

Module 11: Neural Networks

TMA4268 Statistical Learning V2021

Stefanie Muff, Department of Mathematical Sciences, NTNU

April 19 and 20, 2021

Acknowledgements

- The original version of this material stems from Mette Langaas who has put a lot of effort in creating this module. Thanks to Mette for the permission to use the material!

Learning material for this module

Primary material:

These slides

Secondary material:

- Videos on neural networks and back propagation
 - Video 1
 - Video 2
 - Video 3
 - Video 4
- Background material: Chapters 6-8 Goodfellow, Bengio, and Courville (2016) <https://www.deeplearningbook.org>

See also *References and further reading* (last slide), for further reading material.

Neural networks

→ prediction

What will you learn?

- Translating from statistical to neural networks language
 - linear regression
 - logistic regression
 - multiclass (multinomial) regression
- Feedforward networks
- Neural network parts: model – method – algorithm – recent developments
- Deep learning
 - The timeline
 - Keras

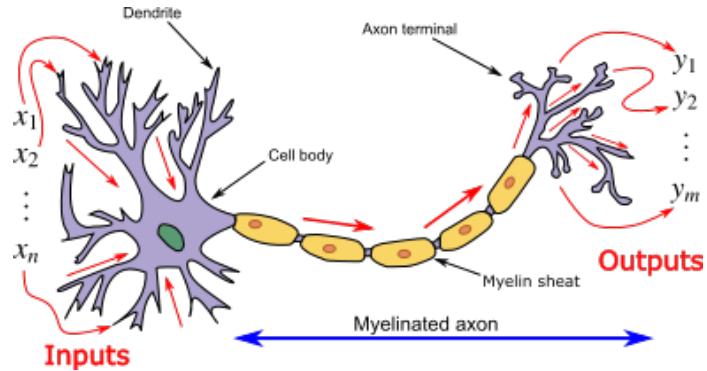
Introduction

- Neural networks (NN) were first introduced in the 1990's.
- Shift from statistics to computer science and machine learning, as they are highly parameterized
- Statisticians were skeptical: “It's just a nonlinear model”.
- After the first hype, NNs were pushed aside by boosting and support vector machines.
- Revival since 2010: The emergence of *Deep learning* as a consequence of improved computer resources, some innovations, and applications to image and video classification, and speech and text processing

Why a module on neural networks?

- Every day you read about the success of AI, machine learning – and in particular *deep learning*.
- In the last five years the field of deep learning has gone from low level performance to excellent performance – particularly in image recognition and speech transcription.
- Deep learning is based on a layered artificial neural network structure.

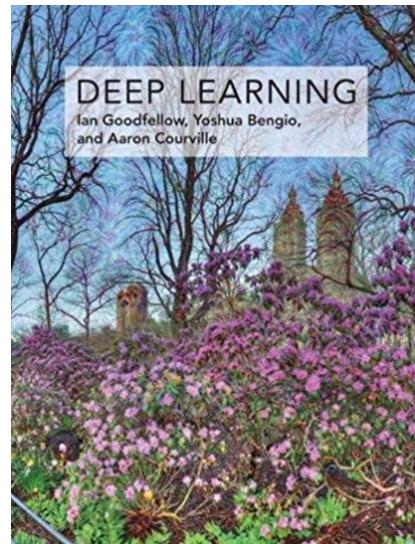
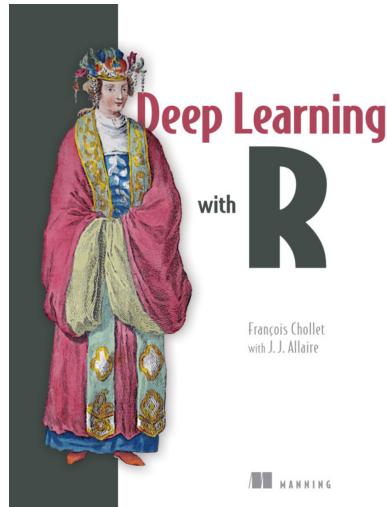
So we first need to understand: what is a *neural network*?



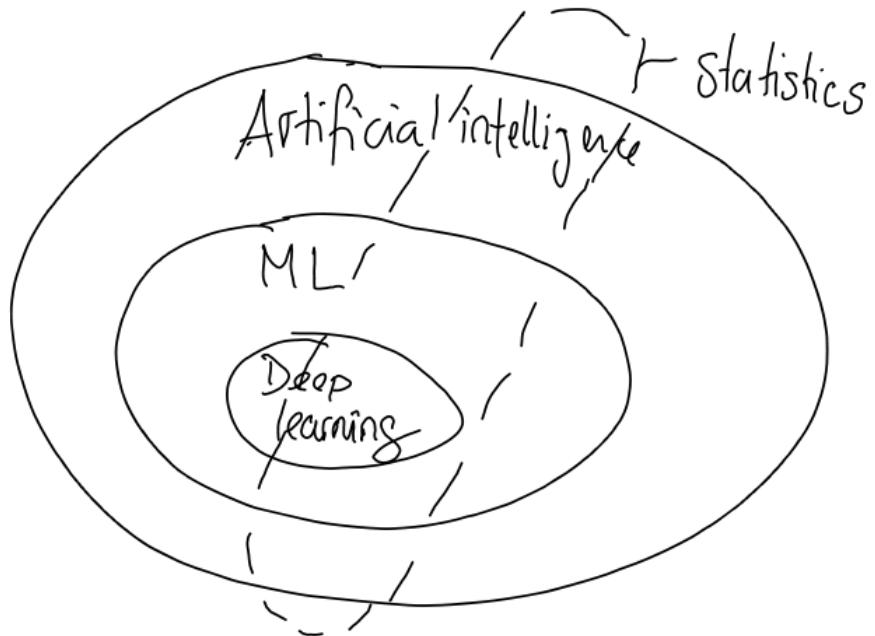
Neuron and myelinated axon, with signal flow from inputs at dendrites to outputs at axon terminals. Image credits: By Egm4313.s12 (Prof. Loc Vu-Quoc)

<https://commons.wikimedia.org/w/index.php?curid=72816083>

- There are several (self-study) learning resources (some listed under ‘further references’) that the students my turn to for further knowledge into deep learning, but this presentation is heavily based on Chollet and Allaire (2018), with added formulas and theory.
- There is a new IT3030 deep learning course at NTNU.



AI, machine learning and statistics



- Artificial intelligence (AI) dates back to the 1950s, and can be seen as *the effort to automate intellectual tasks normally performed by humans* (page 4, Chollet and Allaire (2018)).
- AI was first based on hardcoded rules (like in chess programs), but turned out to be intractable for solving more complex, fuzzy problems.

Machine learning

- With the field of *machine learning* the shift is that a system is *trained* rather than explicitly programmed.
- Machine learning is related to mathematical statistics, but differs in many ways.
- ML deals with much larger and more complex data sets than what is usually done in statistics.
- The focus in ML is oriented towards *engineering*, and ideas are proven *empirically* rather than theoretically (which is the case in mathematical statistics).

According to Chollet and Allaire (2018) (page 19):

Machine learning isn't mathematics or physics, where major advancements can be done with a pen and a piece of paper. It's an engineering science.

Deep learning

- *Deep* does not refer to a deeper understanding.
- Rather, deep refers to the *layers of representation*, for example in a neural network.

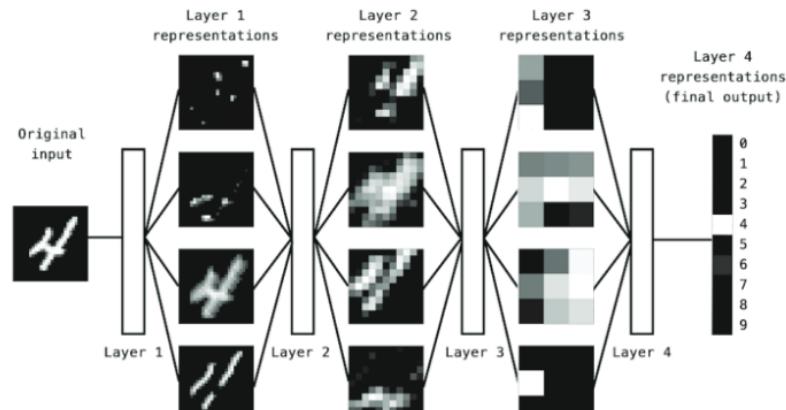


Figure 1.6 Deep representations learned by a digit classification model

From statistics to artificial neural networks

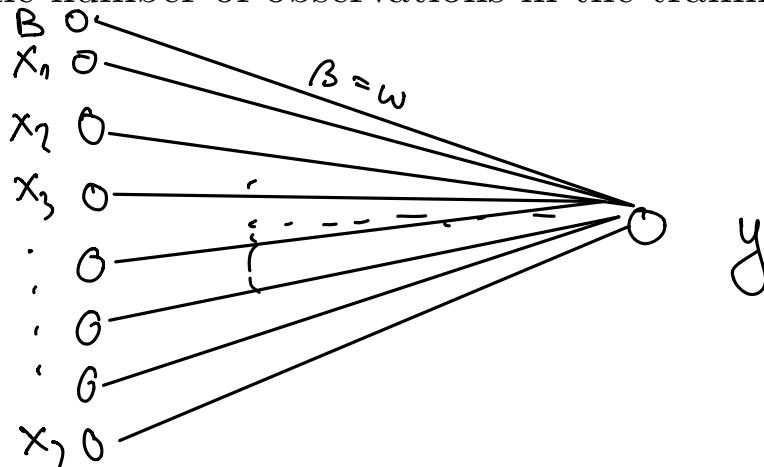
Recapitulate from Module 3 with the bodyfat dataset that contained the following variables.

- y • bodyfat: % of body fat.
- x₁ • age: age of the person.
- x₂ • weight: body weight.
- x₃ • height: body height.
- x₄ • neck: neck thickness.
- x₅ • bmi: bmi.
- x₆ • abdomen: circumference of abdomen.
- x₇ • hip: circumference of hip.

We will now look at modelling the **bodyfat** as response and using all other variables as covariates - this will give us

- seven covariates
- one intercept
- one numerical output

Let n be the number of observations in the training set, here $n = 243$.



Multiple linear regression model (from Module 3)

We assume

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i = x_i^T \beta + \varepsilon_i ,$$

for $i = 1, \dots, n$, where x_{ij} is the value j th predictor for the i th datapoint, and $\beta^\top = (\beta_0, \beta_1, \dots, \beta_p)$ the regression coefficients.

We used the compact matrix notation for all observations $i = 1, \dots, n$ together:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} Y = X\beta + \varepsilon .$$

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

Assumptions:

1. $E(\varepsilon) = 0$.
2. $\text{Cov}(\varepsilon) = E(\varepsilon \varepsilon^T) = \sigma^2 I$.
3. The design matrix has full rank, $\text{rank}(X) = p + 1$. (We assume $n \gg (p + 1)$.)

The classical *normal* linear regression model is obtained if additionally

4. $\varepsilon \sim N_n(0, \sigma^2 I)$ holds. Here N_n denotes the n -dimensional multivariate normal distribution.

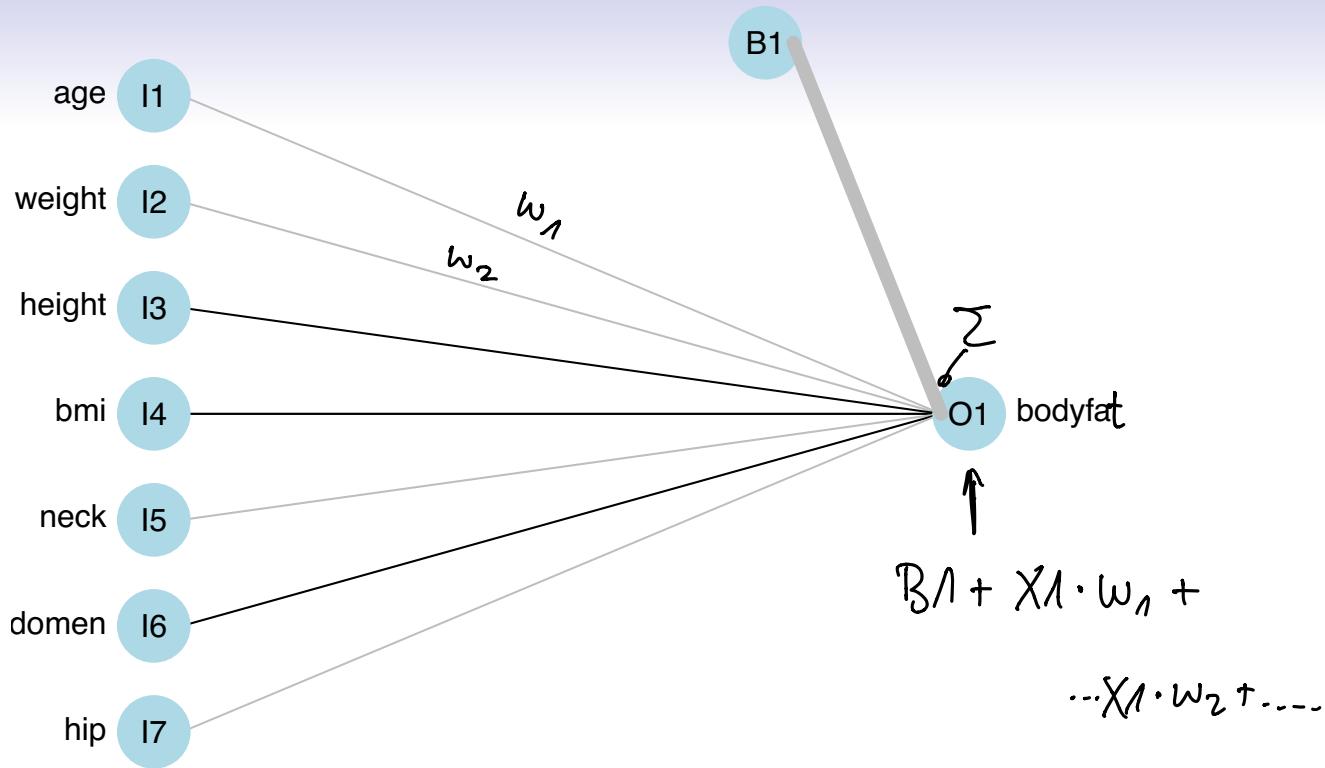


From statistical model to network architecture

How can our statistical model be represented as a network?

We need *new concepts*:

- Covariates → *input nodes* in an *input layer*.
- The intercept → *bias* node.
- The response → *output node* in an *output layer*.
- The regression coefficients → *weights* (often written on the arrows from the inputs to the output layer).
 β w



- All lines going into the output node signifies that we multiply the covariate values in the input nodes with the *weights* (regression coefficients), and then sum.
- This sum can be sent through a so-called *activation function* (here just the identity function).

- *Regression notation:*

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i .$$

- *Neural network notation:*

activation function ↓
 $y_1(x_i) = \phi_o(w_0 + w_1 x_{i1} + \cdots + w_p x_{ip}) ,$

where $\phi_o(x) = x$ (identity function).

In the statistics world

- we would have written $\hat{y}_1(x_i)$ to specify that we are estimating a predicted value of the response for the given covariate value.
- we would have called the $ws \ \hat{\beta}s$ instead.

Statistical parameter estimation

Remember: The estimator $\hat{\beta}$ is found by minimizing the RSS for a multiple linear regression model:

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \dots - \hat{\beta}_p x_{ip})^2 \\ &= \sum_{i=1}^n (y_i - x_i^T \beta)^2 = (Y - X\hat{\beta})^T (Y - X\hat{\beta}) . \end{aligned}$$

Solution:

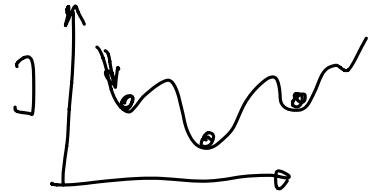
$$\hat{\beta} = (X^T X)^{-1} X^T Y .$$

Neural networks: loss function and gradient descent

We now translate from the statistical into the neural networks world:

1. Replace the parameters β with *network weights* w .
2. Replace the RSS in our *training data set* with the following **loss function** (*mean squared error*)

$$J(w) = \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_1(x_i))^2}_{\text{RSS}} ,$$



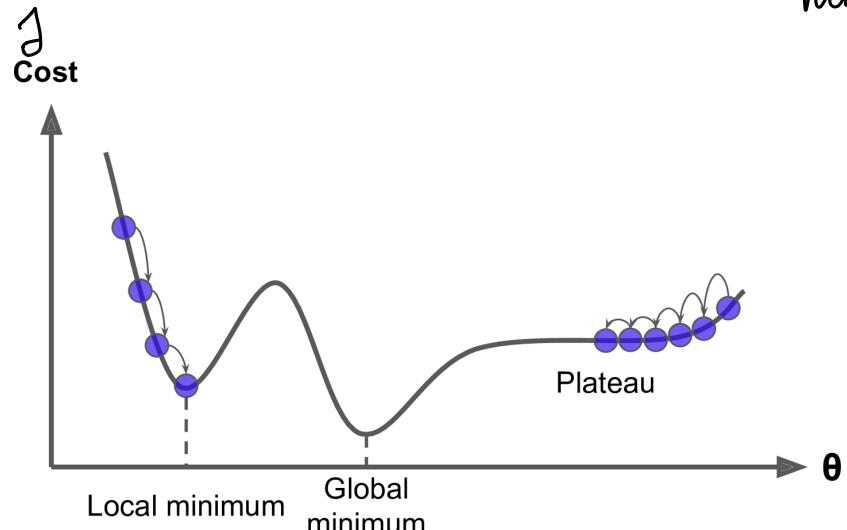
where $J(w)$ indicates that the unknown parameters are the weights w .

3. Replace *minimizing* the loss function (RSS) by **more general minimization procedures** that work also when the loss function does not have a closed form. In particular, use the the **gradient descent algorithm**.

Gradient descent and local minima

Loss function $J(w)$

→ We try to
minimize



(<https://github.com/SoojungHong/MachineLearning/wiki/Gradient-Descent>)

Finding optimal weights: Gradient descent algorithm

1. Let $t = 0$ and denote the given initial values for the weights $w^{(t)}$,
2. Until finding a (local) optimum, repeat a) to e)
 - a) Calculate the predictions $\hat{y}_1(x_i)$.
 - b) Calculate the loss function $J(w^{(t)})$.
 - c) Find the gradient (direction) in the $(p + 1)$ -dimensional space of the weights, and evaluate this at the current weight values
$$\nabla J(w^{(t)}) = \frac{\partial J}{\partial w}(w^{(t)}).$$
 - d) Go with a given step length (learning rate) λ in the direction of the negative of the gradient of the loss function to get

$$\underbrace{w^{(t+1)}}_{\substack{\text{updated} \\ \text{weights}}} = \underbrace{w^{(t)}}_{\substack{\text{old} \\ \text{weights}}} - \lambda \underbrace{\nabla J(w^{(t)})}_{\substack{\text{learning rate}}}. \quad \text{steepest increase}$$

- e) Set $t = t + 1$.
3. The final values of the weights in that $(p + 1)$ dimensional space are our parameter estimates and your network is *trained*.

λ is very important

Backpropagation

- In neural networks the gradient part of the gradient descent algorithm is implemented efficiently in an algorithm called *backpropagation* (see later).

Here we compare

- the MLR solution with `lm`
- the neural network solution with `nnet`¹

¹Uses an improved gradient descent with Hessian information and the BFGS-algorithm. BFGS is a quasi-Newton method (also known as a variable metric algorithm), specifically that published simultaneously in 1970 by Broyden, Fletcher, Goldfarb and Shanno. This uses function values and gradients to build up a picture of the surface to be optimized.

Example 1: Continuous outcome (linear regression)

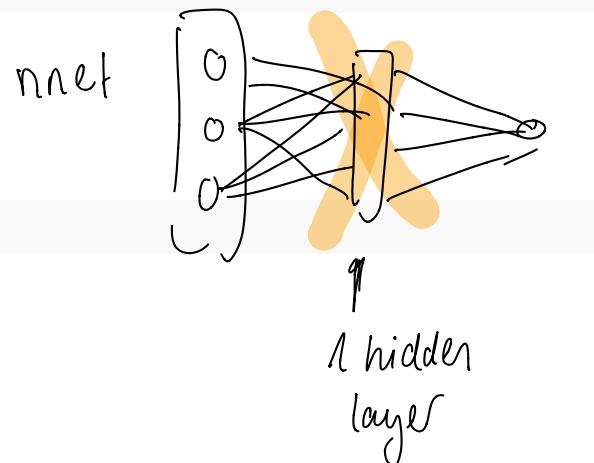
Linear regression vs. neural networks: an example.

```
fit = lm(bodyfat ~ age + weight + height + bmi + neck + abdomen + hip,  
        data = d.bodyfat)  
fitnnet = nnet(bodyfat ~ age + weight + height + bmi + neck + abdomen +  
               hip, data = d.bodyfat, linout = TRUE, size = 0, skip = TRUE, maxit = 1000,  
               entropy = FALSE)  
  
## # weights:  8  
## initial  value 1827258.372491  
## iter  10 value 4470.017765  
## final  value 4415.453729  
## converged  
cbind(fitnnet$wts, fit$coefficients)
```

lm

nnet

	[,1]	[,2]
## (Intercept)	-9.748915e+01	-9.748903e+01
## age	-9.607628e-04	-9.607669e-04
## weight	-6.292828e-01	-6.292820e-01
## height	3.974890e-01	3.974884e-01
## bmi	1.785333e+00	1.785330e+00
## neck	-4.945725e-01	-4.945725e-01
## abdomen	8.945189e-01	8.945189e-01
## hip	-1.255549e-01	-1.255549e-01



Example 2: Binary outcome (logistic regression)

Aim is to predict if a person has diabetes². The data stem from a population of women of Pima Indian heritage in the US, available in the R MASS package. The following information is available for each woman:

- Y
- **diabetes**: 0= not present, 1= present
- X
- npreg: number of pregnancies
 - glu: plasma glucose concentration in an oral glucose tolerance test
 - bp: diastolic blood pressure (mmHg)
 - skin: triceps skin fold thickness (mm)
 - bmi: body mass index (weight in kg/(height in m)²)
 - ped: diabetes pedigree function.
 - age: age in years

²Logistic regression is the "hello world" of machine learning.

The statistical model: Logistic regression

- $i = 1, \dots, n$ observations in the training set. We will use r (instead of p) to be the number of covariates, to avoid confusion with the probability p .
- The binary response $Y_i \in \{1, 0\}$ with

$$P(Y_i = 1) = p_i$$

and

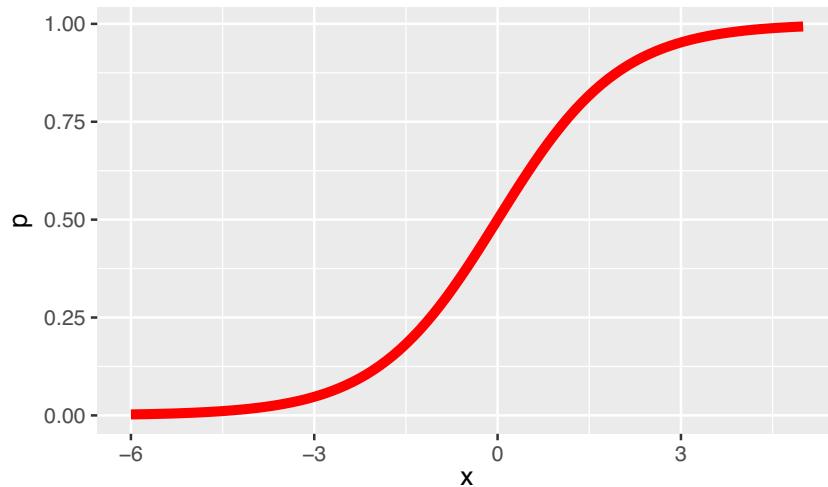
$$\log \left(\frac{p_i}{1 - p_i} \right) = \underbrace{\beta^\top x_i}_{\text{}} ,$$

where $\log(\frac{x}{1-x})$ is the *logistic function*. $\beta_0 + \beta_1 x^1 + \dots + \beta_r x^r$

- Therefore

$$p_i = \frac{\exp(\beta^\top x_i)}{1 + \exp(\beta^\top x_i)} = \frac{1}{1 + \exp(-\beta^\top x_i)}.$$

- This function is S-shaped, and ranges between 0 and 1 (so the p_i is between 0 and 1).



Parameter estimation in the statistical model (Maximum likelihood)

- Given n independent pairs of covariates and responses $\{x_i, y_i\}$, the log-likelihood function of a logistic regression model can be written as:

$$L(\beta) = \prod p_i^{y_i} (1-p_i)^{1-y_i}$$

$$\ln(L(\beta)) = l(\beta) = \sum_{i=1}^n \left(y_i \ln p_i + (1 - y_i) \ln(1 - p_i) \right).$$

- Maximization: set the $r + 1$ partial derivatives (to form the gradient) to 0.
- No closed form solution, thus we use a **gradient-based method** (the Newton-Raphson or **Fisher scoring algorithm**) to find $\hat{\beta}$ and $SD(\hat{\beta})$:

$$\beta^{(t+1)} = \beta^{(t)} + F(\beta^{(t)})^{-1} s(\beta^{(t)}) , \quad g^{(M(-))}$$

where the gradient of the log-likelihood $s(\beta) = \frac{\partial l}{\partial \beta}$ is called the score vector, and here the new quantity $F(\beta^{(t)})^{-1}$ is called the inverse expected Fisher information matrix.

The neural network model: architecture and activation function

- Remember: in the neural network (NN) version of *linear regression*, we had:

$$y_1(x_i) = w_0 + w_1 x_{i1} + \cdots + w_r x_{ir} ,$$

with activation function $\phi_o(x) = x$.

- In the NN version of *logistic regression* we instead have the *sigmoid activation function* $\phi_o(x) = \frac{1}{1+\exp(-x)}$, often denoted as $\sigma(x)$. Again, we prefer to use $\hat{y}_1(x_i)$ and get:

$$\hat{y}_1(x_i) = \sigma(x_i) = \frac{1}{1 + \exp\left(-\underbrace{(w_0 + w_1 x_{i1} + \cdots + w_r x_{ir})}_{\text{linear predictor (in regression)}}\right)} \in (0, 1) .$$

Neural networks: loss function and gradient descent

- For NNs we use *binomial cross-entropy loss*

$$J(w) = -\frac{1}{n} \sum_{i=1}^n (y_i \ln(\hat{y}_1(x_i)) + (1 - y_i) \ln(1 - \hat{y}_1(x_i))) ,$$

which is a scaled version of the negative of the binomial loglikelihood!

- Optimization is done also with *gradient descent*, but we need the chain rule (due to the activation function) to get the partial derivatives for the gradient direction.
- We need the *backpropagation algorithm*, using the activation and loss functions given here.

Parameter estimation vs. network weights

```
fitlogist = glm(diabetes ~ npreg + glu + bp + skin + bmi + ped + age,
                 data = train, family = binomial(link = "logit"))
summary(fitlogist)

##
## Call:
## glm(formula = diabetes ~ npreg + glu + bp + skin + bmi + ped +
##      age, family = binomial(link = "logit"), data = train)
##
## Deviance Residuals:
##       Min      1Q  Median      3Q     Max
## -1.9830 -0.6773 -0.3681  0.6439  2.3154
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -9.773062   1.770386 -5.520 3.38e-08 ***
## npreg        0.103183   0.064694  1.595  0.11073
## glu          0.032117   0.006787  4.732 2.22e-06 ***
## bp           -0.004768   0.018541 -0.257  0.79707
## skin         -0.001917   0.022500 -0.085  0.93211
## bmi          0.083624   0.042827  1.953  0.05087 .
## ped          1.820410   0.665514  2.735  0.00623 **
## age          0.041184   0.022091  1.864  0.06228 .
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 256.41  on 199  degrees of freedom
## Residual deviance: 178.39  on 192  degrees of freedom
## AIC: 194.39
##
## Number of Fisher Scoring iterations: 5
```

```

set.seed(787879)
library(nnet)
fitnnet = nnet(diabetes ~ npreg + glu + bp + skin + bmi + ped + age,
  data = train, linout = FALSE, size = 0, skip = TRUE, maxit = 1000,
  entropy = TRUE, Wts = fitlogist$coefficients + rnorm(8, 0, 0.1))

## # weights:  8
## initial value 213.575955
## iter  10 value 89.511044
## final value 89.195333
## converged

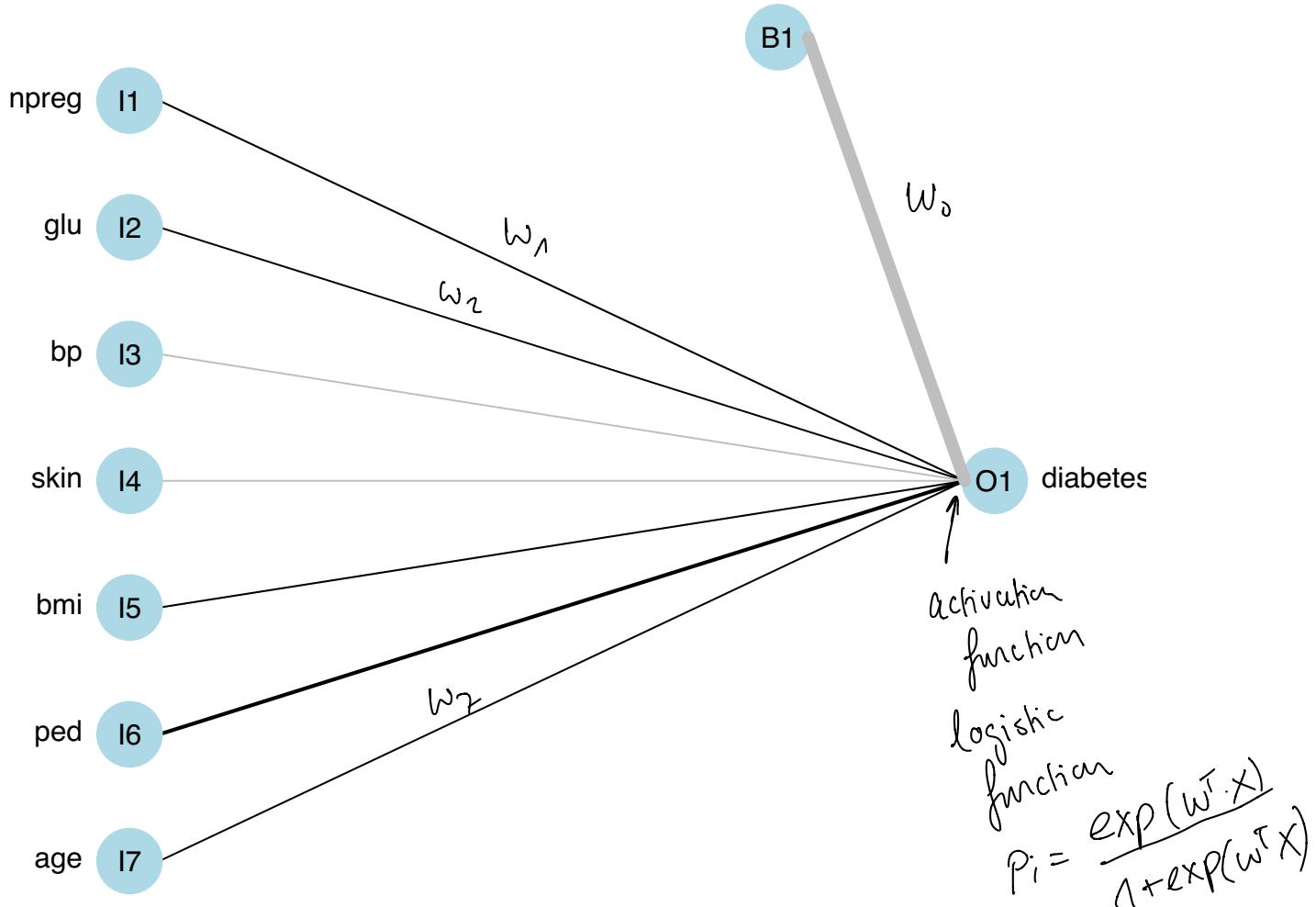
# entropy=TRUE because default is least squares
cbind(fitnnet$wts, fitlogist$coefficients)

```

	[,1]	[,2]
## (Intercept)	-9.773046277	-9.773061533
## npreg	0.103183171	0.103183427
## glu	0.032116832	0.032116823
## bp	-0.004767678	-0.004767542
## skin	-0.001917105	-0.001916632
## bmi	0.083624151	0.083623912
## ped	1.820397792	1.820410367
## age	0.041183744	0.041183529

By setting `entropy=TRUE` we minimize the cross-entropy loss.

plotnet(fitnnnet)



But, there may also exist local minima.

```
set.seed(123)
fitnnet = nnet(diabetes ~ npreg + glu + bp + skin + bmi + ped + age,
  data = train, linout = FALSE, size = 0, skip = TRUE, maxit = 10000,
  entropy = TRUE, Wts = fitlogist$coefficients + rnorm(8, 0, 1))
  ↑   ↑
## # weights:  8
## initial  value 24315.298582
## final    value 12526.062906
## converged
cbind(fitnnet$wts, fitlogist$coefficients)

##          [,1]      [,2]
## (Intercept) -36.733537 -9.773061533
## npreg        -77.126994  0.103183427
## glu         -2984.409175  0.032116823
## bp          -1835.934259 -0.004767542
## skin         -718.072629 -0.001916632
## bmi         -818.561311  0.083623912
## ped          -8.687473  1.820410367
## age         -773.023878  0.041183529
```

only

Why can NN and logistic regression lead to such different results?

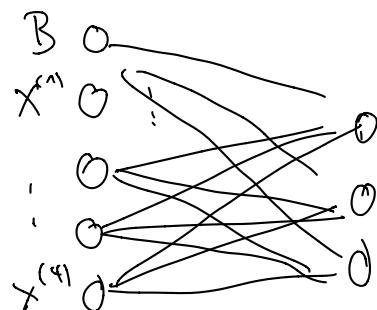
Example 3: Categorical outcome (multiclass regression)

The `iris` flower data set was introduced by the British statistician and biologist Ronald Fisher in 1936.

- Three plant species: {`setosa`, `virginica`, `versicolor`}.
- Four features: Sepal.Length, Sepal.Width, Petal.Length and Petal.Width.

$$\begin{array}{c} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \hline x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ x^{(4)} \end{array}$$

The aim is to predict the species of an iris plant.



$$\begin{aligned} P(y_i^{(1)}=1) + P(y_i^{(2)}=1) + P(y_i^{(3)}=1) \\ = 1 \end{aligned}$$

5 - 3 weights

The statistical model

We only briefly mentioned multiclass regression in module 4.

- Assume we have independent observation pairs (x_i, Y_i) , where the covariate vector x_i consists of the same measurements for each response category.
- Each observation can only belong to one response class, $Y_i \in \{1, \dots, C\}$.
- *Dummy variable coding* of the response in a C -dimensional vector:
 $y_i = (0, 0, \dots, 0, 1, 0, \dots, 0)$ with a value of 1 in the c^{th} element of y_i if the class is c .

$$\begin{array}{c} \nearrow \\ \begin{matrix} & 3 & 2 & 1 \\ \left(\begin{matrix} 0 \\ 0 \\ \vdots \end{matrix} \right) & \left(\begin{matrix} 0 \\ 1 \\ 0 \end{matrix} \right) & \left(\begin{matrix} 1 \\ 0 \\ 0 \end{matrix} \right) \\ \brace{ } & \brace{ } & \brace{ } \end{matrix} \end{array}$$

- Probabilities that the response is category c for subject i

$$\underline{p_{ic}} = P(Y_i = c) ,$$

where $\sum_{c=1}^C p_{ic} = 1$. In statistics we do not model p_{i1} , because $p_{i1} = 1 - \sum_{c=2}^C p_{ic}$, and $\beta_1 = 0$.

- Generalization of the logistic regression model:

$$p_{ic} = P(Y_i = c) = \frac{\exp(x_i^T \beta_c)}{1 + \sum_{s=2}^C \exp(x_i^T \beta_s)}$$

- Classification to the class with the highest probability,
 $\underset{c}{\operatorname{argmax}}(p_{ic})$.

Parameter estimation in the statistical model

- The likelihood of the multinomial regression model can be written as

$$\ln(L(\beta)) \propto \sum_{i=1}^n \sum_{c=1}^C y_{ic} \ln(p_{ic}) ,$$

where $p_{iC} = 1 - p_{i1} - p_{i2} - \dots - p_{i,C-1}$, and the regression parameters enter via the p_{ic} s.

Not
really
useful

- Parameter estimation is done in the same way as for the logistic regression, with the Fisher scoring algorithm.
- However (and this might be confusing), an efficient function in R also relies on neural networks for optimization (see below).

Neural network architecture and activation function

- Builds an output layer with C nodes and corresponding 0/1 targets (responses) using the dummy variable coding of the responses, called one-hot coding.
- The activation function for the output layer is called softmax. For each class $c = 1, \dots, C$ it is given as

$$P_{ic} = \hat{y}_c(x_i) = \frac{\exp(x_i^T w_c)}{\sum_{s=1}^C \exp(x_i^T w_s)}, \quad \in (0, 1)$$

weight of category c
sum of weight of all categories

where each w_s is a $r + 1$ dimensional vector of weights.

- Note: there is some redundancy here, since $\sum_{c=1}^C \hat{y}_c(x_i) = 1$, so we could have had $C - 1$ output nodes, but this is not done.
- The focus of neural networks is not to interpret the weights, and there is no need to assume full rank of a matrix with output nodes.

Neural networks: loss function and gradient descent

- For neural networks the negative of the multinomial loglikelihood is a scaled version of the *categorical cross-entropy loss*

$$J(w) = -\frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{C} \sum_{c=1}^C (y_{ic} \ln(\hat{y}_c(x_i)))}_{\text{log-likelihood}} .$$

- The optimization is done using gradient descent, with minor changes from what was done for the logistic regression due to the added sum and the change in the activation function.

Fitting multinomial regression vs a neural network

First select a training sample

```
library(nnet)
set.seed(123)
train = sample(1:150, 50)
iris_train = ird[train, ]
iris_test = ird[-train, ]
```

Then fit the `nnet()` (by default using the softmax activation function)

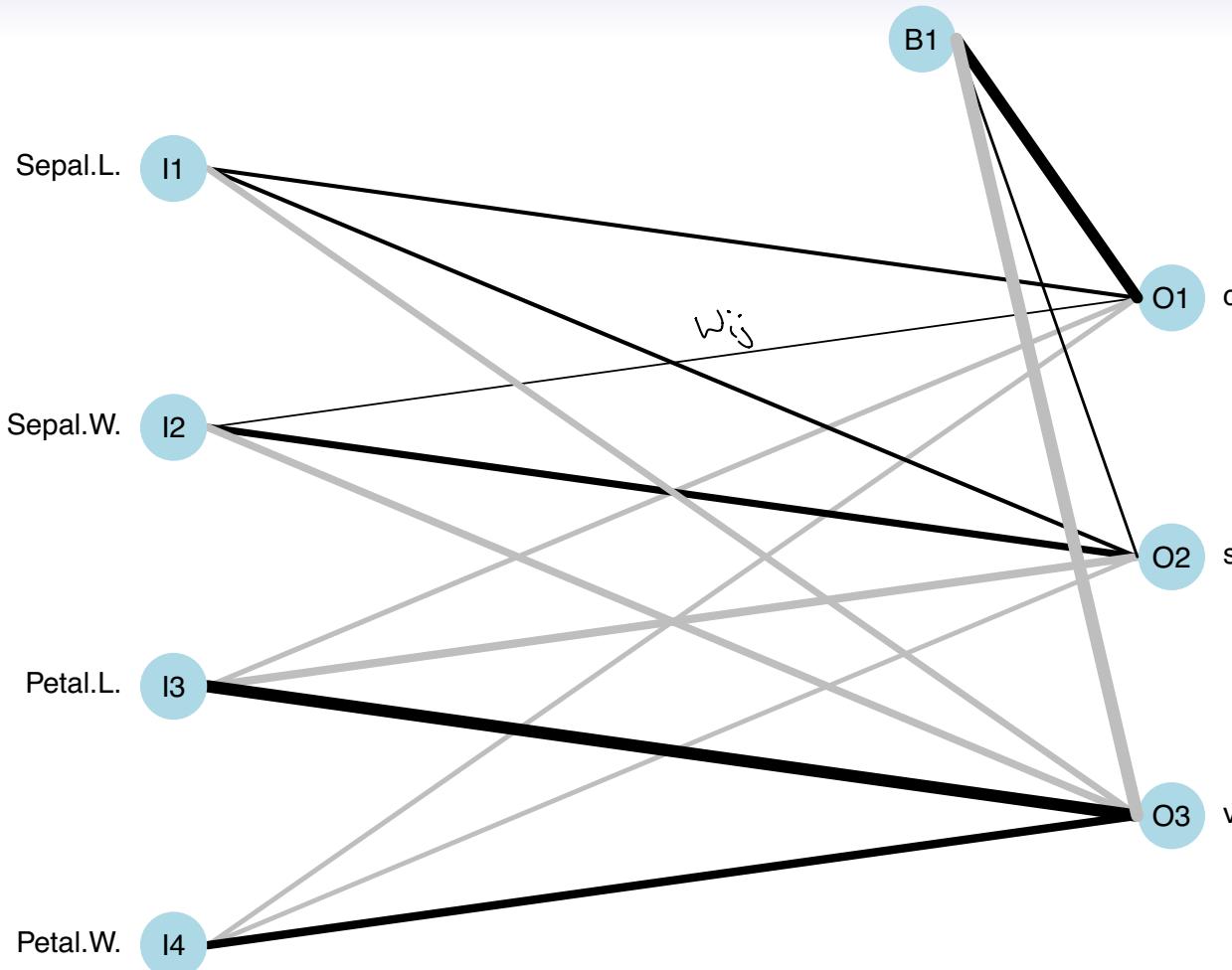
```
set.seed(1234)
iris.nnet <- nnet(species ~ ., data = ird, subset = train, size = 0,
                  skip = TRUE, maxit = 100)
```

```
## # weights:  15
## initial  value 105.595764
## iter  10 value 1.050064
## iter  20 value 0.018814
## iter  30 value 0.003937
## iter  40 value 0.002062
## iter  50 value 0.001460
## iter  60 value 0.000150
## iter  70 value 0.000125
## iter  80 value 0.000110
## final  value 0.000096
## converged
```

- How many weights (parameters) have been estimated?
- What does the graph look like?

```
summary(iris.nnet)

## a 4-0-3 network with 15 weights
## options were - skip-layer connections softmax modelling
## b->o1 i1->o1 i2->o1 i3->o1 i4->o1
## 36.79  9.69  1.26 -14.33 -13.24
## b->o2 i1->o2 i2->o2 i3->o2 i4->o2
##  5.16  10.10  21.33 -25.47 -12.48
## b->o3 i1->o3 i2->o3 i3->o3 i4->o3
## -42.03 -20.25 -23.12  40.80  25.96
```



Multinomial regression is also done with nnet, but using a wrapper `multinom` (we use default settings, so results are not necessarily the same as above).

```
library(caret)
fit = multinom(species ~ ., family = multinomial, data = iris_train,
               trace = FALSE)
coef(fit)

##   (Intercept) Sepal.L. Sepal.W. Petal.L. Petal.W.
## s  -15.67359 -8.679545 42.40492 -25.94868 -7.132379
## v  -157.32804 -19.694488 -17.03445  46.05297 59.135341
fit$wts

## [1] 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000
## [7] 0.000000 -15.673595 -8.679545 42.404924 -25.948681 -7.132379
## [13] 0.000000 -157.328044 -19.694488 -17.034447  46.052968 59.135341
```

↑
For identifiability reasons!

Summing up

1. *Multiple linear regression*

- NN with one input layer and one node in the output layer,
- linear activation function,
- mean squared loss.

2. *Logistic regression* (2 classes)

- NN with one input layer and one node in the output layer,
- sigmoid activation function,
- and binary cross-entropy loss.

3. *Multinomial regression* ($C > 2$ classes)

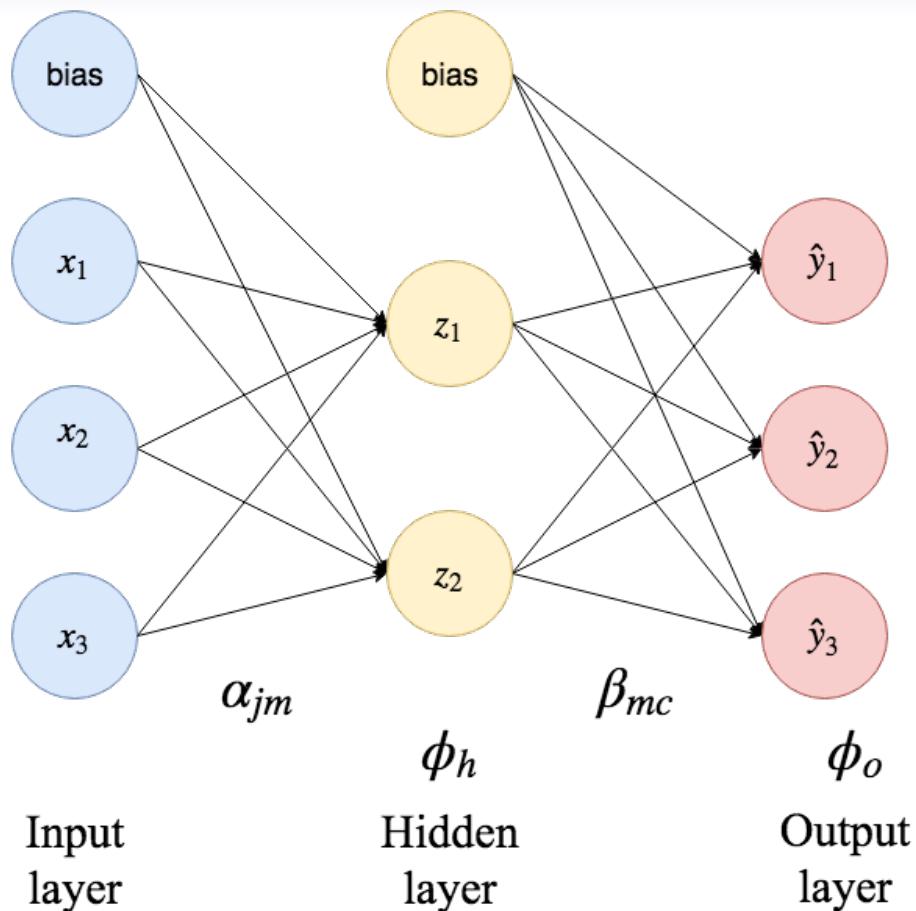
- NN with one input layer and C nodes in the output layer,
- softmax activation function,
- categorical cross-entropy loss.

But:

- These are only linear models (linear boundaries).
- Parameters (weights) found using gradient descent algorithms where the learning rate (step length) must be set.

Feedforward networks

- Connections are only forward in the network, but no feedback connections that sends the output of the model back into the network.
- Examples: Linear, logistic and multinomial regression with or without any *hidden layers* (between the input and output layers).
- We may have between zero and very many hidden layers.
- Adding *hidden layers* with *non-linear activation functions* between the input and output layer will make nonlinear statistical models.
- The number of hidden layers is called the *depth* of the network, and the number of nodes in a layer is called the *width* of the layer.



$$4 \cdot 2 + 3 \cdot 3 = 17$$

The single hidden layer feedforward network

The nodes are also called *neurons*.

Notation

1. Inputs (input layer nodes), $j = 1, \dots, p$: x_1, x_2, \dots, x_p , or as a vector x .
2. The nodes z_m in the hidden layer, $m = 1, \dots, M$; as vector $z^\top = (z_1, \dots, z_M)$, and the hidden layer activation function ϕ_h .

$$z_m(x) = \underbrace{\phi_h}_{\text{activation function}} \left(\underbrace{\alpha_{0m} + \sum_{j=1}^p \alpha_{jm} x_j}_{\text{linear combination of inputs}} \right)$$

where α_{jm} is the weight³ from input j to hidden node m , and α_{0m} is the bias term for the m th hidden node.

³We stick with greek letters α and β for parameters, but call them weights.

3. The node(s) in the output layer, $c = 1, \dots, C$: y_1, y_2, \dots, y_C , or as vector y , and output layer activation function ϕ_o .

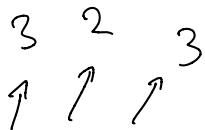
$$\hat{y}_c(x) = \phi_o(\beta_{0c} + \sum_{m=1}^M \underbrace{\beta_{mc} z_m(x)}_{\text{from hidden neuron } m \text{ to output node } c})$$

where β_{mc} is from hidden neuron m to output node c , and β_{0c} is the bias term for the c th output node.

4. Taken together

$$\hat{y}_c(x) = \phi_o(\beta_{0c} + \sum_{m=1}^M \beta_{mc} z_m) = \phi_o(\beta_{0c} + \sum_{m=1}^M \underbrace{\beta_{mc} \phi_h(\alpha_{0m} + \sum_{j=1}^p \alpha_{jm} x_j)}_{\text{from input } x \text{ through hidden layer}}))$$

Hands on:



- Identify p, M, C in the network figure above, and relate that to the $y_c(x)$ equation.
- How many parameters need to be estimated for this network?

$$4 \cdot 2 + 3 \cdot 3 = 17$$

- What determines the values of $\underbrace{p}_{\text{Data}}$ and $\underbrace{C}_{\text{Data}}$?

↳ Data

- How is M determined?



User - determined

Special case: linear activation function for the hidden layer

If we assume that $\phi_h(z) = z$ (linear or identity activation):

$$\hat{y}_c(x) = \underbrace{\phi_o}_{\text{output function}} \left(\beta_{0c} + \sum_{m=1}^M \beta_{mc} \left(\alpha_{0m} + \sum_{j=1}^p \alpha_{jm} x_j \right) \right)$$

Q: Does this look like something you have seen before?

A: Module 6, 10 ... PC - regression

Universal approximation property

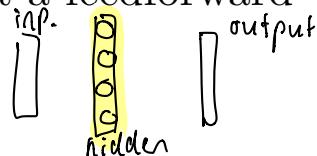
- Think of the goal of a feedforward network to approximate some function f , mapping our input vector x to an output value y .
- What type of mathematical function can a feedforward neural network with one hidden layer and linear output activation represent?

⁴Goodfellow et al 2016, Section 6.4.1, <https://www.deeplearningbook.org>

Universal approximation property

- Think of the goal of a feedforward network to approximate some function f , mapping our input vector x to an output value y .
- What type of mathematical function can a feedforward neural network with one hidden layer and linear output activation represent?

The *universal approximation theorem*⁴ says that a feedforward network with



- a *linear output layer*
- at least one hidden layer with a “squashing” activation function and “enough” hidden units

can approximate any (Borel measurable) function from one finite-dimensional space (our input layer) to another (our output layer) with any desired non-zero amount of error.

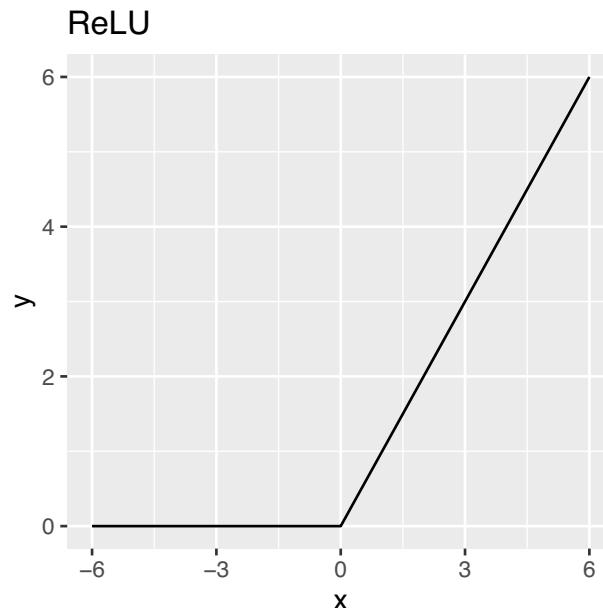
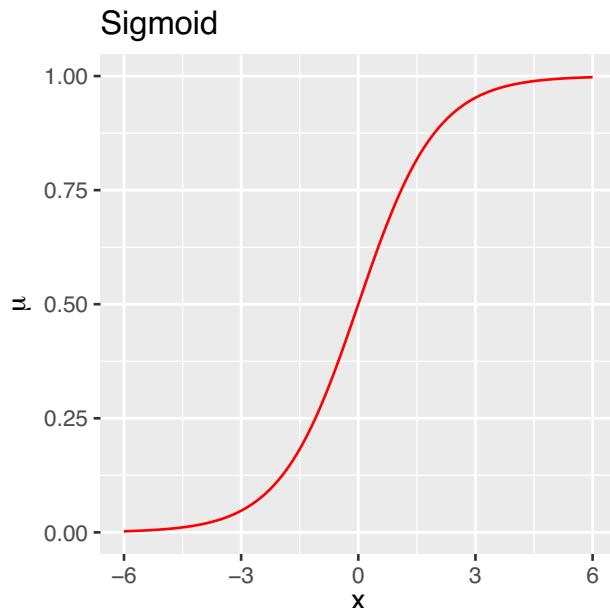
⁴Goodfellow et al 2016, Section 6.4.1, <https://www.deeplearningbook.org>

"Squashing functions"

In particular, the universal approximation theorem holds for

- The **sigmoid** $\phi(a) = 1/(1 + \exp(-a))$ (logistic) activation functions.
- The **rectified linear unit (ReLU)** $\phi_h(a) = \max(0, a)$ activation functions

in the hidden layer.



- The ReLU activation function has replaced the sigmoid function as default in the hidden layer(s) of a feedforward network.
- Even though a large feedforward network with one hidden layer may be able to *represent* a desired function, we may not be able to *learn* (i.e., estimate the) parameters of the function. Why?
 - wrong number of nodes in the hidden layer.
 - our optimization routine may fail / may find a local minimum.
 - we may overfit / underfit the training data.
- Alternative: networks with more than one hidden layer, but fewer total number of nodes but more layers. A network with *many hidden layers* is called a *deep network*.

*only 1
hidden layer*

way to go

The **nnet** and **keras** R packages

- We will use both the rather simple **nnet** R package by Brian Ripley and the currently very popular **keras** package for deep learning (the **keras** package will be presented later).
- **nnet** fits *one hidden layer* with *sigmoid activation function*. The implementation is not gradient descent, but instead BFGS using **optim**.
- Type `?nnet()` into your R-console to see the arguments of **nnet()**.
- If the response in formula is a factor, an appropriate classification network is constructed; this has one output, sigmoid activation and binary entropy loss for a binary response, and a number of outputs equal to the number of classes, softmax activation and categorical cross-entropy loss for more levels.

An example

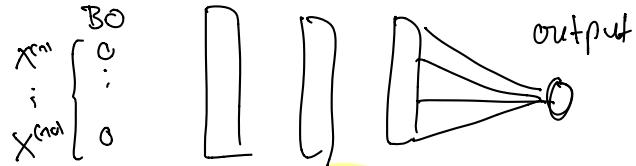
continuous output:
house price

Boston house prices

Objective: To predict the median price of owner-occupied homes in a given Boston suburb in the mid-1970s using 10 input variables.
This data set is both available in the MASS and keras R package.

Preparing the data

- Only 506, split between 404 training samples and 102 test samples (this split already done in the keras library)
- Each feature in the input data (for example, the crime rate) has a different scale, some values are proportions, which take values between 0 and 1; others take values between 1 and 12, others between 0 and 100, and so on.



```
library(keras)
dataset <- dataset_boston_housing()
c(c(train_data, train_targets), c(test_data, test_targets)) %<-% dataset
str(train_targets)

## num [1:404(1d)] 15.2 42.3 50 21.1 17.7 18.5 11.3 15.6 15.6 14.4 ...
head(train_data)

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12]
## [1,] 1.23247 0.0 8.14 0 0.538 6.142 91.7 3.9769 4 307 21.0 396.90
## [2,] 0.02177 82.5 2.03 0 0.415 7.610 15.7 6.2700 2 348 14.7 395.38
## [3,] 4.89822 0.0 18.10 0 0.631 4.970 100.0 1.3325 24 666 20.2 375.52
## [4,] 0.03961 0.0 5.19 0 0.515 6.037 34.5 5.9853 5 224 20.2 396.90
## [5,] 3.69311 0.0 18.10 0 0.713 6.376 88.4 2.5671 24 666 20.2 391.43
## [6,] 0.28392 0.0 7.38 0 0.493 5.708 74.3 4.7211 5 287 19.6 391.13
##      [,13]
## [1,] 18.72
## [2,] 3.11
## [3,] 3.26
## [4,] 8.01
## [5,] 14.65
## [6,] 11.74
```

The column names are missing (we could get them by using the Boston dataset loaded from the MASS library, but they are not relevant here).

- To make the optimization easier with gradient based methods do *feature-wise normalization*.

```
org_train = train_data  
mean <- apply(train_data, 2, mean)  
std <- apply(train_data, 2, sd)  
train_data <- scale(train_data, center = mean, scale = std)  
test_data <- scale(test_data, center = mean, scale = std)
```

From the training data } scaling by the some mean and sd

- **Note:** the quantities used for normalizing the test data are computed using the training data. You should never use in your workflow any quantity computed on the test data, even for something as simple as data normalization.

Just checking out one hidden layer with 5 units to get going.

```
library(nnet)
fit5 <- nnet(train_targets ~ ., data = train_data, size = 5, linout = TRUE,
    maxit = 1000, trace = F)

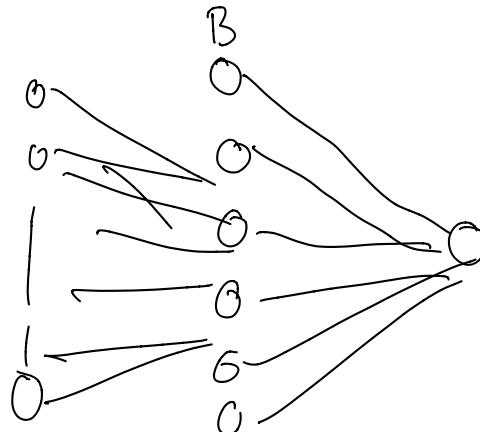
pred = predict(fit5, newdata = test_data, type = "raw")
sqrt(mean((pred[, 1] - test_targets)^2))

## [1] 6.167934 MSE
mean(abs(pred[, 1] - test_targets))

## [1] 4.271899
```

How do we know?

plotnet(fit5)



How to find best number of hidden nodes?

→ Cross-validation! (Time-consuming)

Here we compare to a model with 50 nodes:

```
library(nnet)
fit50 <- nnet(train_targets ~ ., data = train_data, size = 50, linout = TRUE,
               maxit = 5000, trace = F)
# head(summary(fit50))
pred = predict(fit50, newdata = test_data, type = "raw")
sqrt(mean((pred[, 1] - test_targets)^2))

## [1] 5.885913
mae = mean(abs(pred[, 1] - test_targets))
mae
## [1] 4.333412
```

Note: Deeper networks would presumably improve the error.

See here: https://www.math.ntnu.no/emner/TMA4268/2018v/11NN/11-neural_networks_boston_housing.html for running two hidden layers of 64 nodes each with keras. (This gave a mean absolute error on the test set of 2.68 after early stopping.)

Neural network parts

We now focus on the different elements of neural networks.

- 1) Output layer activation
- 2) Hidden layer activation
- 3) Network architecture
- 4) Loss function
- 5) Optimizers

1) Output layer activation → See yesterday

These choices have been guided by solutions in statistics (multiple linear regression, logistic regression, multiclass regression)

- *Linear activation*: for *continuous outcome* (regression problems)
- *Sigmoid activation*: for *binary outcome* (two-class classification problems)
- *Softmax*: for *multinomial/categorical outcome* (multi-class classification problems)

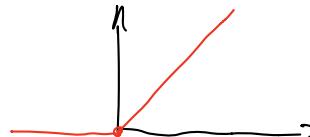
Remark: it is important that the output activation is matched with an appropriate **loss function** (see 4).

cross-entropy
loss

2) Hidden layer activation

(See chapter 6.3 in Goodfellow, Bengio, and Courville (2016))

Very common:

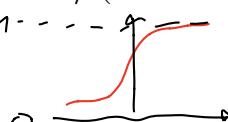


→ Choice nowadays!

- ReLU ($\phi_h(a) = \max(0, a)$) is standard choice for deep networks (many hidden layers and many nodes) today.
- Sigmoid ($\phi_h(a) = \sigma(a) = 1/(1 + \exp(-a))$) or hyperbolic tangent ($\phi_h(a) = \tanh(z)$).

A graph of the sigmoid function. The horizontal axis is labeled with a red arrow pointing right, and the vertical axis has a red arrow pointing up. The curve is an S-shape that passes through the point (0, 0.5). It approaches the horizontal asymptotes at y=0 and y=1 as x goes to negative infinity and positive infinity respectively.

Less common:



- Radial basis functions: as we looked at in Module 9.

- Softplus: $\phi_h(a) = \ln(1 + \exp(a))$



- Hard tanh: $\phi_h(a) = \max(-1, \min(1, a))$



Among all the possibilities, ReLU is nowadays the most popular one.
Why?

- The function is piecewise linear, but *in total non-linear*.
- Replacing sigmoid with ReLU is reported to be one of the major changes that have improved the performance of the feedforward networks⁵.
- Easy to use with gradient descent – even though the function is not differentiable at 0. As we will touch upon later, we don't expect to train a network until the gradient is 0.⁶.

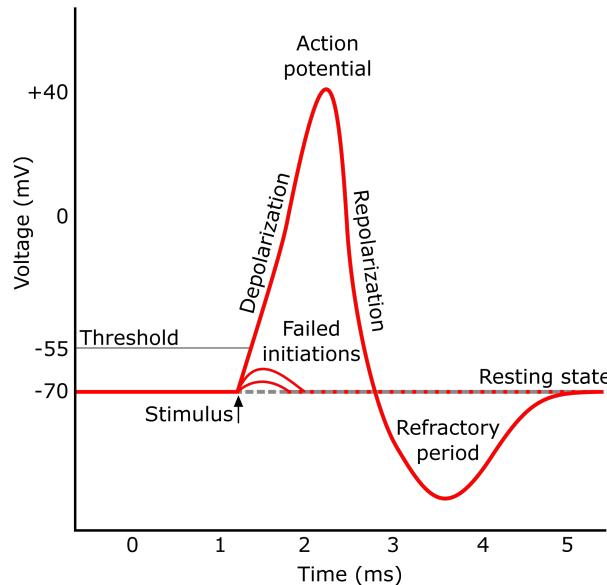
⁵Goodfellow et al, Section 6.6

⁶See Goodfellow et al, Section 6.3 for a discussion on this topic.

ReLU can also be motivated from biology.

- For some inputs a biological neuron can be completely inactive
- For some inputs a biological neuron output can be proportional to the input
- But, most of the time a biological neuron is inactive.

According to Goodfellow, Bengio, and Courville (2016) (Section 6.3), hidden unit design is an *active area of research*.



https://commons.wikimedia.org/wiki/File:Action_potential.svg

Q: Why can we not just use linear activation function in all hidden layers?

A: Then each layer would only be able to do linear transformations of the input data and a deep stack of linear layers would still implement a linear operation. The activation functions sigmoid and relu add non-linearity to the model. And, the **universal approximation** property is dependent on a squashing type activation function.

3) Network architecture

Network architecture contains three components:

- *Width*: How many nodes are in each layer of the network?
- *Depth*: How deep is the network (how many hidden layers)?
- *Connectivity*: How are the nodes connected to each other?

This depends on the problem, and here experience is important.

- We will only consider *feedforward networks*, where all nodes in one layer are connected to all the nodes in the next layer. The layers are then *fully connected* and *dense*.

However, the recent practice is to

Similarity: trees
Start by building a too big tree
Then: Pruning \leadsto regularization

- choose a too large network (too many nodes and/or too many layers) so that if trained until convergence (optimum) then this would result in overfitting, and
- then use other means to avoid this (various variants of regularization and hyperparameter optimization).

This simplifies the choice of network architecture to *choose a large enough network.*

See e.g. Chollet and Allaire (2018), Section 4.5.6/7 and Goodfellow, Bengio, and Courville (2016), Section 7

4) Loss function (“Method”)

- The choice of the loss function is closely related to the output layer activation function.
- Most popular problem types, output activation and loss functions:

Problem	Output activation	Loss function
Regression	linear	$\text{mse} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
Classification (C=2)	sigmoid	binary_crossentropy
Classification (C>2)	softmax	categorical_crossentropy

The loss functions chosen “need” to be (see below):

- differentiable
- possible to compute for each single training data point (or a mini-batch – to be explained soon)

5) Optimizers

Let the unknown parameters be denoted θ (what we have previously denotes as α s and β s), and the loss function to be minimized $J(\theta)$.

- Gradient descent
 - Mini-batch stochastic gradient descent (SGD) and true SGD
 - Backpropagation
- all input-output pairs to calculate $J(\theta)$*

Gradient descent

Remember:

Given the gradient $\nabla J(\theta^{(t)})^7$ of the loss function evaluated at the current estimate $\theta^{(t)}$, then the algorithm estimates the parameter at the next step as:

$$\theta^{(t+1)} = \theta^{(t)} - \lambda \nabla_{\theta} J(\theta^{(t)}) ,$$

where λ is the *learning rate* (usually a small value). In `keras` the default learning rate is 0.01.

- { Q: Why are we moving in the direction of the negative of the gradient? Why not the positive?
A:

⁷Remember that the gradient is the vector of partial derivatives of the loss function with respect to each of the parameter in the network.

- In *full gradient descent*, the loss function is computed as a mean over all training examples.

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{J(x_i, y_i)} .$$

- The gradient is *an average over many individual gradients* from the training example. You can think of this as an estimator for an expectation.

$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} J(x_i, y_i) .$$

- To build a network that generalizes well, it is important to have many training examples, but that would make us spend a lot of time and computer resources at calculating each gradient descent step.

Mini-batch stochastic gradient descent (SGD)

Crucial idea:

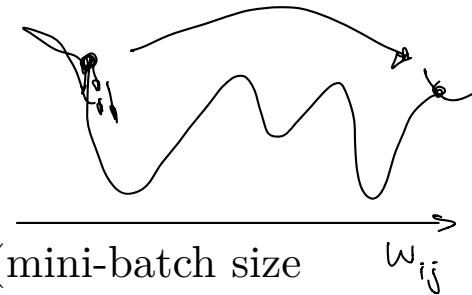
The expectation can be approximated by the average gradient over just a *mini-batch* (random sample) of the observations.

Advantages:

- The optimizer will converge much faster if it can rapidly compute approximate estimates of the gradient, instead of slowly computing the exact gradient (using all training data).
- Mini-batches may be processed *in parallel*, and the batch size is often a power of 2 (32 or 256).
- Small batches also bring in a *regularization effect*, maybe due to the variability they bring to the optimization process.

Mini-batch stochastic gradient descent

1. Divide all the training samples randomly into *mini-batches*.
2. Until convergence, repeat a) to d)
 - a) For each mini-batch: Make predictions of the responses in the mini-batch in a *forward pass*.
 - b) Compute the loss for the training data in this batch.
 - c) Compute the gradient of the loss with regard to the model's parameters (*backward pass*) based on the training data in the batch. $\nabla_{\theta}^* J(\theta^{(t)})$
 - d) Update all weights, but just using the *average gradient* from the mini-batch $\theta^{(t+1)} = \theta^{(t)} - \lambda \nabla_{\theta}^* J(\theta^{(t)})$
3. Network is *trained*; return parameter estimates.



Special case: *True SGD* involves only *one sample* (mini-batch size 1). → Mini-batch SGD is a compromise between SGD (one sample per iteration) and full gradient descent (full dataset per iteration)

In the 3rd video (on backpropagation) from 3Blue1Brown there is nice example of one trajectory from gradient decent and one from SGD (10:10 minutes into the video):

https://www.youtube.com/watch?v=Ilg3gGewQ5U&list=PLZHQBObOWTQDNU6R1_67000Dx_ZCJB-3pi&index=3

Backpropagation algorithm

- *Backpropagation* is a simple and inexpensive way *to calculate the gradient*.
- Computing the analytical expression for the gradient ∇J is not difficult, but *the numerical evaluation may be expensive*.
- The *chain rule* is used to compute derivatives of functions of other functions where the derivatives are known. This is efficiently done with backpropagation.
- Backpropagation starts with the value of the loss function and works backward from the top layers to the bottom layers, applying the chain rule to compute the contribution that each parameter has in the loss value.

More background:

- Mathematical details: Goodfellow, Bengio, and Courville (2016) Section 6.5.
- 3Blue1Brown videos:
<https://www.youtube.com/watch?v=Ilg3gGewQ5U> and
<https://www.youtube.com/watch?v=tIeHLnjs5U8>

Variations of SGD — with adaptive learning rates

The learning rate λ is difficult to set. Some ideas for adaptive learning rates (Goodfellow et al, Chapters 8.3 and 8.5).

- *Momentum term*: previous gradients are allowed to contribute.
- *AdaGrad*: individually adapt the learning rates (λ) of all model parameters. The parameters with the largest partial derivative of the loss (over all historical squared values) have a correspondingly rapid decrease in their learning rate.
- *RMSprop*: modification to AdaGrad in non-convex setting. Scales with exponentially weighted moving average instead of all historical squared gradient values. This helps to “forget” information from the extreme past.

Regularization

To avoid overfitting

- Goodfellow, Bengio, and Courville (2016), Chapter 7, define regularization as *any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.*
- Remember (module 6): The aim of regularization was to trade *increased bias for reduced variance*. The idea was to add a penalty to the loss function.
- The penalties we looked at were of type absolute value of parameter (L_1 , lasso, where we looked at this as model selection) and square value of parameter (L_2 , ridge regression). This can also be done for neural networks.

Regularization in neural networks

- *Weight decay*: In neural networks this means adding a L_2 -penalty to the loss function to *penalize large weights*⁸:

$$\tilde{J}(w) = \underbrace{\frac{\alpha}{2} w^\top w}_{\text{"Ridge idea"} + J(w).$$

Also : $L_1 : \alpha \cdot |w|$

- *Dataset augmentation*: Adding fake data to the dataset, in order that the trained model will generalize better. For some learning tasks it is straightforward to create fake data. For image data this can be done by rotating and scaling the images.
- *Label smoothing*: Motivated by the fact that the training data may contain errors in the responses recorded, and replaced the one-hot coding for C classes with $\epsilon/(C - 1)$ and $1 - \epsilon$ for some small ϵ .
 - ① $y = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0.05 \\ 0.05 \end{pmatrix}$
 - ② $y = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0.05 \\ 0.05 \end{pmatrix}$
 - ③ $y = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0.05 \\ 0.05 \end{pmatrix}$
- *Early stopping*
- *Dropout*

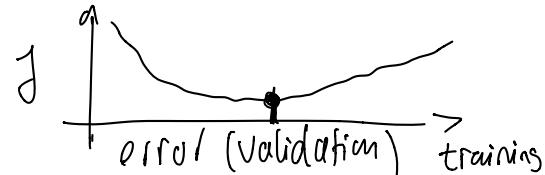
⁸see chapter 7.1.1 in Goodfellow et al

Early stopping

Based on Goodfellow, Bengio, and Courville (2016), Section 7.8

How many training iterations are needed?

- The most commonly used form of regularization is *early stopping*.
- If we have chosen a sufficiently large model with the capacity to overfit the training data, we would observe that the training error decreases steadily during training, but the error on the validation set at some point begins to increase.
- If we stop the learning early and return the parameters giving the test performance on the validation set, this model would hopefully be a better model than if we trained the model until convergence.
- It is possible to think of the number of *training steps* as a *hyperparameter*. This hyperparameter can easily be set, and the cost is just the performance on the validation set during training. Alternatively, cross-validation can be used.



Dropout

e.g. 20% : set output of nodes in
hidden layers = 0

Based on Goodfellow, Bengio, and Courville (2016), Section 7.12, and Chollet and Allaire (2018) 4.4.3

Dropout was developed by Geoff Hinton and his students.

- During training: randomly *dropout* (set to zero) some outputs in a given layer at each iteration. Drop-out rates may be chosen between 0.2 and 0.5.
- During test: no dropout, but scale down the layer output values by a factor equal to the drop-out rate (since now more units are active than we had during training).
- Alternatively, the drop-out and scaling (now upscaling) can be done during training.

Ways to avoid overfitting

Many **hyperparameters** when building and fitting a neural network, like the network architecture, the number of batches to run before terminating the optimization, the drop-out rate, etc.

To avoid overfit, we have some strategies:

- Reduce network size.
- Collect more observations.
- Regularization. → See previous slides

It is important that the hyperparameters are chosen on a validation set or by cross-validation.

However, we may run into *validation-set overfitting*: when using the validation set to decide many hyperparameters, so many that you may effectively overfit the validation set.

Deep learning

Deep Learning is an algorithm which has no theoretical limitations of what it can learn; the more data you give and the more computational time you provide, the better it is.

Geoffrey Hinton (Google)

Timeline

(based on Chollet and Allaire (2018))

- Neural networks were investigated in “toy form” in the 1950s.
- The first big step was taken in the 1980s when the backpropagation algorithm was rediscovered.
- In 1989 (Bell Labs, Yann LeCun) used convolutional neural networks to classifying handwritten digits.
- Not so much activity in the neural network field in the 2000s.

- In 2011 neural networks with many layers (and trained with GPUs) were performing well on image classification tasks.
- The *ImageNet* classification challenge (classify high resolution colour images into 1k different categories after training on 1.4M images) was won by solutions with deep convolutional neural networks (convnets). In 2011 the accuracy was 74.3%, in 2012 83.6% and in 2015 96.4%.
- From 2012, convnets is the general solution for computer vision tasks. Other application areas are natural language processing.

Deep?

- Deep learning does not mean a deeper understanding, but refers to successive layers of representations - where the number of layers gives the *depth* of the model. Often tens to hundreds of layers are used.
- Deep neural networks are not seen as models of the brain, and are not related to neurobiology.
- A deep network can be seen as many stages of *information-distillation*, where each stage performs a simple data transformation. These transformations are not curated by the data analyst, but is estimated in the network.
- In contrast, in statistics we first select a set of inputs, then look at how these inputs should be transformed, before we apply some statistical methods (*feature engineering*).

- The success of deep learning is dependent upon the breakthroughs
 - in *hardware* development, especially with faster CPUs and massively parallel graphical processing units (GPUs).
 - *datasets* and benchmarks (internet/tech data).
 - improvements of the *algorithms*.
- Achievements of deep learning includes high quality (near-human to super human) image classification, speech recognition, handwriting transcription, machine translation, digital assistants, autonomous driving, advertising targeting, web searches, playing Go and chess.

The R keras package

- Earlier good programming skills (especially in C++) was essential to work in deep learning.
- With the launch of the Keras library now users may only need basic skills in Python or R.
- **Keras** can be seen as a way to use LEGO bricks in deep learning.
To quote the web-page:

Keras is a high-level neural networks API developed with a focus on enabling fast experimentation. Being able to go from idea to result with the least possible delay is key to doing good research.

More information on the R solution: <https://keras.rstudio.com/>

Cheat-sheet for R Keras:

<https://github.com/rstudio/cheatsheets/raw/master/keras.pdf>

MNIST dataset

- This is a larger image version of the handwritten digits data set (a different version, ZIP-codes is found under Recommended exercises).
- This data analysis is based on
https://www.math.ntnu.no/emner/TMA4268/2018v/11NN/8-neural_networks_mnist.html and the R `keras` cheat sheet.

See problem 4 in RecEx 11
(Recommended!)

$^{28} \quad ^{28}$

Objective: classify the digit contained in an image (~~128~~ \times ~~128~~ greyscale).



Labels for the training data:

```
## train_labels
##    0    1    2    3    4    5    6    7    8    9
## 5923 6742 5958 6131 5842 5421 5918 6265 5851 5949
```

1. Training and test data

60 000 images for training and 10 000 images for testing.

```
# Training data
train_images <- mnist$train$x → str(train_images)
train_labels <- mnist$train$y ← label between 0 and 9

# Test data
test_images <- mnist$test$x
test_labels <- mnist$test$y
org_testlabels <- test_labels
```

The `train_images` is a tensor (generalization of a matrix) with 3 axis, (`samples, height, width`).

2. Defining the model

Using `keras_model_sequential()` we build a model with a stack of layers. `layer_dense()` adds densely connected layers on top of the input layer. Each sample contains $28 \times 28 = 784$ pixels (= input nodes) and 10 features (=output nodes). Adding a bias term (intercept) is default for `layer_dense`.

```
network <- keras_model_sequential() %>% layer_dense(units = 512, activation = "relu",
  input_shape = c(28 * 28)) %>% layer_dense(units = 10, activation = "softmax")
summary(network)

## Model: "sequential"
## -----
## Layer (type)          Output Shape         Param #
## =====
## dense (Dense)        (None, 512)           401920
## -----
## dense_1 (Dense)       (None, 10)            5130
## -----
## Total params: 407,050
## Trainable params: 407,050
## Non-trainable params: 0
## -----
```

As activation function we use `relu` for the hidden layer, and `softmax` for the output layer - since we have 10 classes (where one is correct each time).

3. Configure the learning process

We choose

- The loss function we will use, which is cross-entropy.
- The version of the optimizer, which is RMSprop.
- Which measure (metrics) do we want to monitor in our training phase? Here we choose accuracy (=percentage correctly classified).

```
network %>% compile(optimizer = "rmsprop", loss = "categorical_crossentropy",
  metrics = c("accuracy"))
```

4. Preprocessing to match with model inputs and outputs

- The training data is stored in an array of dimension (60000, 28, 28) of type integer with values in the [0, 255] interval.
- The data must be reshaped into the shape the network expects ($28 \cdot 28$). In addition the grey scale values are scaled to be in the [0, 1] interval.
- Also, the response must be transformed from 0-10 to a vector of 0s and 1s (dummy variable coding) aka *one-hot-coding*.

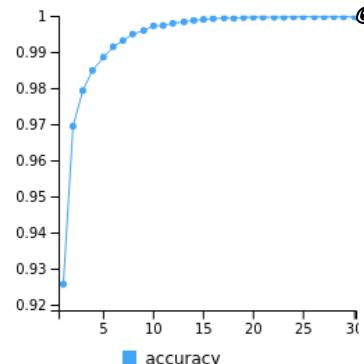
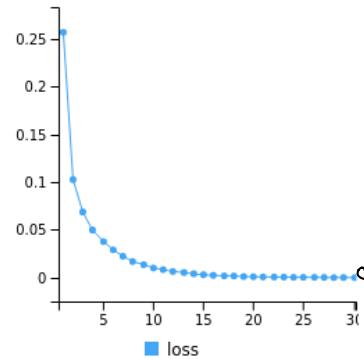
```
train_images <- array_reshape(train_images, c(60000, 28 * 28))
train_images <- train_images/255
train_labels <- to_categorical(train_labels)

test_images <- array_reshape(test_images, c(10000, 28 * 28))
test_images <- test_images/255
test_labels <- to_categorical(test_labels)
```

5. Fit the model

```
fitted<-network %>% fit(train_images, train_labels,  
  epochs = 30, batch_size = 128)  
library(ggplot2)  
plot(fitted)+ggtitle("Fitted model")
```

} slow part!



6. Evaluation and prediction

→ now it is trained

```
network %>% evaluate(test_images, test_labels)
testres = network %>% predict_classes(test_images)
# $loss [1] 0.1194063 $acc [1] 0.9827
confusionMatrix(factor(testres), factor(org_testlabels))
```

Confusion Matrix and Statistics

Reference											
Prediction	0	1	2	3	4	5	6	7	8	9	
0	971	0	3	0	1	2	5	1	1	1	
1	1	1128	2	0	0	0	2	2	2	3	
2	1	1	1009	3	3	0	2	7	2	0	
3	0	1	2	997	0	11	1	3	4	3	
4	1	0	2	0	969	1	4	2	3	7	
5	0	1	0	1	0	871	3	0	1	4	
6	3	2	2	0	3	4	940	0	1	0	
7	1	1	4	2	1	0	0	1002	2	3	
8	2	1	7	0	0	2	1	4	953	1	
9	0	0	1	7	5	1	0	7	5	987	

Overall Statistics

Accuracy : 0.9827
95% CI : (0.9799, 0.9852)

No Information Rate : 0.1135
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.9808

McNemar's Test P-Value : NA

Kaggle

- Kaggle has hosted machine-learning competitions since 2010, and by looking at solutions to competitions it is possible to get an overview of what works.
- In 2016-2017 gradient boosting methods won the competitions with structured data (“shallow” learning problems), while deeplearning won perceptual problems (as image classification).
- Kaggle has helped (and is helping) the rise in deep learning.
- Very timely topic: Analysis of COVID-19 data!

Other analyses

Boston housing price

taken from Chollet and Allaire (2018):

https://www.math.ntnu.no/emner/TMA4268/2018v/11NN/11-neural_networks_boston_housing.html

Movie data base

https://www.math.ntnu.no/emner/TMA4268/2018v/11NN/9-neural_networks_imdb.html

Reuters data

https://www.math.ntnu.no/emner/TMA4268/2018v/11NN/10-neural_networks_reuters.html

Summing up

- Feedforward network architecture: mathematical formula - layers of multivariate transformed (`relu`, `linear`, `sigmoid`) inner products - sequentially connected.
- What is the number of parameters that need to be estimated? Intercept term (for each layer) is possible and is referred to as “bias term”.
- Loss function to minimize (on output layer): regression (mean squared), classification binary (binary crossentropy), classification multiple classes (categorical crossentropy) — and remember to connect to the correct choice of output activation function: mean squared loss goes with linear activation, binary crossentropy with sigmoid, categorical crossentropy with softmax.
- How to minimize the loss function: gradient based (chain rule) back-propagation - many variants.
- Technicalities: `nnet` and `keras` in R; Piping sequential layers, piping to estimation and then to evaluation (metrics).

References and further reading

- <https://youtu.be/aircAruvnKk> from 3BLUE1BROWN - 4 videos
 - using the MNIST-data set as the running example
- Look at how the hidden layer behave:
<https://playground.tensorflow.org>
- Friedman, Hastie, and Tibshirani (2001), Chapter 11: Neural Networks
- Efron and Hastie (2016), Chapter 18: Neural Networks and Deep Learning
- Chollet and Allaire (2018)
- Goodfellow, Bengio, and Courville (2016) (used in IT3030)
<https://www.deeplearningbook.org/>
- Explaining backpropagation
<http://neuralnetworksanddeeplearning.com/chap2.html>
- Slides from MA8701 (Thiago Martins) <https://www.math.ntnu.no/emner/MA8701/2019v/DeepLearning/>

Acknowledgements

Chollet, François, and J. J. Allaire. 2018. *Deep Learning with R*. Manning Press.

<https://www.manning.com/books/deep-learning-with-r>.

Efron, Bradley, and Trevor Hastie. 2016. *Computer Age Statistical Inference - Algorithms, Evidence, and Data Science*. Cambridge University Press.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2001. *The Elements of Statistical Learning*. Vol. 1. Springer series in statistics New York.

Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.