases very slowly  $\Rightarrow$  increase the learning rate.
- ...oscillates widely  $\Rightarrow$  reduce the learning rate.



# Computational challenge 1 - $\dim(\theta)$ is big

---

At each optimization step we need to compute the gradient

$$\mathbf{g}_t = \nabla_{\theta} J(\theta_t) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta_t).$$

**Computational challenge 1 -  $\dim(\theta)$  big:** A neural network contains a lot of parameters. Computing the gradient is costly.

**Solution:** A NN is a composition of multiple layers. Hence, each term  $\nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta)$  can be computed efficiently by repeatedly applying the chain rule. This is called the **back-propagation algorithm**. Not part of the course.

## Computational challenge 2 - $n$ is big

---

At each optimization step we need to compute the gradient

$$\mathbf{g}_t = \nabla_{\theta} J(\theta_t) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta_t).$$

**Computational challenge 2 -  $n$  big:** We typically use a lot of training data  $n$  for training the neural network. Computing the gradient is costly.

**Solution:** For each iteration, we only use a small part of the data set to compute the gradient  $\mathbf{g}_t$ . This is called the **stochastic gradient descent**.

# Stochastic gradient descent

A big data set is often redundant = many data points are similar.

Training data

$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$	$\mathbf{x}_5$	$\mathbf{x}_6$	$\mathbf{x}_7$	$\mathbf{x}_8$	$\mathbf{x}_9$	$\mathbf{x}_{10}$	$\mathbf{x}_{11}$	$\mathbf{x}_{12}$	$\mathbf{x}_{13}$	$\mathbf{x}_{14}$	$\mathbf{x}_{15}$	$\mathbf{x}_{16}$	$\mathbf{x}_{17}$	$\mathbf{x}_{18}$	$\mathbf{x}_{19}$	$\mathbf{x}_{20}$
$\mathbf{y}_1$	$\mathbf{y}_2$	$\mathbf{y}_3$	$\mathbf{y}_4$	$\mathbf{y}_5$	$\mathbf{y}_6$	$\mathbf{y}_7$	$\mathbf{y}_8$	$\mathbf{y}_9$	$\mathbf{y}_{10}$	$\mathbf{y}_{11}$	$\mathbf{y}_{12}$	$\mathbf{y}_{13}$	$\mathbf{y}_{14}$	$\mathbf{y}_{15}$	$\mathbf{y}_{16}$	$\mathbf{y}_{17}$	$\mathbf{y}_{18}$	$\mathbf{y}_{19}$	$\mathbf{y}_{20}$

If the training data is big

$$\nabla_{\theta} J(\theta) \approx \sum_{i=1}^{\frac{n}{2}} \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta) \quad \text{and}$$

$$\nabla_{\theta} J(\theta) \approx \sum_{i=\frac{n}{2}+1}^n \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta).$$

We can do the update with only half the computation cost!

$$\theta_{t+1} = \theta_t - \gamma \frac{1}{n/2} \sum_{i=1}^{\frac{n}{2}} \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta_t),$$

$$\theta_{t+2} = \theta_{t+1} - \gamma \frac{1}{n/2} \sum_{i=\frac{n}{2}+1}^n \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta_{t+1}).$$

# Stochastic gradient descent

Training data																			
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$	$x_{17}$	$x_{18}$	$x_{19}$	$x_{20}$
$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$	$y_8$	$y_9$	$y_{10}$	$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$y_{15}$	$y_{16}$	$y_{17}$	$y_{18}$	$y_{19}$	$y_{20}$

Mini-batch

$$\theta_3 = \theta_2 - \gamma \frac{1}{5} \sum_{i=11}^{15} \nabla_{\theta} L(x_i, y_i, \theta_2)$$

- The extreme version of this strategy is to use only one data point at each training step (called **online learning**)
- We typically do something in between (not one data point, and not all data). We use a smaller set called **mini-batch**.
- One pass through the training data is called an **epoch**.

# Stochastic gradient descent

$\mathbf{x}_7$	$\mathbf{x}_{10}$	$\mathbf{x}_3$	$\mathbf{x}_{20}$	$\mathbf{x}_{16}$	$\mathbf{x}_2$	$\mathbf{x}_1$	$\mathbf{x}_{18}$	$\mathbf{x}_{19}$	$\mathbf{x}_{12}$	$\mathbf{x}_6$	$\mathbf{x}_{11}$	$\mathbf{x}_{17}$	$\mathbf{x}_{15}$	$\mathbf{x}_5$	$\mathbf{x}_{14}$	$\mathbf{x}_4$	$\mathbf{x}_9$	$\mathbf{x}_{13}$	$\mathbf{x}_8$
$\mathbf{y}_7$	$\mathbf{y}_{10}$	$\mathbf{y}_3$	$\mathbf{y}_{20}$	$\mathbf{y}_{16}$	$\mathbf{y}_2$	$\mathbf{y}_1$	$\mathbf{y}_{18}$	$\mathbf{y}_{19}$	$\mathbf{y}_{12}$	$\mathbf{y}_6$	$\mathbf{y}_{11}$	$\mathbf{y}_{17}$	$\mathbf{y}_{15}$	$\mathbf{y}_5$	$\mathbf{y}_{14}$	$\mathbf{y}_4$	$\mathbf{y}_9$	$\mathbf{y}_{13}$	$\mathbf{y}_8$

Iteration: 3

Epoch: 1

- If we pick the mini-batches in order, they might be unbalanced and not representative for the whole data set.
- Therefore, we pick data points **at random** from the training data to form a mini-batch.
- One implementation is to randomly reshuffle the data before dividing it into mini-batches.
- After each epoch we do another reshuffling and another pass through the data set.

# Mini-batch gradient descent

The full **stochastic gradient descent** algorithm (a.k.a **mini-batch gradient descent**) is as follows

1. Initialize  $\theta_0$ , set  $t \leftarrow 1$ , choose batch size  $n_b$  and number of epochs  $E$ .
2. For  $i = 1$  to  $E$ 
  - (a) Randomly shuffle the training data  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ .
  - (b) For  $j = 1$  to  $\frac{n}{n_b}$ 
    - (i) Approximate the gradient of the loss function using the mini-batch  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=(j-1)n_b+1}^{jn_b}$ ,
$$\hat{\mathbf{g}}_t = \frac{1}{n_b} \sum_{i=(j-1)n_b+1}^{jn_b} \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta) \Big|_{\theta=\theta_t}.$$
    - (ii) Do a gradient step  $\theta_{t+1} = \theta_t - \gamma \hat{\mathbf{g}}_t$ .
    - (iii) Update the iteration index  $t \leftarrow t + 1$ .

At each time we get a stochastic approximation of the true gradient  $\hat{\mathbf{g}}_t \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta) \Big|_{\theta=\theta_t}$ , hence the name.

# A few concepts to summarize lecture 9

---

**Convolutional neural network (CNN):** A NN with a particular structure tailored for input data with a grid-like structure, like for example images.

**Kernel:** (a.k.a filter) A set of parameters that is convolved with a hidden layer. Each kernel produces a new channel.

**Channel:** A set of hidden units produced by the same kernel. Each hidden layer consists of one or more channels.

**Stride:** A positive integer deciding how many steps to move the kernel during the convolution.

**Tensor:** A generalization of matrices to arbitrary order.

**Gradient descent:** An iterative optimization algorithm where we at iteration take a step proportional to the negative gradient.

**Learning rate:** (a.k.a step length). A scalar tuning parameter deciding the length of each gradient step in gradient descent.

**Stochastic gradient descent (SGD):** A version of gradient descent where we at each iteration only use a small part of the training data (a mini-batch).

**Mini-batch:** The group of training data that we use at each iteration in SGD

**Batch size:** The number of data points in one mini-batch