# Deep Feedforward Networks

Master-Seminar: Introduction to Deep Learning

---

Bodo Burger

2016-01-13

Ludwig-Maximilians-Universität

## Overview

1. Basics of artificial neural networks
   - Historical notes

2. Architecture design
   - Output units
   - Hidden units

3. Gradient-based learning
   - Back-propagation
   - Other differentiation algorithms

4. Application
   - Odd/Even
   - "multiple of"

# Basics of Artificial Neural Networks

What are deep feedforward networks?

- *the* basic deep learning model
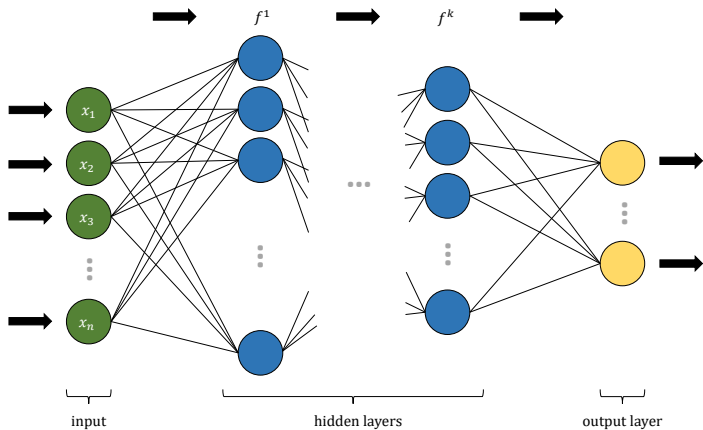- names: feedforward neural networks, multilayer perceptrons (MLPs)

- composing together many different functions

$$f(\boldsymbol{x}) = f^{(3)}(f^{(2)}(f^{(1)}(\boldsymbol{x})))$$

$f^{(1)}$ would be the first layer, $f^{(2)}$ the second layer, and so on …

- length of the chain: depth of the network
-          input: $\boldsymbol{x}$
  output layer: last layer
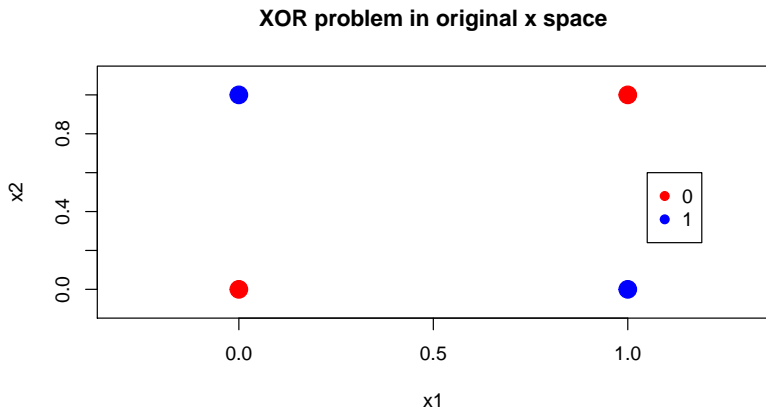  hidden layers: layers in between

- *given:* noisy example $\boldsymbol{x}$, corresponding label $y$
- The goal of the neural network is to approximate a function $f^*(\boldsymbol{x})$ ($y \approx f^*(\boldsymbol{x})$).
- Find optimal parameters $\boldsymbol{\theta}$ so that a function $f$ maps $(\boldsymbol{x}, \boldsymbol{\theta})$ on $\boldsymbol{y}$.

**XOR problem in original x space**

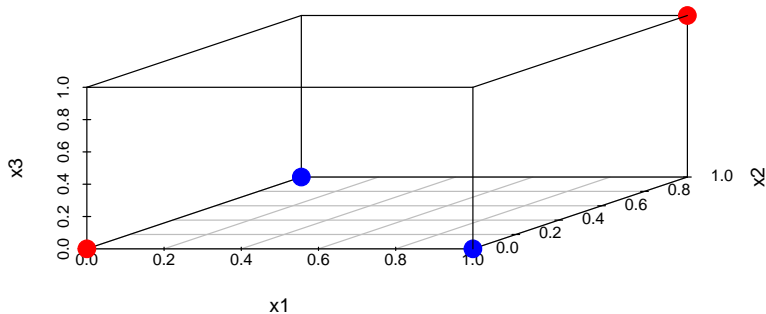Transform the input $x$ by a nonlinear transformation $\phi$.

Possible solutions for $\phi$:

1. kernel trick, SVM
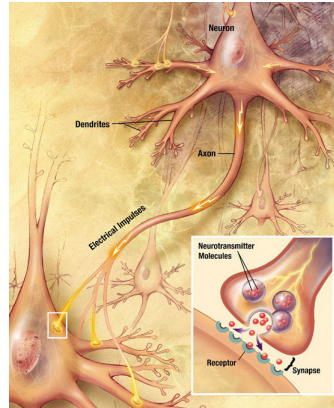2. manually engineer $\phi$
3. learn $\phi$

Adding a third input $x_3 = x_1 \cdot x_2$:
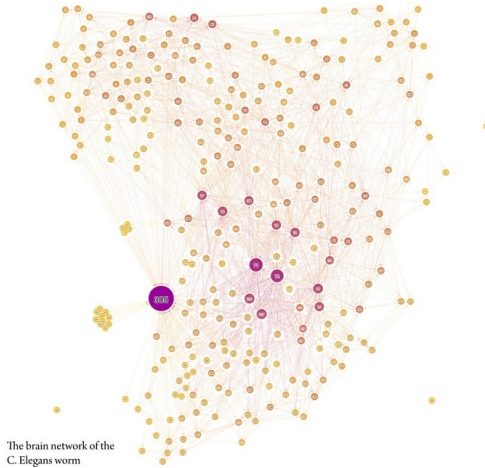


**XOR problem in new space**

**Human brain:** around 85 billion neurons



left: `https://commons.wikimedia.org/wiki/File:Purkinje_cell_by_Cajal.png`

right: `https://commons.wikimedia.org/wiki/File:Chemical_synapse_schema_cropped.jpg`

**C. elegans:** worm with around 300 neurons



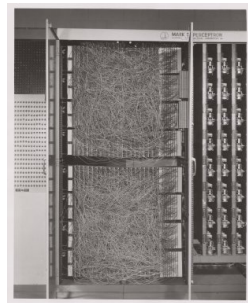The brain network of the
C. Elegans worm

- **More precise:** artificial neural networks (ANN)
- architecture is inspired by biological neural networks
- **but:** not an exact model of the function of a brain
- development of ANN shaped by maths / computer science

# Basics of Artificial Neural Networks

Historical Notes

- 1943: McCulloch / Pitts create first computational model for neural networks
- late 1940s: Hebbian learning
- 1951: Minsky develops SNARC, possibly the first ANN
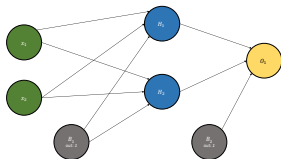- 1958: Rosenblatt creates the Mark I perceptron



https://upload.wikimedia.org/wikipedia/en/5/52/Mark_I_perceptron.jpeg

- 1969: Minsky/Papert "Perceptrons: An Introduction to Computational Geometry"
- limited processing power of computers at the time
- period of defunding of AI research
- 1986: Rumelhart, Hinton and Williams on MLP and back-propagation
- Schmidhuber: "Deep learning since 1991"

- since 2000s: many breakthroughs because of increasing processing power and massive data availability
- Google Translate
- Voice recognition in smartphones
- Amazon Echo / Alexa speech recognition
- DeepMind's AlphaGo
- Self-driving cars

# Architecture Design

What is the structure of the neural network?

- How many units?
- Which units?
- How do you connect these units?

## Chain structure

Typical design approach: multiple layers that consist of multiple units are linked in a chain structure.

First layer is a function of the inputs:

$$\boldsymbol{h}^{(1)} = g^{(1)}\left(\boldsymbol{W}^{(1)T}\boldsymbol{x} + \boldsymbol{b}^{(1)}\right)$$

Second layer is a function of the first layer's output:

$$\boldsymbol{h}^{(2)} = g^{(2)}\left(\boldsymbol{W}^{(2)T}\boldsymbol{h}^{(1)} + \boldsymbol{b}^{(2)}\right)$$

...

$g^{()}$ is called activation function.

Bias puts the decision boundary at the correct position in the learned space.

**without bias**

**with bias**

Depth of the network

number of layers of the network

Width of a layer

number of units in a layer

Deeper networks often work with thiner layers which means less parameters.

**Ideal architecture has to be found via experimentation.**

# Universal approximation theorem

## Universal approximation theorem

"In the mathematical theory of artificial neural networks, the universal approximation theorem states that a feed-forward network with a single hidden layer containing a finite number of neurons (i.e., a multilayer perceptron), can approximate continuous functions on compact subsets of $\mathbb{R}^n$, under mild assumptions on the activation function. The theorem thus states that simple neural networks can represent a wide variety of interesting functions when given appropriate parameters; however, it does not touch upon the algorithmic learnability of those parameters."

https://en.wikipedia.org/wiki/Universal_approximation_theorem

# Why may learning a function fail?

1. optimization algorithm may not find the parameters of the function
2. overfitting

### "No free lunch" theorem

"There is no strategy - applied on all possible problems - that is better than pure guessing."

https://de.wikipedia.org/wiki/No-free-Lunch-Theoreme

Depth and height do not need to be the only parameters of the architecture of a neural network.

**Possible modifications:**

- Add connections that skip a layer and connect layer $i$ to layer $i + 2$
- Do not connect every input of a layer with every output of the preceding layer.

The choice of modifications is very dependent on the actual application!

# Alternative network designs: examples

- Convolutional networks (for computer vision)
- Recurrent neural networks (for sequence processing)
    - Long short-term memories (Google voice, Siri)
- Autoencoder (dimensionality reduction)
- Deep belief network
- …

# Architecture Design

Output Units

The output unit defines what happens to the output $h$ of the last hidden layer. So it is the last transformation on $x$.

The choice of the output unit also determines the cost function.

$$\hat{y} = z = W^T h + b$$



z
one-dimensional linear unit

$$\hat{y} = z = W^T h + b$$

- used for predicting the mean of a Gaussion distribution
- easy to handle for optimization algorithms

# Sigmoid units

Sigmoid function: $\sigma(z) = \frac{1}{1+e^{-z}}$.

Sigmoid function: $\sigma(z) = \frac{1}{1+e^{-z}}$.

$$\hat{y} = \sigma\left(\boldsymbol{w}^T\boldsymbol{h} + b\right)$$

- applying the sigmoid function on a linear unit
- used for predicting binary outputs
- problematic: saturation for large absolute values of z

Signum function

# Softmax units

Softmax function: $\text{softmax}(\boldsymbol{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$.

- generalization of sigmoid function for multiclass problem
- used for predicting discrete output with k possible values

# Architecture Design

Hidden Units

- active area of research, not many definitive theoretical guidelines
- Try out different kinds of units and look which work best.
- good default choice: Rectified linear units
- Hidden units do not have to be differentiable.

$$g(z) = \max\{0, z\}$$

- typically used with a bias $b$ (set to small, positive value)
- not differentiable at zero, but some generalizations of a ReLU are

There are various generalizations of ReLU:

- absolute value rectification
- leaky ReLU
- parametric ReLU (PReLU)

$$h_i = g(\boldsymbol{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

Absolute value rectification: $\alpha_i = -1 \Rightarrow g(z) = |z|$.

$$h_i = g(\boldsymbol{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

leaky ReLU: $\alpha_i = 0.01$.

$$h_i = g(\boldsymbol{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

**parametric ReLU:** treat $\alpha_i$ as learnable parameter.

$$g(\boldsymbol{z})_i = \max z_j$$

- divide $\boldsymbol{z}$ into groups of $k$ values
- learning a piecewise linear function
- Maxout units avoid catastrophic forgetting.

## Logistic sigmoid and hyperbolic tangent

Sigmoid function:

$$g(z) = \sigma(z)$$

Hyperbolic tangent function:

$$g(z) = tanh(z)$$

Closely related, because $tanh(z) = 2\sigma(2z) - 1$.

Used e.g. for recurrent networks or autoencoders.

- radial basis function
- sofplus
- hard tanh

two inputs, one hidden layer with two units



sigmoid activation function for hidden layer and output

initialising random weights for all connections

$x_1 = 1, x_2 = 0$:



$H_1$: in $= -0.07 \cdot 1 + 0.22 \cdot 0 + -0.46 \cdot 1 = -0.53$

# Forward propagation

$x_1 = 1, x_2 = 0$:



$H_1$: out $= \frac{1}{1-e^{-0.53}}$

$x_1 = 1, x_2 = 0$:



$H_2$: in $= -0.07 \cdot 1 + 0.22 \cdot 0 + -0.46 \cdot 1 = -0.53$

$x_1 = 1, x_2 = 0$:

$H_2$: out $= \frac{1}{1 - e^{-0.53}}$

# Forward propagation

$x_1 = 1, x_2 = 0$:



$O_1$: in $= -0.07 \cdot 1 + 0.22 \cdot 0 + -0.46 \cdot 1 = -0.53$

$x_1 = 1, x_2 = 0$:



$O_1$: out $= \frac{1}{1-e^{-0.53}}$

Error: $E = f(\boldsymbol{x}) - y = \hat{y} - y$

$E = 0.76 - 1 = -0.24$

Total error can be defined as

$$E = \frac{1}{2} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2.$$

For our example:

$$E = \frac{1}{2} \left[ (0.73 - 1)^2 + (0.74 - 1)^2 + (0.76 - 1)^2 + (0.76 - 1)^2 \right]$$
$$= 0.12785$$

# Gradient-Based Learning

## Why gradient-based learning?

- Nonlinearity of neural nets causes loss functions to become non-convex.
- Convex optimization algorithms do not work anymore.
- Use iterative, gradient-based optimization.
- Does not guarantee convergence and results may depend heavily on initial parameters.
- For *very large* data sets, it does make sense to train a linear model or SVM by gradient descent, too.

# Gradient-Based Learning

Back-Propagation

### Forward propagation

The information of the inputs $x$ flows through the hidden units to finally produce $\hat{y}$ and the cost $J(\theta)$.

### Back-propagation

The information of the cost $J(\theta)$ flows backwards through the hidden units to calculate the gradient.

## Calculating the gradient (output layer)

Total error is defined as: $E = \frac{1}{2} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2$

**Gradient of connection from unit i to output $O_1$:**

$$\frac{\partial E}{\partial w_{O_1,i}} = \frac{\partial E}{\partial \text{out}_{O_1}} \cdot \frac{\partial \text{out}_{O_1}}{\partial \text{in}_{O_1}} \cdot \frac{\partial \text{in}_{O_1}}{\partial w_{O_1,i}}$$

If we define $\delta_{O_1} = (\hat{y}_i - y_i) \cdot \sigma'_{O_1}$, then we get

$$\frac{\partial E}{\partial w_{O_1,i}} = \delta_{O_1} \text{out}_i,$$

with $\text{out}_{O_1} = \hat{y}_i$ and $\frac{\partial \text{out}_{O_1}}{\partial \text{in}_{O_1}} = \sigma'_{O_1}$.

Gradient of connection from unit i to unit j:

$$\frac{\partial E}{\partial w_{j,i}} = \frac{\partial E}{\partial out_j} \cdot \frac{\partial \text{out}_j}{\partial \text{in}_j} \cdot \frac{\partial \text{in}_j}{\partial w_{j,i}}$$

Node delta is now defined as $\delta_j = \sigma_j' \sum_k w_{k,j} \cdot \delta_k$. We get

$$\frac{\partial E}{\partial w_{j,i}} = \delta_j \text{out}_i.$$

Delta of output node:



$$\delta_{O_1} = -(-0.24) \cdot d\sigma(1.13) = 0.24 \cdot \sigma(1.13)(1 - \sigma(1.13)) = 0.0452$$

47

Delta of hidden node 1:



$$\delta_{H_1} = d\sigma(-0.53) \cdot (-0.22 \cdot 0.0452) = -0.0023$$

Delta of hidden node 2:



$$\delta_{H_2} = d\sigma(1.04) \cdot (0.58 \cdot 0.0452) = 0.0051$$

Gradient from $H_1$ to $O_1$:



$$G_{O_1,H_1} = 0.37 \cdot 0.0452 = 0.0168$$

Gradient from $H_2$ to $O_1$:



$G_{O_1, H_2} = 0.74 \cdot 0.0452 = 0.0334$

Gradient from $B_2$ to $O_1$:



$G_{O_1,B_2} = 1 \cdot 0.0452 = 0.0452$

Gradient from $x_1, x_2, B_1$ to $H_1$:



$$G_{H_1, \cdot} = \begin{pmatrix} -0.0023 & 0 & -0.0023 \end{pmatrix}$$

Gradient from $x_1, x_2, B_1$ to $H_2$:



$\delta_{H_1} = -0.0023$

in: *-0.53*
$H_1$
out: *0.37*

w: -0.07
G: -0.0023

$x_1$
= 1.0

w: 0.94
G: 0.0051

w: -0.22
G: 0.0168

$\delta_{O_1} = 0.0452$

in: *1.13*
$O_1$
out: *0.76*

w: 0.22
G: 0

$\delta_{H_2} = 0.0051$

in: *1.04*
$H_2$
out: *0.74*

w: 0.58
G: 0.0334

$x_2$
= 0.0

w: 0.46
G: 0

w: -0.46
G: -0.0023

w: 0.10
G: 0.0051

w: 0.78
G: 0.0452

$B_1$
out: *1*

$B_2$
out: *1*

$$G_{H_2,\cdot} = \begin{pmatrix} 0.0051 & 0 & 0.0051 \end{pmatrix}$$

$$\Delta w^{(t)} = -\epsilon \frac{\partial E}{\partial w^{(t)}} + \alpha \Delta w^{(t-1)}$$

- learning rate $\epsilon$
    - too high: may fail to converge
    - too low: very slow learning

- momentum $\alpha$
    - $0 \leq \alpha < 1$
    - can speed up learning
    - can help to avoid local minima

- $\epsilon$ and $\alpha$ have to be determined by trial and error or by experience

- stochastic learning:
  - random select a training set element, modify weights and repeat
  - can help to avoid local minima, but slow
- batch learning:
  - gradients for each training set element are summed up, then weights are updated
  - faster as it incorporates the average error of all elements
- mini-batch: select several training set elements and compute a weight update

New weights $w_{O_1,\cdot}$:



$$w_{O_1,\cdot} = \begin{pmatrix} -0.22 & 0.58 & 0.78 \end{pmatrix} + 0.7 \cdot \begin{pmatrix} 0.0168 & 0.0334 & 0.0452 \end{pmatrix} + 0.3 \cdot 0$$

# Weight update

New weights $w_{H_1,\cdot}$:



$$w_{H_1,\cdot} = \begin{pmatrix} -0.07 & 0.22 & -0.46 \end{pmatrix} + 0.7 \cdot \begin{pmatrix} -0.0023 & 0 & -0.0023 \end{pmatrix} + 0.3 \cdot 0$$

New weights $w_{H_2,\cdot}$:



$$w_{H_2,\cdot} = \begin{pmatrix} 0.94 & 0.46 & 0.10 \end{pmatrix} + 0.7 \cdot \begin{pmatrix} 0.0051 & 0 & 0.0051 \end{pmatrix} + 0.3 \cdot 0$$

## Total error for new weights

Total error before adjusting weights was 0.12785.

Now:

$$E = \frac{1}{2} \left[ (0.74 - 1)^2 + (0.75 - 1)^2 + (0.77 - 1)^2 + (0.77 - 1)^2 \right]$$
$$= 0.11795$$

Cross-entropy error function:

$$J = -\frac{1}{n} \sum_{i=1}^{n} [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

The node delta for the output layer becomes: $\delta_O = \hat{y} - y$.

Cross-entropy error function often performs better than MSE for learning a classification task.

# Gradient-Based Learning

Other Algorithms

## Resilient back-propagation (Rprop)

Rprop weight change:

$$\Delta w_{j,i}^{(t)} = \begin{cases} -\Delta_{j,i}^{(t)}, \text{if } c > 0 \\ +\Delta_{j,i}^{(t)}, \text{if } c < 0 \\ 0, otherwise \end{cases}$$

Modify update values:

$$\Delta_{j,i}^{(t)} = \begin{cases} \eta^+ \cdot \Delta_{j,i}^{(t-1)}, \text{if } c > 0 \\ \eta^- \cdot \Delta_{j,i}^{(t-1)}, \text{if } c < 0 \\ \Delta_{j,i}^{(t-1)}, otherwise \end{cases}$$

LMA is a very efficient training method for neural networks.

It combines the Gauss-Newton algorithm and the method of gradient descent.

$$W_{min} = W_0 - (H + \lambda I)^{-1} g$$

# Example Odd/Even

```
print(X[100,1:11])
print(X[100,12:22])
print(X[100,23:33])
print(X[100,34:44])
print(X[100,45:55])

## [1] 0 1 0 0 1 1 1 0 1 1 1
## [1] 0 1 0 0 1 0 1 0 1 0 1
## [1] 0 1 0 0 1 0 1 0 1 0 1
## [1] 0 1 0 0 1 0 1 0 1 0 1
## [1] 0 1 0 0 1 1 1 0 1 1 1
```

```r
y <- rep(c(-1,1), 500)[-1000]
y[5]
y[56]
y[998]
y[999]

## [1] -1
## [1] 1
## [1] 1
## [1] -1
```

```
perceptron(X = X, y = y, w = w, eta = 0.1,
           max_reps = 80, max_error_rate = 0.05)
```

- a single neuron
- uses signum function
- stochastic learning
- weight adjustment: $w \leftarrow w + \eta \cdot x \cdot y$

# Perceptron function

```
perceptron(X = X, y = y, w = w, eta = 0.1,
          max_reps = 80, max_error_rate = 0.05)
```

- X: feature matrix
- y: target vector
- w: weights vector
- eta: learning rate
- max_reps: maximum weight update iterations
- max_error_rate: stop learning, if missclassification rate ≤ max_error_rate

```
# number of features and number of obs
m = ncol(X)

# initial weights:
w = rnorm(m, 1)/10
w = rep(0.01, m)
```

```
set.seed(909)
fit = perceptron(X, y, w, eta = 0.1,
      max_reps =  80, max_error_rate = 0)

## [1] "At least 100% of all points correctly class
## [1] "Number of iterations:"
## [1] 70
```

Initial weights:

1 iteration:



Used number: 21

2 iterations:



Used numbers: 21, 984

# Weight changes over iterations

3 iterations:



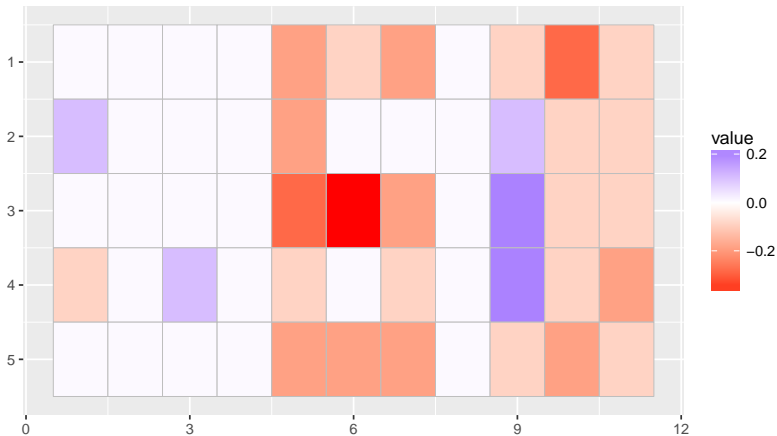Used numbers: 21, 984, 49

# Weight changes over iterations

4 iterations:



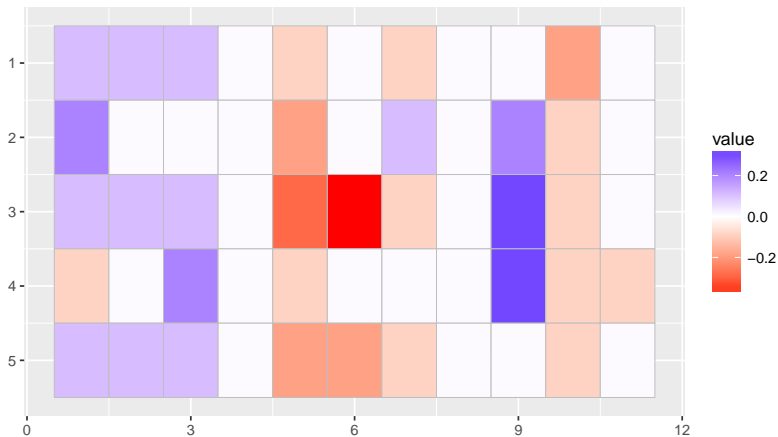Used numbers: 21, 984, 49, 293

# Weight changes over iterations
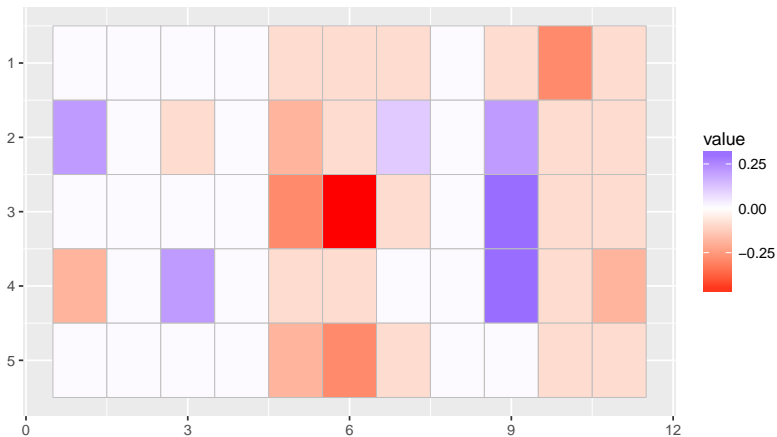
5 iterations:



Used numbers: 21, 984, 49, 293, 372

# Weight changes over iterations

6 iterations:



Used numbers: 21, 984, 49, 293, 372, 508

7 iterations:

Used numbers: 21, 984, 49, 293, 372, 508, 363

# Weight changes over iterations

8 iterations:



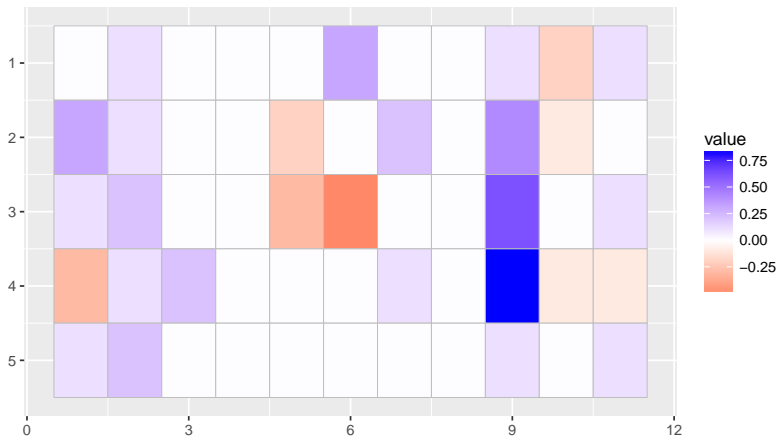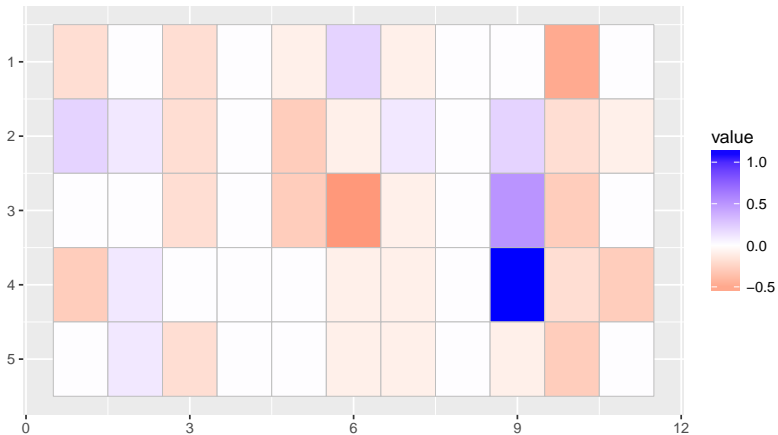Used numbers: 21, 984, 49, 293, 372, 508, 363, 567

# Weight changes over iterations

9 iterations:



Used numbers: 21, 984, 49, 293, 372, 508, 363, 567, 570

10 iterations:



Used numbers: 21, 984, 49, 293, 372, 508, 363, 567, 570, 217

# Weight changes over iterations

20 iterations:

Used numbers: 21, 984, 49, 293, 372, 508, 363, 567, 570, 217, 516, 745, 938, 349, 100, 533, 972, 649, 96, 344
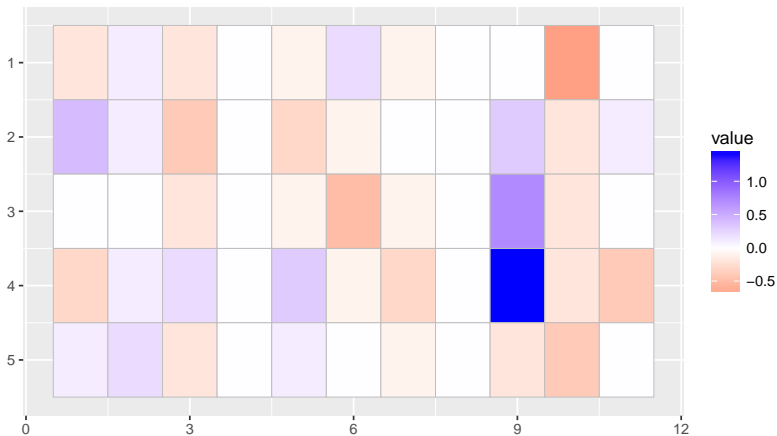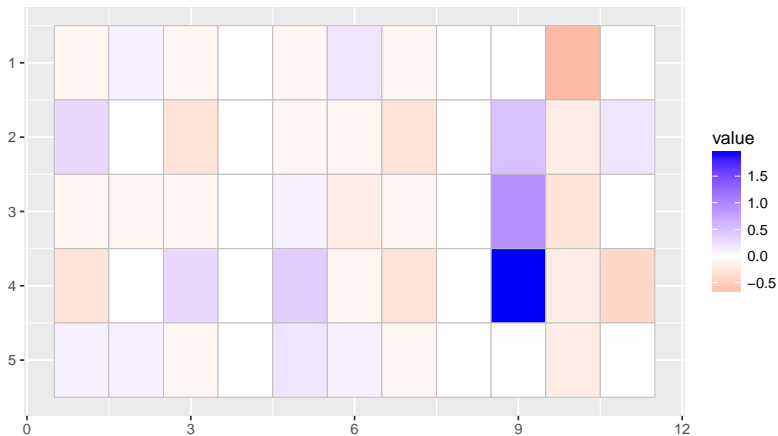
# Weight changes over iterations
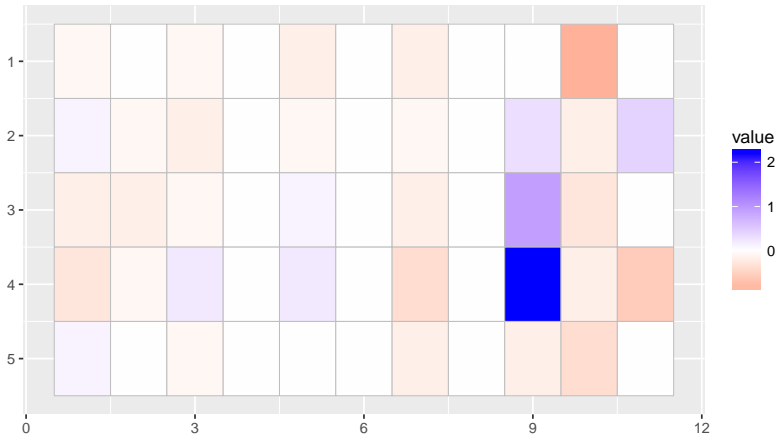
30 iterations:

40 iterations:
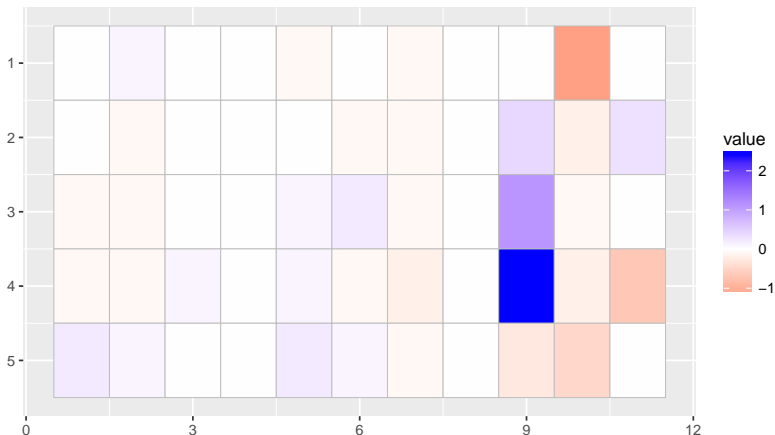
50 iterations:

60 iterations:

# Weight changes over iterations

Final weights after 70 iterations:

- "multiple of 5"
- "multiple of 4"
- "multiple of 3"
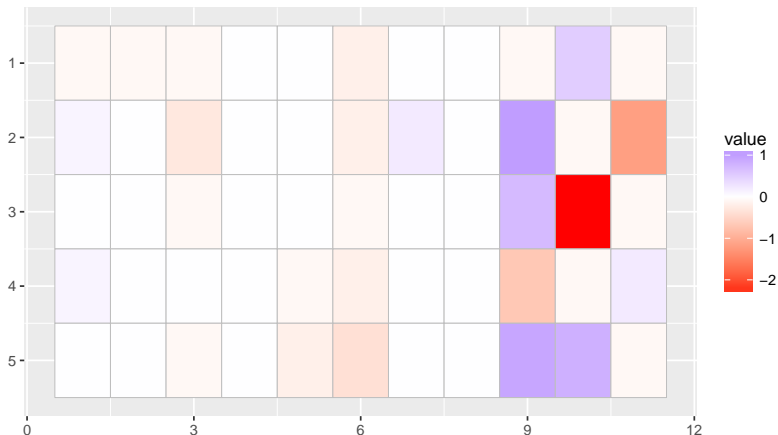
## New target: multiple of 5

```
y = rep(-1, 999)
y[(1:999) %% 5 == 0] = 1

set.seed(909)
fit = perceptron(X, y, w, eta = 0.1,
     max_reps =  300, max_error_rate = 0.05)

## [1] "At least 95% of all points correctly classi
## [1] "Number of iterations:"
## [1] 112
```

Final weights for "multiple of 5":

# References

TODO