



Deep Learning Training Course

Introduction to Deep Learning Models

Prof. Dr. – Ing. Morris Riedel

Adjunct Associated Professor

School of Engineering and Natural Sciences, University of Iceland

Research Group Leader, Juelich Supercomputing Centre, Germany

LECTURE 3

Neural Network

May 22th, 2019

Juelich Supercomputing Centre, Juelich, Germany



UNIVERSITY OF ICELAND
SCHOOL OF ENGINEERING AND NATURAL SCIENCES
FACULTY OF INDUSTRIAL ENGINEERING,
MECHANICAL ENGINEERING AND COMPUTER SCIENCE



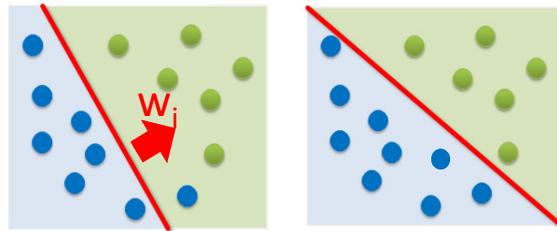
JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE



HELMHOLTZ
ARTIFICIAL INTELLIGENCE
COOPERATION UNIT

Review of Lecture 2 – Overview of Deep Learning



[2] F. Rosenblatt, 1957

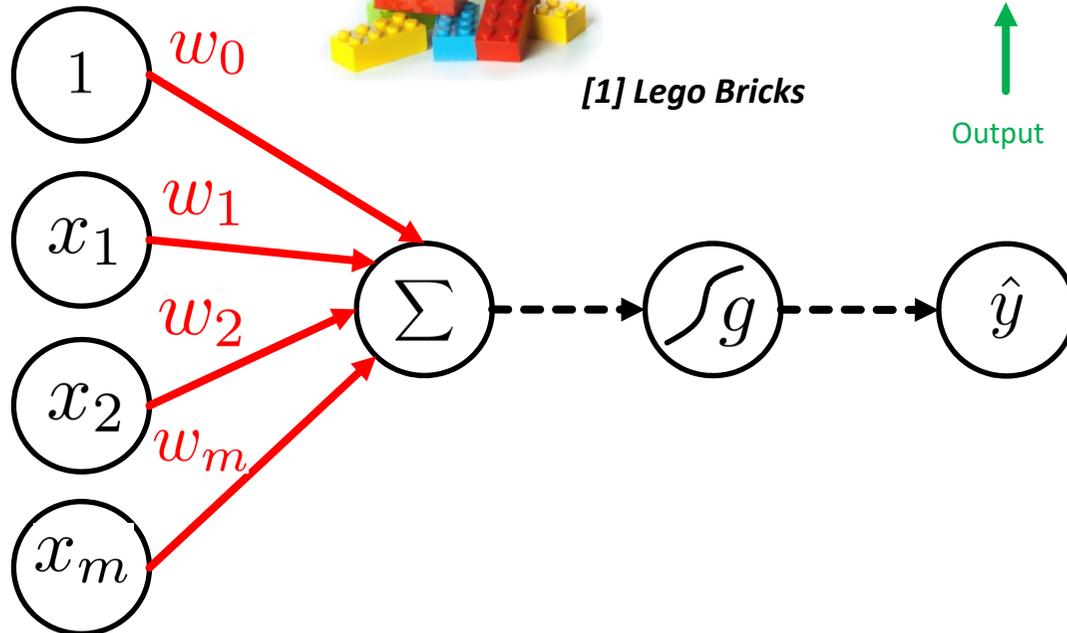


[3] Keras Library



[1] Lego Bricks

Constants



Input Data

Trainable Weights

Sum

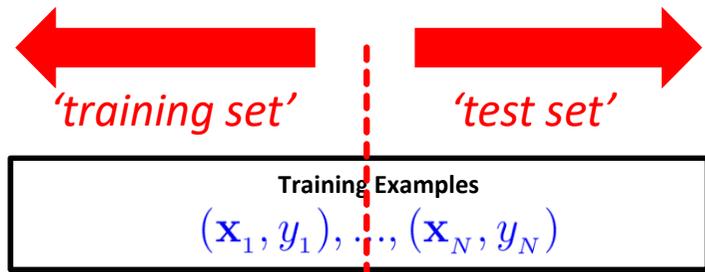
Non-Linearity

Output

$$\hat{y} = g \left(1 * w_0 + \sum_{i=1}^m x_i * w_i \right)$$

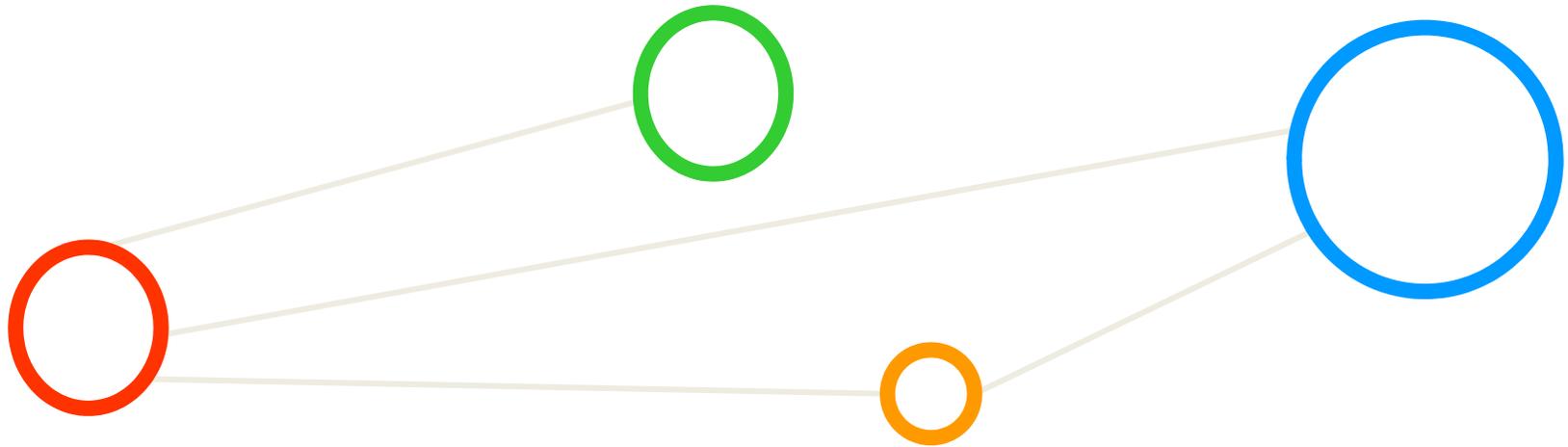
Annotations for the equation:

- ↑ Output (pointing to \hat{y})
- ↑ Bias (pointing to $1 * w_0$)
- ↑ Sum (pointing to the summation symbol \sum)
- ↑ linear combination of input data (pointing to $x_i * w_i$)
- ↓ non-linear activation function (pointing to g)



(historical records, groundtruth data, examples)

Outline



Outline of the Course

1. Deep Learning driven by HPC & Jupyter
2. Overview of Deep Learning
3. Neural Network
4. Convolutional Neural Network
5. Introduction to Deep Learning for Remote Sensing & 1D/2D CNNs for Hyperspectral Images Classification
6. 3D CNNs for Hyperspectral Images Classification, Training Set Selection and Performance Evaluation
7. Recurrent Neural Network
8. LSTM Neural Network
9. Deep Reinforcement Learning
10. Course Summary & Lessons Learned

Theoretical Lectures

Practical Lectures

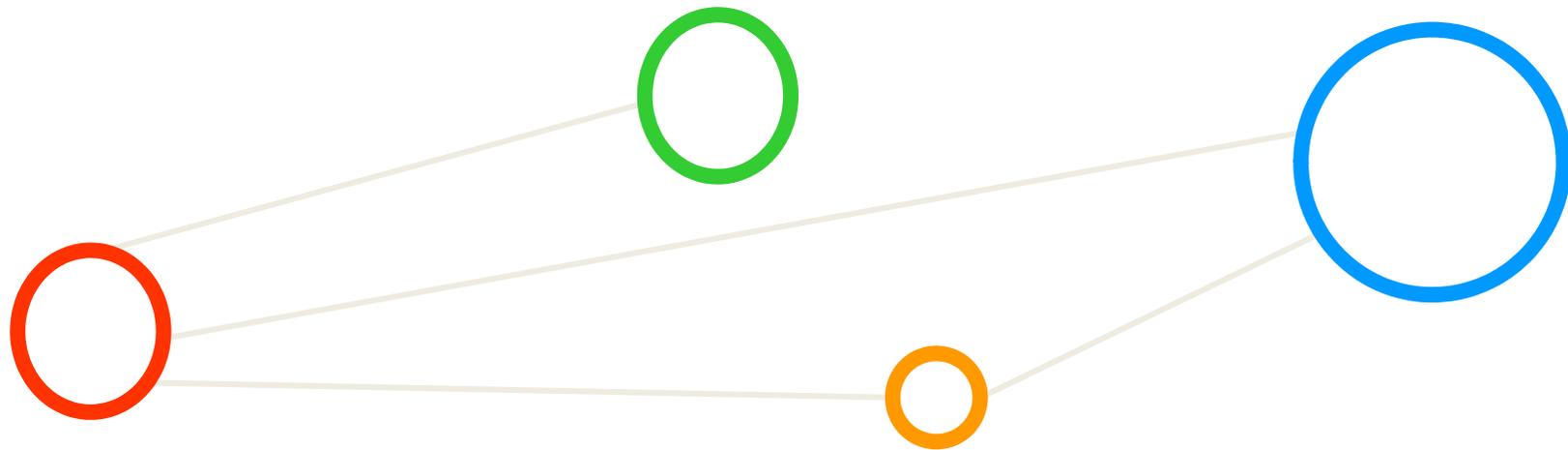


Outline

- Learning from Data with the Perceptron Model
 - Hand-written Character Recognition Problem
 - Data Exploration using Jupyter & NumPy
 - Multi-Class Classification Problem
 - Multi-Output Perceptron Model
 - Using Keras & TensorFlow in Jupyter
- Artificial Neural Networks (ANNs)
 - Creating ANN Network Topologies
 - Overfitting Reasoning & Validation
 - Validation Datasets & Splits
 - Many Parameters & Hidden Layers
 - Regularization Techniques & Examples



Learning from Data with the Perceptron Model

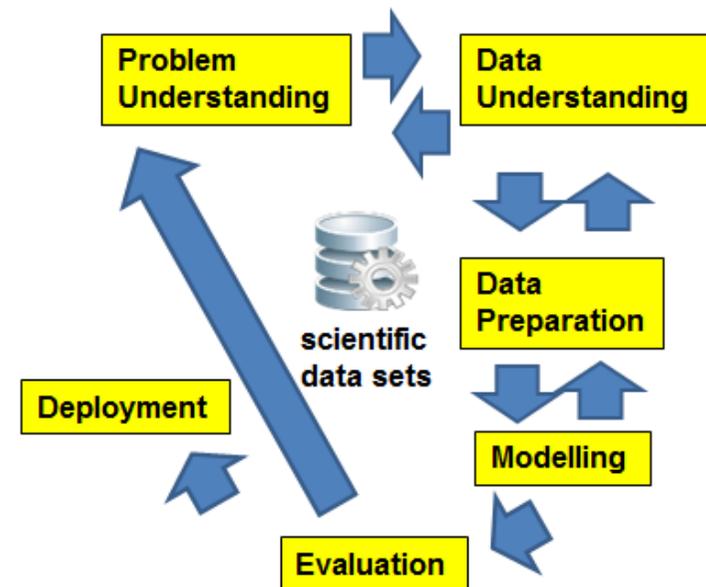


Systematic Process to Support Learning From Data

- Systematic data analysis guided by a ‘standard process’
 - Cross-Industry Standard Process for Data Mining (CRISP-DM)

- A data mining project is guided by these six phases:
 - (1) Problem Understanding;
 - (2) Data Understanding;
 - (3) Data Preparation;
 - (4) Modeling;
 - (5) Evaluation;
 - (6) Deployment

(learning takes place)



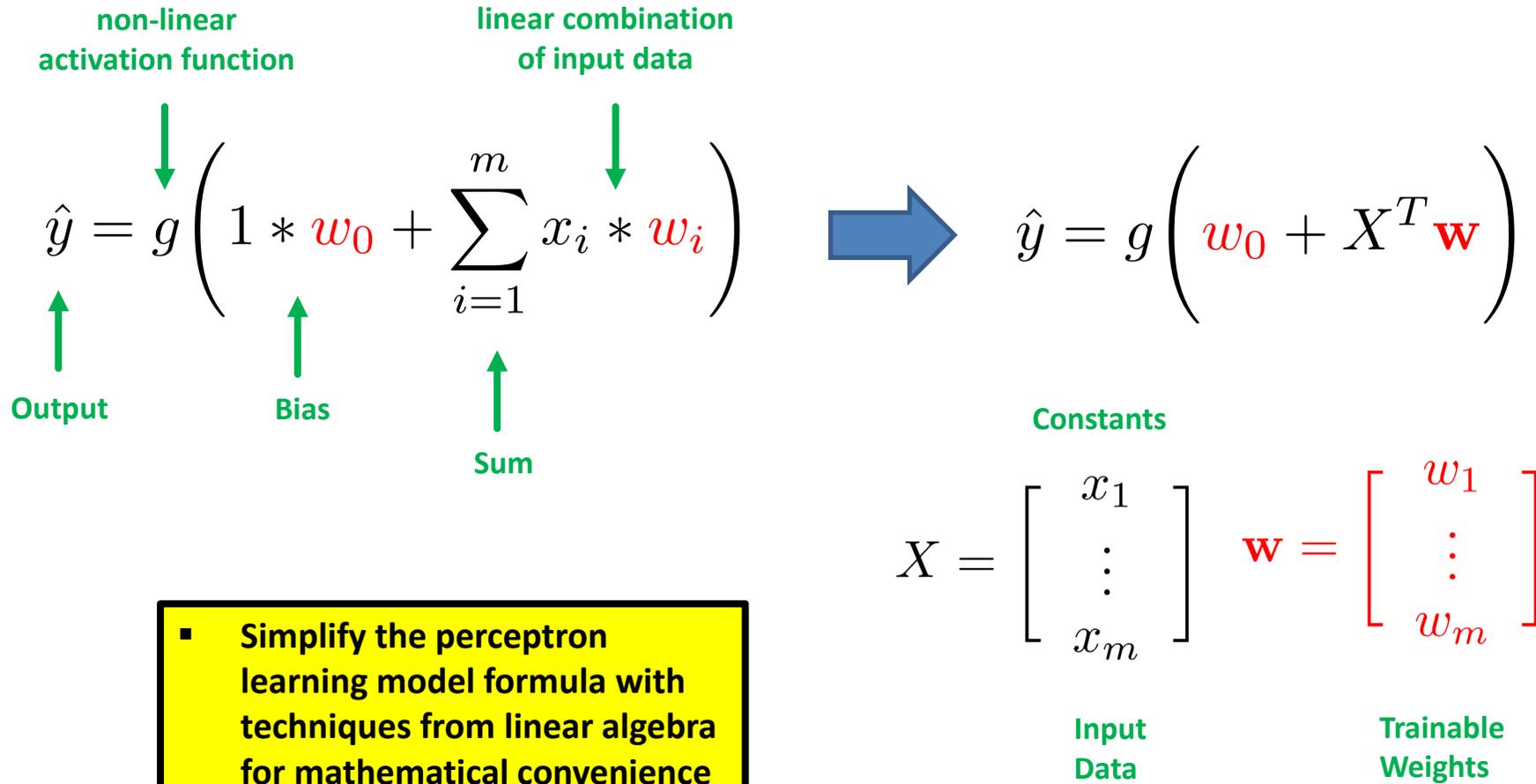
- Lessons Learned from Practice

- Go back and forth between the different six phases

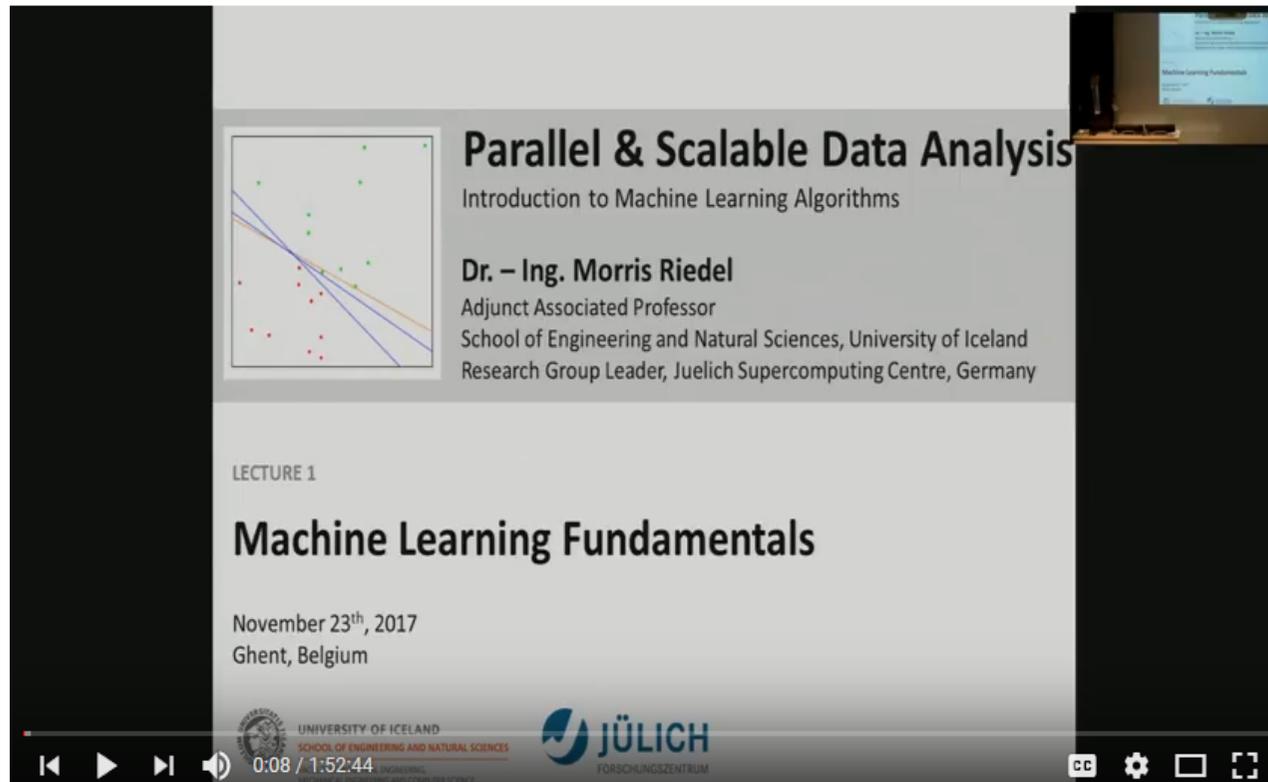
[7] C. Shearer, CRISP-DM model, Journal Data Warehousing, 5:13

➤ A more detailed description of all six CRISP-DM phases is in the Appendix A of the slideset

Perceptron Model – Mathematical Notation for one Neuron



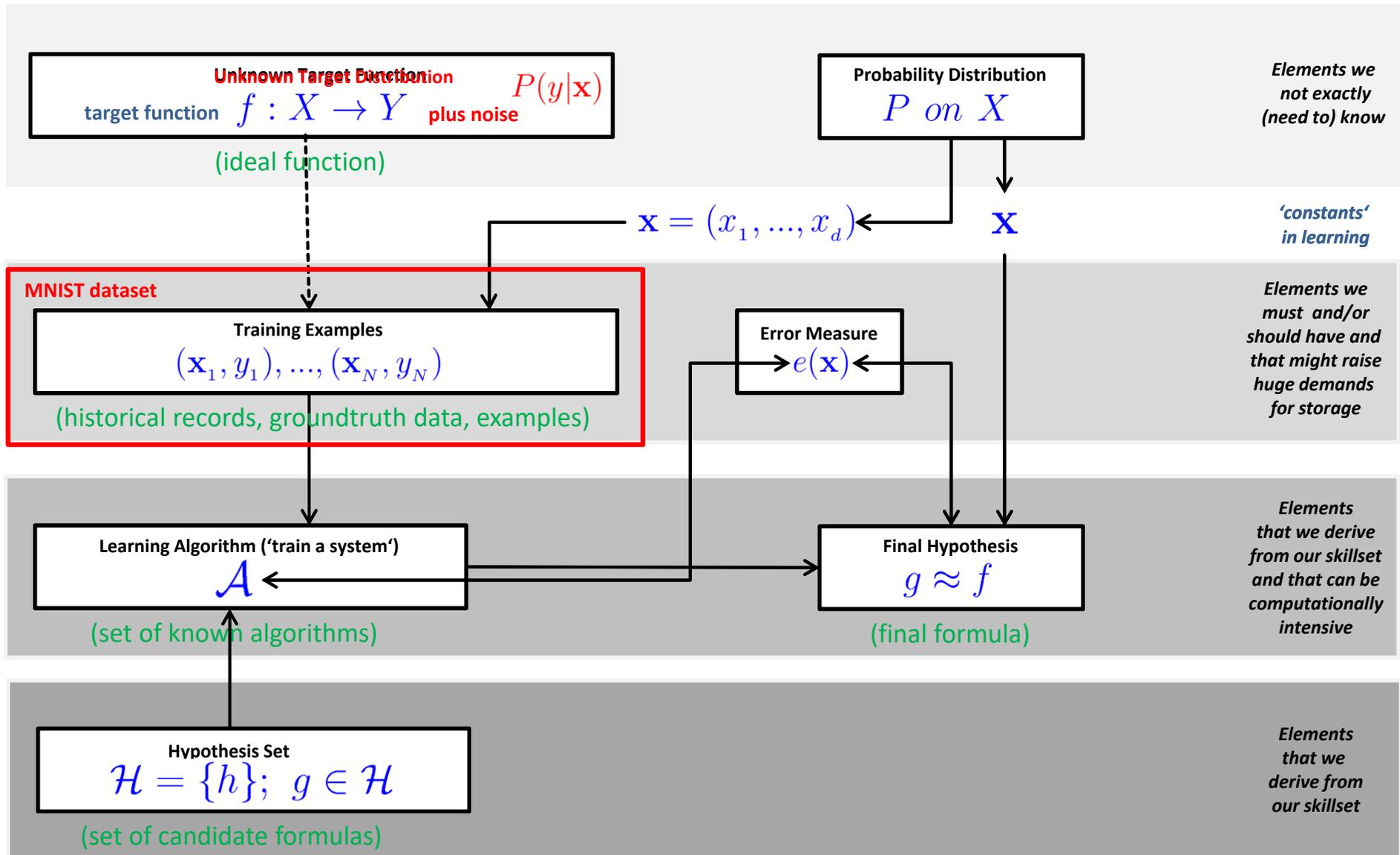
[YouTube Lectures] More Details about Perceptron Model



[8] Morris Riedel, 'Introduction to Machine Learning Algorithms', Invited YouTube Lecture, six lectures, University of Ghent, 2017

➤ Note that this course is not a full machine learning course but rather focusses on deep learning

Supervised Learning – MNIST Example Overview



Handwritten Character Recognition MNIST Dataset

(1) Problem Understanding Phase

- Metadata
 - Subset of a larger dataset from US National Institute of Standards (NIST)
 - Handwritten digits including corresponding labels with values 0 to 9
 - All digits have been size-normalized to 28 * 28 pixels and are centered in a fixed-size image for direct processing
 - Not very challenging dataset, but good for experiments / tutorials

Dataset Samples

(10 class classification problem)

- Labelled data (10 classes)
- Two separate files for training and test
- 60000 training samples (~47 MB)
- 10000 test samples (~7.8 MB)



(2) Data Understanding Phase

MNIST Dataset – Data Access

- When working with the dataset

(2) Data Understanding Phase

- Dataset is not in any standard image format like jpg, bmp, or gif (i.e. file format not known to a graphics viewer)
- Data samples are stored in a simple file format that is designed for storing vectors and multidimensional matrices (i.e. numpy arrays)
- The pixels of the handwritten digit images are organized row-wise with pixel values ranging from 0 (white background) to 255 (black foreground)
- Images contain grey levels as a result of an anti-aliasing technique used by the normalization algorithm that generated this dataset

- Available for the tutorial

- Easy download via Keras from an Amazon Web Services (AWS) cloud

(downloads data into `~home/.keras/datasets` as NPZ file format of numpy that provides storage of array data using gzip compression)

```
import numpy as np
from keras.datasets import mnist
```

```
# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```
[riedell@juron1-adm datasets]$ pwd
/p/home/jusers/riedell/juron/.keras/datasets
[riedell@juron1-adm datasets]$ ls -al
total 11234
drwxr-xr-x 2 riedell jusers 4096 Jan 20 22:05 .
drwxr-xr-x 3 riedell jusers 4096 Jan 20 22:03 ..
-rw-r--r-- 1 riedell jusers 11490434 Jan 20 22:05 mnist.npz
```



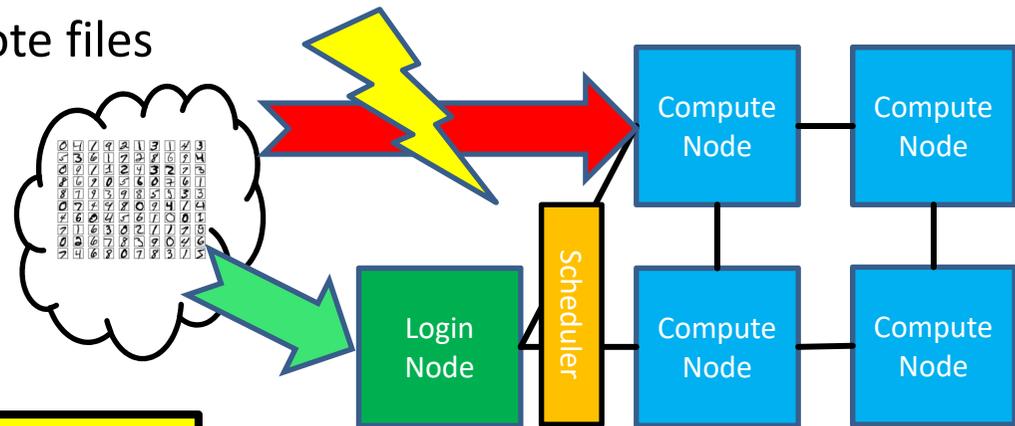
MNIST Dataset – Data Access & HPC Challenges

- Warning for HPC environments
 - Note that **HPC batch nodes** often do not allow for download of remote files

(2) Data Understanding Phase

```
# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()

Downloading data from https://s3.amazonaws.com/img-datasets/mnist.npz
11493376/11490434 [=====] - 6s 1us/step
```



- A useful workaround for download remotely stored datasets and files is to start the Keras script on the login node and after data download stop the script for a proper execution on batch nodes for training & inference

```
import numpy as np
from keras.datasets import mnist
```

```
# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```
[riedell@juron1-adm datasets]$ pwd
/p/home/jusers/riedell/juron/.keras/datasets
[riedell@juron1-adm datasets]$ ls -al
total 11234
drwxr-xr-x 2 riedell jusers 4096 Jan 20 22:05 .
drwxr-xr-x 3 riedell jusers 4096 Jan 20 22:03 ..
-rw-r--r-- 1 riedell jusers 11490434 Jan 20 22:05 mnist.npz
```



MNIST Dataset – Training and Testing Datasets

- Different Phases in Learning

(3) Data Preparation Phase

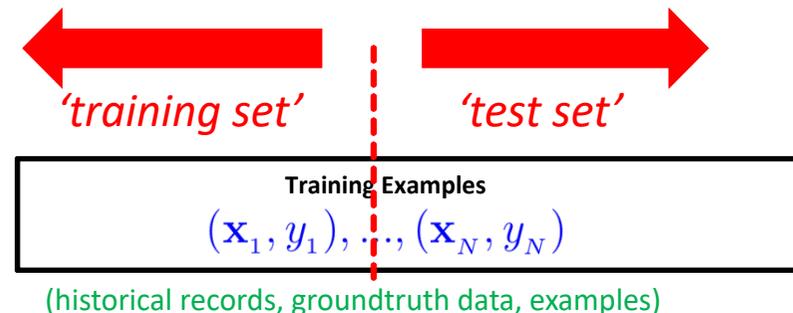
- **Training** phase is a hypothesis search
- **Testing** phase checks if we are on right track (once the hypothesis clear)
- **Validation** phase for model selection (more details later in tutorial)

- Start Work on Two disjoint datasets

- One **for training only (i.e. training set)**, one **for testing only (i.e. test set)**
- Exact separation is **rule of thumb per use case** (e.g. 10 % training, 90% test)
- Practice: If you get a dataset take immediately test data away (**‘throw it into the corner and forget about it during modelling’**)
- Once we learned from training data it has an **‘optimistic bias’**

```
import numpy as np
from keras.datasets import mnist

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```



MNIST Dataset – Exploration – One Character Encoding

```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 3 18 18 18 126 136 175 26 166 255 247 127 0 0 0 0
0 0 0 0 0 0 0 0 30 36 94 154 170 253 253 253 253 253 225 172 253 242 195 64 0 0 0 0
0 0 0 0 0 0 0 49 238 253 253 253 253 253 253 253 251 93 82 82 56 39 0 0 0 0 0 0
0 0 0 0 0 0 0 18 219 253 253 253 253 253 198 182 247 241 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 80 156 107 253 253 205 11 0 43 154 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 14 1 154 253 90 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 139 253 190 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 11 190 253 70 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 35 241 225 160 108 1 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 81 240 253 253 119 25 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 45 186 253 253 150 27 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 16 93 252 253 187 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 249 253 249 64 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 46 130 183 253 253 207 2 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 39 148 229 253 253 253 250 182 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 24 114 221 253 253 253 253 201 78 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 23 66 213 253 253 253 253 198 81 2 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 18 171 219 253 253 253 253 195 80 9 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 55 172 226 253 253 253 253 244 133 11 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 136 253 253 253 212 135 132 16 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
```

Label:
5

MNIST Dataset – Data Exploration Script Training Data

```
import numpy as np
from keras.datasets import mnist

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```
# function to explore one hand-written character
def character_show(character):
    for y in character:
        row = ""
        for x in y:
            row += '{0: <4}'.format(x)
        print(row)
```

```
# view first 10 hand-written characters
for i in range (0,9):
    character_show(X_train[i])
    print("\n")
    print("Label:")
    print(y_train[i])
    print("\n")
```

- Loading MNIST training datasets (X) with labels (Y) stored in a binary numpy format

- Format is 28 x 28 pixel values with grey level from 0 (white background) to 255 (black foreground)

- Small helper function that prints row-wise one 'hand-written' character with the grey levels stored in training dataset
- Should reveal the nature of the number (aka label)

- Example: loop of the training dataset (e.g. first 10 characters as shown here)
- At each loop interval the 'hand-written' character (X) is printed in 'matrix notation' & label (Y)

Startup Jupyter – Remember Kernel ‘dl_tutorial_jureca’

The screenshot shows the JupyterLab web interface. The browser address bar displays the URL: `https://jupyter-jsc.fz-juelich.de/user/uid%3Dm.riedel@fz-juelich.de/lab?redirects=1`. The interface includes a menu bar (File, Edit, View, Run, Kernel, Hub, Tabs, Settings, Help), a file browser on the left, and a central code editor for 'Untitled1.ipynb'. A 'Select Kernel' dialog box is open, showing a list of kernels. The 'dl_tutorial_jureca' kernel is highlighted with a red box. To the right, a red-bordered box highlights the kernel status area, which shows 'dl_tutorial_jureca' with a grey circle and a red arrow pointing to it. Below this, the text '(wait until becomes white)' is written in green. At the bottom right, the text 'Logout Kernel when done!' is written in red. The system tray at the bottom shows the time as 06:20 on 22.05.2019.

Select Kernel

Select kernel for: "Untitled1.ipynb"

- Python 3
- Start Preferred Kernel
 - Python 3
- Use No Kernel
 - No Kernel
- Start Other Kernel
 - Bash
 - C++11
 - C++14
 - dl_tutorial_jureca**
 - Javascript (Node.js)
 - ml_tutorial_jureca
- Use Kernel from Preferred Session
- Use Kernel from Other Session

dl_tutorial_jureca ●

(wait until becomes white)

dl_tutorial_jureca ○

Logout Kernel when done!

Exercises – Explore Testing Data



Exercises – Explore Testing Data – Solution

```
In [1]: import numpy as np
        from keras.datasets import mnist

        Using TensorFlow backend.

In [2]: # download and shuffled as training and testing set
        (X_train, y_train), (X_test, y_test) = mnist.load_data()

In [3]: # function to explore one hand-written character
        def character_show(character):
            for y in character:
                row = ""
                for x in y:
                    row += '{0: <4}'.format(x)
                print(row)

In [4]: # view first 10 hand-written characters
        for i in range(0, 9):
            character_show(X_test[i])
            print("\n")
            print(y_test[i])
            print("\n")

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 84 185 159 151 60 36 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 222 254 254 254 254 241 198 198 198 198 198 198 198 170 52 0 0 0 0 0
0 0 0 0 0 0 0 67 114 72 114 163 227 254 225 254 254 254 250 229 254 254 140 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 17 66 14 67 67 67 59 21 236 254 106 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 83 253 209 18 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 22 233 255 83 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 129 254 238 44 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 59 249 254 62 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 133 254 187 5 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 205 248 58 0 0 0 0 0 0 0 0 0 0
```


MNIST Dataset – Reshape & Normalization

```
import numpy as np
from keras.datasets import mnist

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```
# reshape for input to perceptron 28 x 28 = 784
RESHAPED = 784
X_train = X_train.reshape(60000, RESHAPED)
X_test = X_test.reshape(10000, RESHAPED)
# float32 for GPU execution
X_train = X_train.astype('float32')
X_test = X_test.astype('float32')
```

```
# normalization
X_train /= 255
X_test /= 255
```

```
# data exploration: number of samples
print(X_train.shape[0], 'train samples')
print(X_test.shape[0], 'test samples')
```

```
# data exploration: number of values / samples
print(X_train.shape[1], 'input pixel values per train samples')
print(X_test.shape[1], 'input pixel values per test samples')
```

```
# data output: vectorized character
print(X_train[0])
```

(3) Data Preparation Phase

- Loading MNIST training datasets (X) and testing datasets (Y) stored in a binary numpy format with labels for X and Y
- Format is 28 x 28 pixel values with grey level from 0 (white background) to 255 (black foreground)
- Reshape from 28 x 28 matrix of pixels to 784 pixel values considered to be the input for the neural networks later
- Normalization is added for mathematical convenience since computing with numbers get easier (not too large)

Exercises – Perform Data Reshaping & Normalization

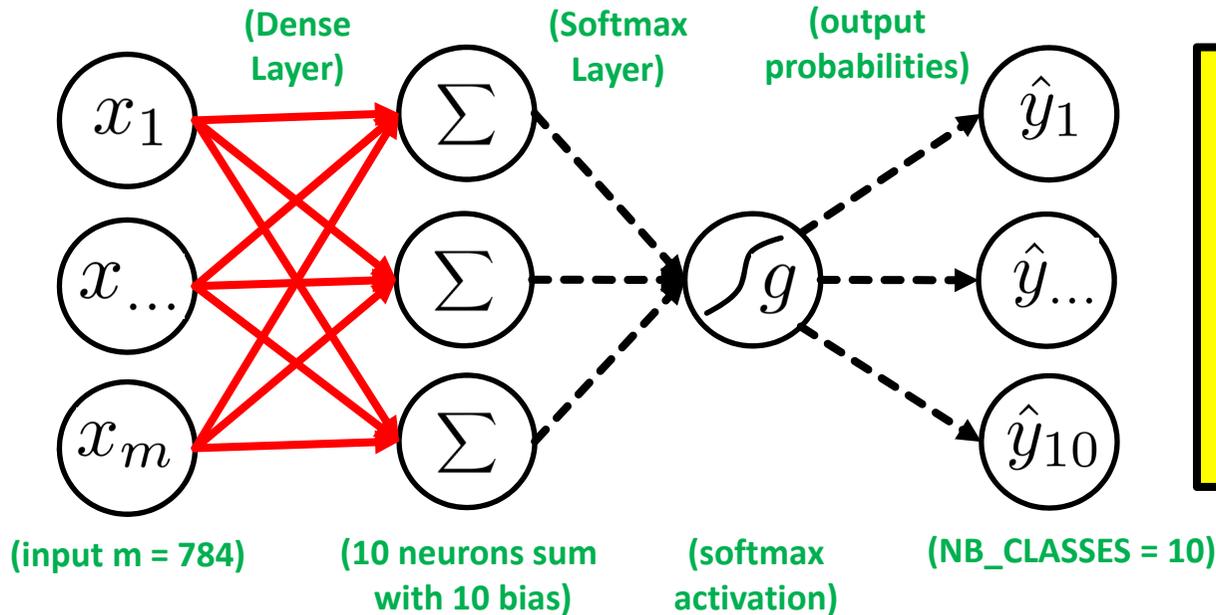


MNIST Dataset & Multi Output Perceptron Model

- 10 Class Classification Problem

(4) Modeling Phase

- Use 10 Perceptrons for 10 outputs with softmax activation function



- Note that the output units are independent among each other in contrast to neural networks with one hidden layer
- The output of softmax gives class probabilities

```
from keras.models import Sequential
from keras.layers.core import Dense, Activation

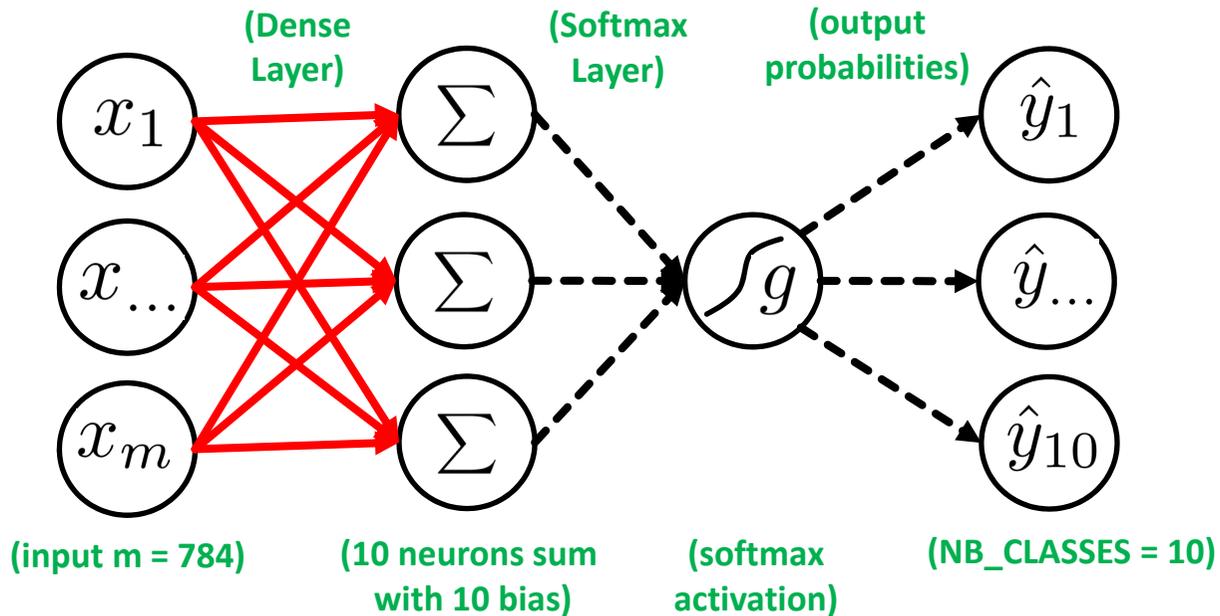
# model Keras sequential
model = Sequential()

# add fully connected layer - input with output
model.add(Dense(NB_CLASSES, input_shape=(RESHAPED,)))
```



MNIST Dataset & Activation Function Softmax

- Activation Function Softmax
 - Softmax enables probabilities for 10 classes



```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
```

```
# add activation function layer to get class probabilities
model.add(Activation('softmax'))
```

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{h=1}^K e^{z_h}}$$

(4) Modeling Phase

The non-linear Activation function 'softmax' represents a generalization of the sigmoid function – it squashes an n-dimensional vector of arbitrary real values into a n-dimensional vector of real values in the range of 0 and 1 – here it aggregates 10 answers provided by the Dense layer with 10 neurons



AUDIENCE QUESTION

How many parameters we have to learn and why?

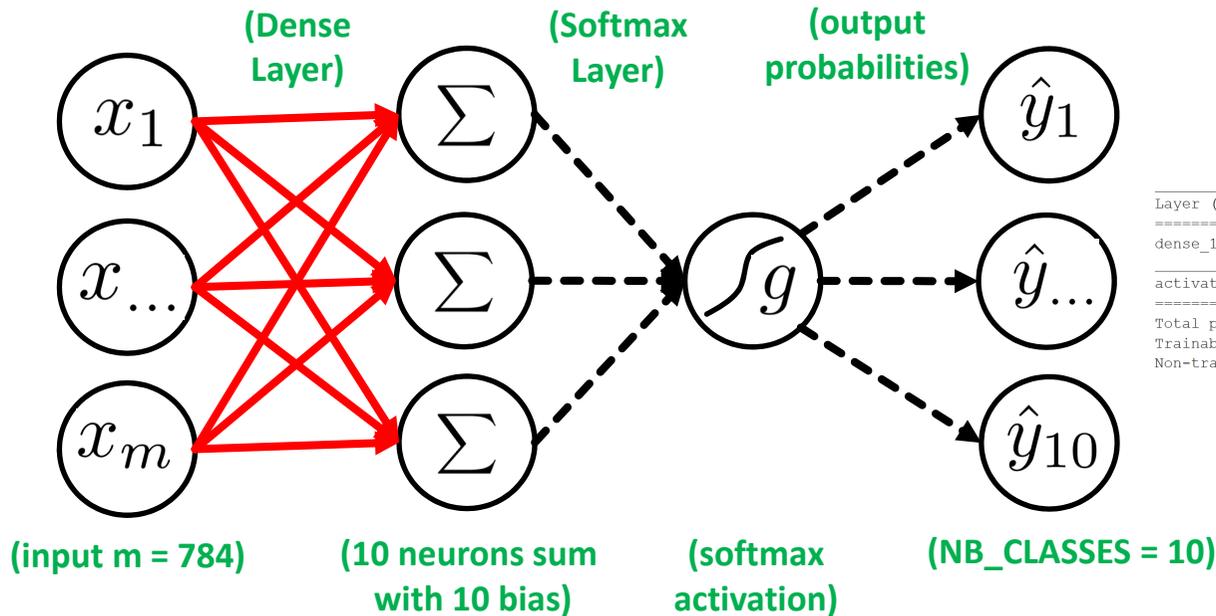


MNIST Dataset & Model Summary & Parameters

- Activation Function Softmax

(4) Modeling Phase

- Softmax enables probabilities for 10 classes



Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 10)	7850
activation_1 (Activation)	(None, 10)	0
Total params: 7,850		
Trainable params: 7,850		
Non-trainable params: 0		

(parameters = $784 * 10 + 10$ bias = 7850)



```
# printout a summary of the model to understand model complexity
model.summary()
```

MNIST Dataset & Compile the Model

- **Compile** the model

(4) Modeling Phase

- Choose **optimizer** as algorithm used to update weights while training the model
- Specify **loss function** (i.e. objective function) that is used by the optimizer to navigate the space of weights (note: process of optimization is also called loss minimization)
- Indicate **metric** for model evaluation

■ **Compile the model to be executed by the Keras backend (e.g. TensorFlow)**

```
In [1]: import numpy as np
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD
from keras.utils import np_utils
```

Using TensorFlow backend.

- Specify **loss function**

- Compare prediction vs. Given class label
- E.g. **categorical crossentropy**

- **Loss function is a multi-class logarithmic loss: target is $t_{i,j}$ and prediction is $p_{i,j}$**
- **Categorical crossentropy is very suitable for multiclass label predictions (default with softmax)**



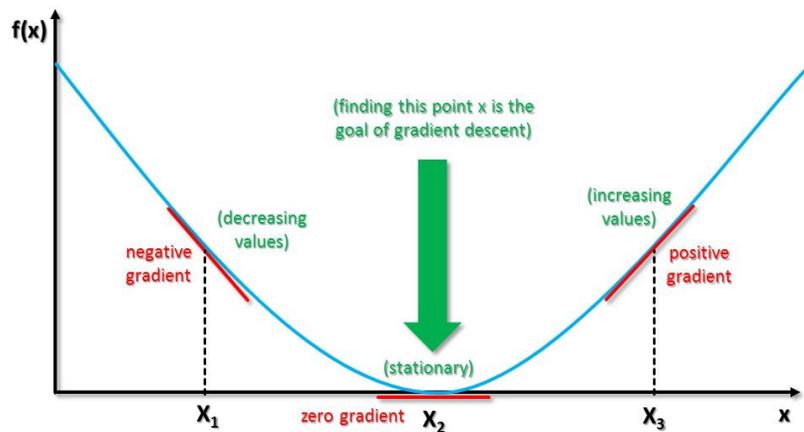
```
# specify loss, optimizer and metric
model.compile(loss='categorical_crossentropy', optimizer=OPTIMIZER, metrics=['accuracy'])
```

$$L_i = -\sum_j t_{i,j} \log(p_{i,j})$$

MNIST Dataset & Optimization / Learning Approach

- Choosing an **Optimizer**
 - Example: **Stochastic Gradient Descent (SGD)**

(4) Modeling Phase



(minimization: subtract gradient term because we move towards local minima)

$$b = a - \gamma \nabla f(a)$$

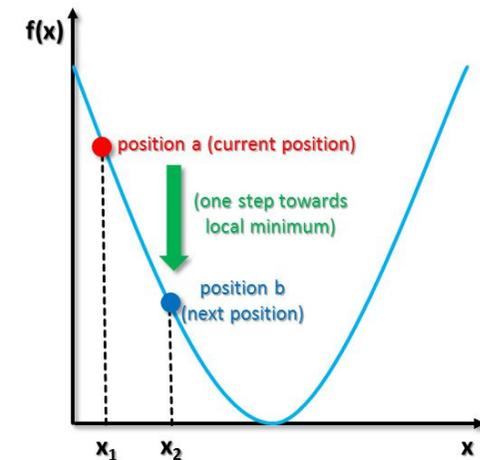
(the derivative of f with respect to a)

(old position before the step)

(new position after the step)

(weighting factor known as step-size, can change at every iteration, also called learning rate)

(gradient term is steepest ascent)



[4] Big Data Tips, Gradient Descent



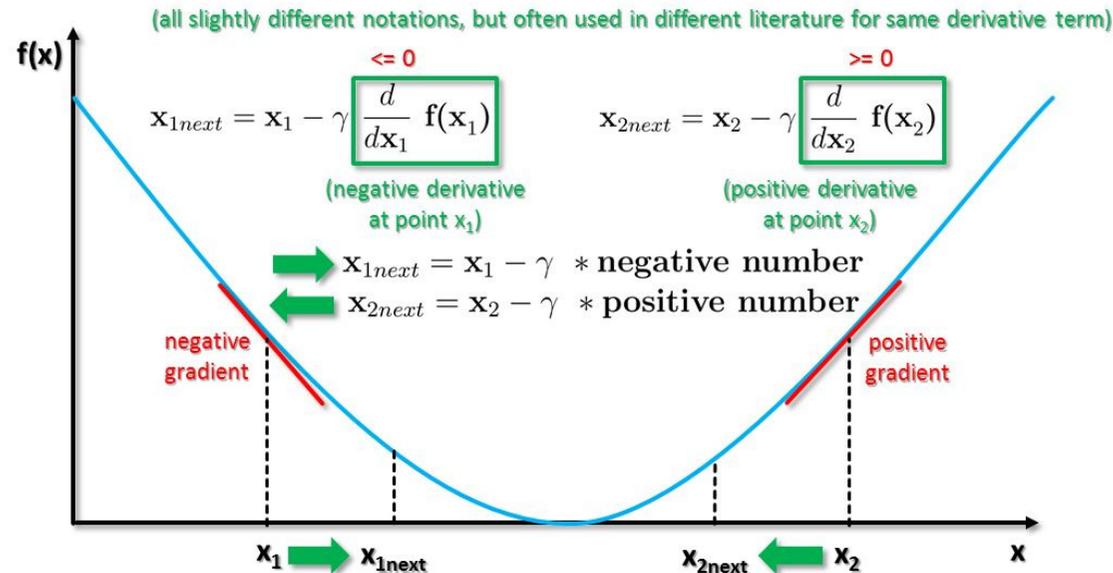
```
from keras.optimizers import SGD
```

```
OPTIMIZER = SGD() # optimization technique
```

MNIST Dataset & Stochastic Gradient Descent Method

- Gradient Descent (GD) uses all the training samples available for a step within a iteration
- Stochastic Gradient Descent (SGD) converges faster: only one training samples used per iteration

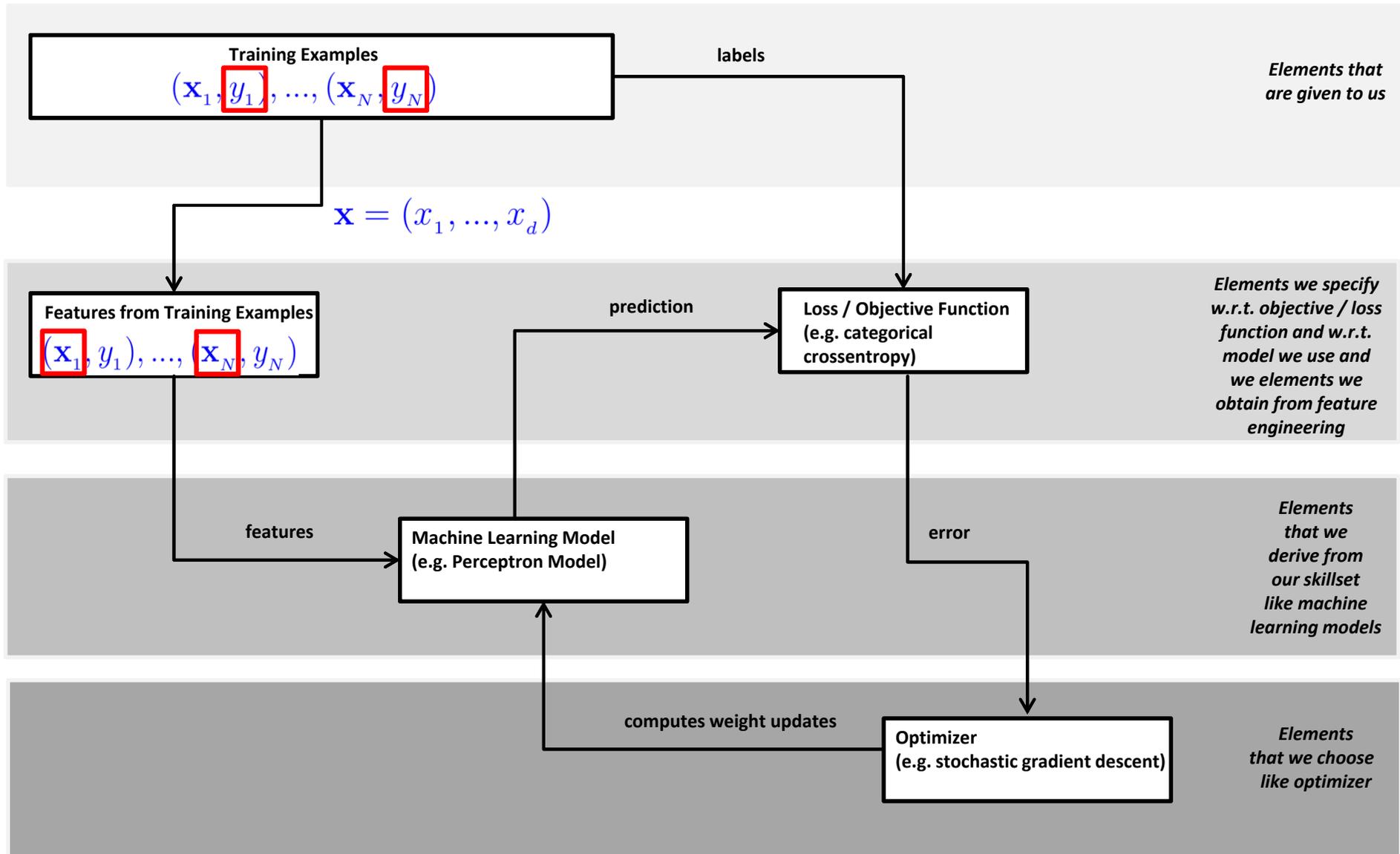
$$b = a - \gamma \nabla f(a) \quad b = a - \gamma \frac{\partial}{\partial a} f(a) \quad b = a - \gamma \frac{d}{da} f(a)$$



```
from keras.optimizers import SGD  
  
OPTIMIZER = SGD() # optimization technique
```

[4] Big Data Tips, Gradient Descent

Optimizer – Effect to the Model – Overview & SGD Example



MNIST Dataset – Model Parameters & Data Normalization

```
import numpy as np
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD
from keras.utils import np_utils
```

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
```

```
# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# X_train is 60000 rows of 28x28 values --> reshaped in 60000 x 784
RESHAPED = 784
X_train = X_train.reshape(60000, RESHAPED)
X_test = X_test.reshape(10000, RESHAPED)
X_train = X_train.astype('float32')
X_test = X_test.astype('float32')

# normalize
X_train /= 255
X_test /= 255
```

```
# output number of samples
print(X_train.shape[0], 'train samples')
print(X_test.shape[0], 'test samples')
```

- **NB_CLASSES: 10 Class Problem**
- **NB_EPOCH: number of times the model is exposed to the overall training set – at each iteration the optimizer adjusts the weights so that the objective function is minimized**
- **BATCH_SIZE: number of training instances taken into account before the optimizer performs a weight update to the model**
- **OPTIMIZER: Stochastic Gradient Descent ('SGD') – only one training sample/iteration**

- **Data load shuffled between training and testing set in files**
- **Data preparation, e.g. X_train is 60000 samples / rows of 28 x 28 pixel values that are reshaped in 60000 x 784 including type specification (i.e. float32)**
- **Data normalization: divide by 255 – the max intensity value to obtain values in range [0,1]**

MNIST Dataset – A Multi Output Perceptron Model

- The Sequential() Keras model is a linear pipeline (aka 'a stack') of various neural network layers including Activation functions of different types (e.g. softmax)

- Dense() represents a fully connected layer used in ANNs that means that each neuron in a layer is connected to all neurons located in the previous layer

- The non-linear activation function 'softmax' is a generalization of the sigmoid function – it squashes an n-dimensional vector of arbitrary real values into a n-dimensional vector of real values in the range of 0 and 1 – here it aggregates 10 answers provided by the Dense layer with 10 neurons

```
# convert class label vectors using one hot encoding
Y_train = np_utils.to_categorical(y_train, NB_CLASSES)
Y_test = np_utils.to_categorical(y_test, NB_CLASSES)
```

```
# model Keras sequential
model = Sequential()
```

```
# add fully connected layer - input with output
model.add(Dense(NB_CLASSES, input_shape=(RESHAPED,)))
```

```
# add activation function layer to get class probabilities
model.add(Activation('softmax'))
```

```
# printout a summary of the model to understand model complexity
model.summary()
```

```
# specify loss, optimizer and metric
model.compile(loss='categorical_crossentropy', optimizer=OPTIMIZER, metrics=['accuracy'])
```

```
# model training
history = model.fit(X_train, Y_train, batch_size=BATCH_SIZE, epochs=NB_EPOCH, verbose=VERBOSE)
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

- Loss function is a multi-class logarithmic loss: target is $t_{i,j}$ and prediction is $p_{i,j}$

$$L_i = -\sum_j t_{i,j} \log(p_{i,j})$$

- Train the model ('fit')

Exercises – Multi Output Perceptron Model & 20 Epochs



Model Evaluation – Testing Phase & Confusion Matrix

- Model is fixed
 - Model is just used with the testset
 - Parameters are set
- Evaluation of model performance
 - Counts of test records that are incorrectly predicted
 - Counts of test records that are correctly predicted
 - E.g. create **confusion matrix** for a two class problem

Counting per sample		Predicted Class	
		Class = 1	Class = 0
Actual Class	Class = 1	f_{11}	f_{10}
	Class = 0	f_{01}	f_{00}

(serves as a basis for further performance metrics usually used)

Model Evaluation – Testing Phase & Performance Metrics

Counting per sample		Predicted Class	
		Class = 1	Class = 0
Actual Class	Class = 1	f_{11}	f_{10}
	Class = 0	f_{01}	f_{00}

(100% accuracy in learning often points to problems using machine learning methods in practice)

- Accuracy (usually in %)

$$Accuracy = \frac{\text{number of correct predictions}}{\text{total number of predictions}}$$

- Error rate

$$Error\ rate = \frac{\text{number of wrong predictions}}{\text{total number of predictions}}$$

Using JURECA & GPUs – Remember use GPU Partition

JURECA

Hardware Characteristics of the Cluster Module

- 1872 compute nodes
 - Two Intel Xeon E5-2680 v3 Haswell CPUs per node
 - 2 x 12 cores, 2.5 GHz
 - Intel Hyperthreading Technology (Simultaneous Multithreading)
 - AVX 2.0 ISA extension
 - 75 compute nodes equipped with two NVIDIA K80 GPUs (four visible devices per node)
 - 2 x 4992 CUDA cores
 - 2 x 24 GiB GDDR5 memory
 - DDR4 memory technology (2133 MHz)
 - 1605 compute nodes with 128 GiB memory
 - 128 compute nodes with 256 GiB memory
 - 64 compute nodes with 512 GiB memory
- 12 visualization nodes
 - Two Intel Xeon E5-2680 v3 Haswell CPUs per node
 - Two NVIDIA K40 GPUs per node
 - 2 x 12 GiB GDDR5 memory
 - 10 nodes with 512 GiB memory
 - 2 nodes with 1024 GiB memory
- Login nodes with 256 GiB memory per node
- 45,216 CPU cores
- 1.8 (CPU) + 0.44 (GPU) Petaflop per second peak performance
- Based on the T-Platforms V-class server architecture
- Mellanox EDR InfiniBand high-speed network with non-blocking fat tree topology
- 100 GiB per second storage connection to [JUST](#)

Exercises – Evaluate Multi Output Perceptron Model



MNIST Dataset – A Multi Output Perceptron Model – Output

```
Epoch 7/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4419 - acc: 0.8838
Epoch 8/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4271 - acc: 0.8866
Epoch 9/20
60000/60000 [=====] - 2s 25us/step - loss: 0.4151 - acc: 0.8888
Epoch 10/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4052 - acc: 0.8910
Epoch 11/20
60000/60000 [=====] - 2s 26us/step - loss: 0.3968 - acc: 0.8924
Epoch 12/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3896 - acc: 0.8944
Epoch 13/20
60000/60000 [=====] - 2s 26us/step - loss: 0.3832 - acc: 0.8956
Epoch 14/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3777 - acc: 0.8969
Epoch 15/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3727 - acc: 0.8982
Epoch 16/20
60000/60000 [=====] - 1s 24us/step - loss: 0.3682 - acc: 0.8989
Epoch 17/20
60000/60000 [=====] - 1s 25us/step - loss: 0.3641 - acc: 0.9001
Epoch 18/20
60000/60000 [=====] - 1s 25us/step - loss: 0.3604 - acc: 0.9007
Epoch 19/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3570 - acc: 0.9016
Epoch 20/20
60000/60000 [=====] - 1s 24us/step - loss: 0.3538 - acc: 0.9023
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 41us/step
Test score: 0.33423959468007086
Test accuracy: 0.9101
```

AUDIENCE QUESTION

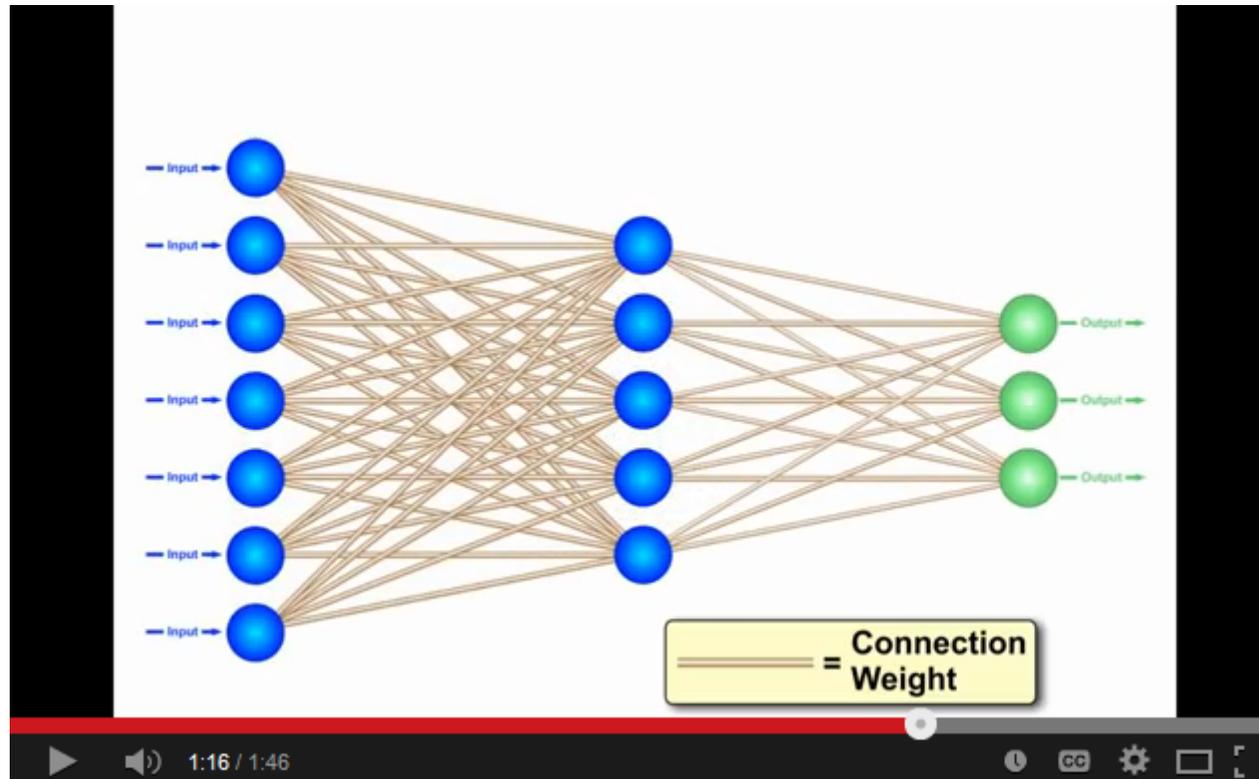
What would you change to get better accuracy?



Exercises – Multi Output Perceptron Model & 50 Epochs

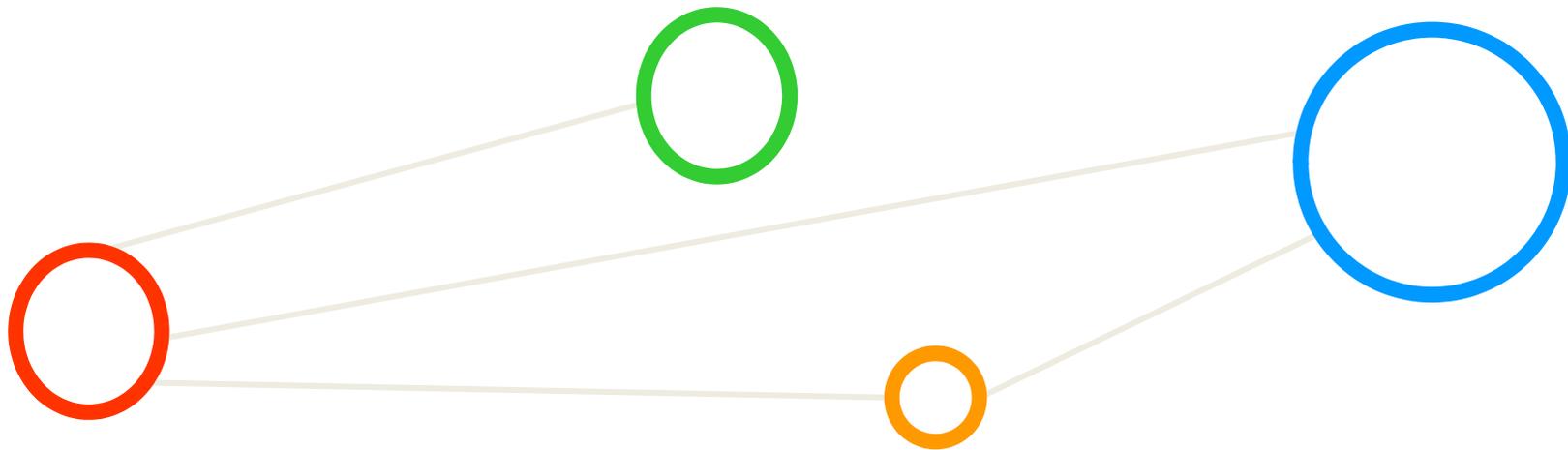


[Video] Towards Multi-Layer Perceptrons



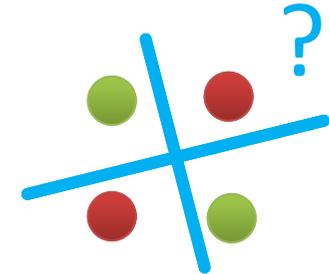
[5] YouTube Video, Neural Networks – A Simple Explanation

Artificial Neural Network (ANNs)



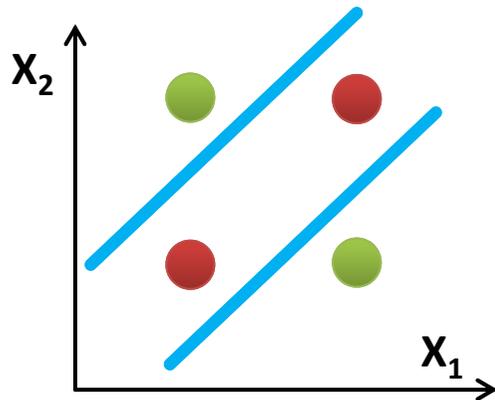
Artificial Neural Network (ANN)

- Simple perceptrons fail: 'not linearly seperable'



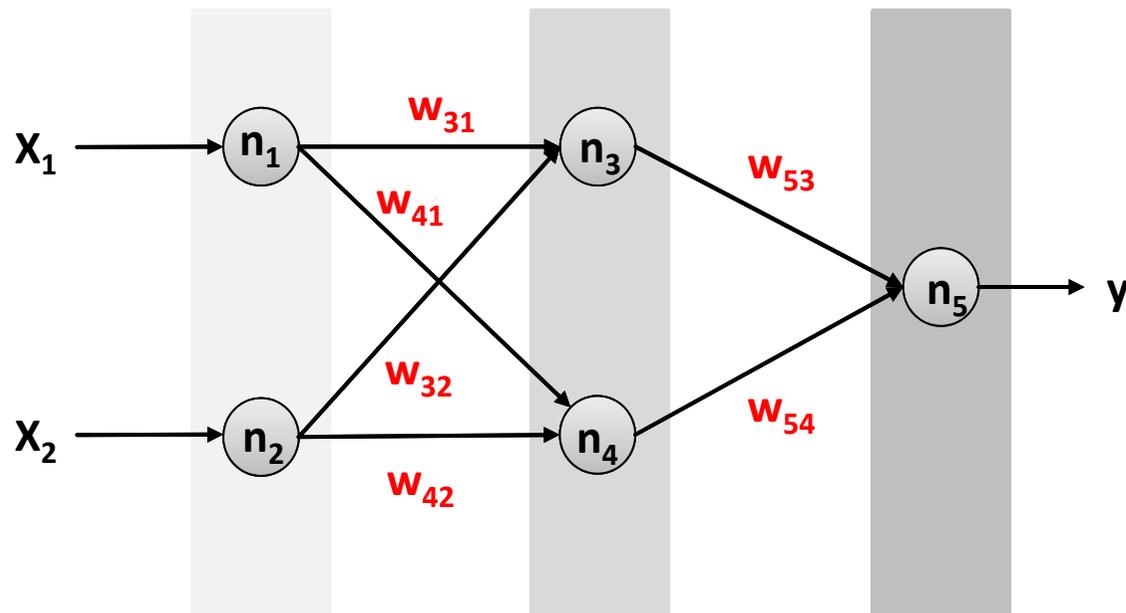
X_1	X_2	Y
0	0	-1
1	0	1
0	1	1
1	1	-1

Labelled Data Table



Decision Boundary

(Idea: instances can be classified using two lines at once to model XOR)



Two-Layer, feed-forward Artificial Neural Network topology

Machine Learning Challenges – Problem of Overfitting

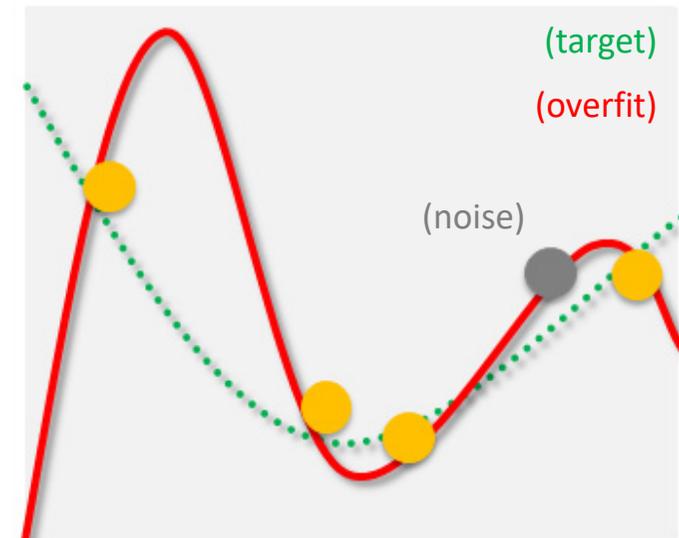
- Overfitting refers to fit the data too well – more than is warranted – thus may misguide the learning
- Overfitting is not just ‘bad generalization’ - e.g. the VC dimension covers noiseless & noise targets
- Theory of Regularization are approaches against overfitting and prevent it using different methods

- Key problem: noise in the target function leads to overfitting

- Effect: ‘noisy target function’ and its noise misguides the fit in learning
- There is always ‘some noise’ in the data
- Consequence: poor target function (‘distribution’) approximation

- Example: Target functions is second order polynomial (i.e. parabola)

- Using a higher-order polynomial fit
- Perfect fit: low $E_{in}(g)$, but large $E_{out}(g)$



(but simple polynomial works good enough)

(‘over’: here meant as 4th order, a 3rd order would be better, 2nd best)

Problem of Overfitting – Clarifying Terms

- A good model must have low training error (E_{in}) and low generalization error (E_{out})
- Model overfitting is if a model fits the data too well (E_{in}) with a poorer generalization error (E_{out}) than another model with a higher training error (E_{in})

- **Overfitting & Errors**

- $E_{in}(g)$ goes **down**

- $E_{out}(g)$ goes **up**

- **'Bad generalization area' ends**

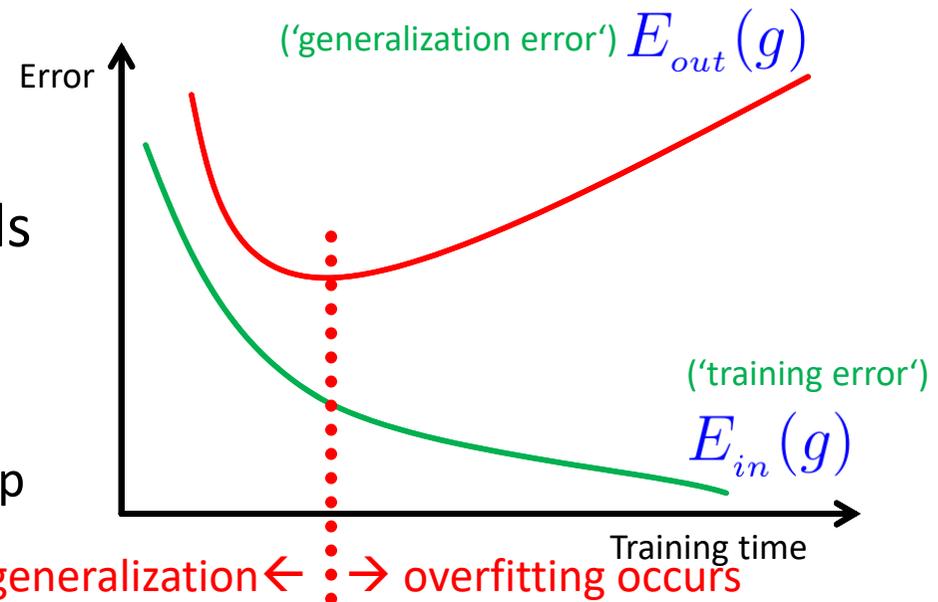
- Good to reduce $E_{in}(g)$

- **'Overfitting area' starts**

- Reducing $E_{in}(g)$ does not help

- Reason **'fitting the noise'**

bad generalization ← → overfitting occurs



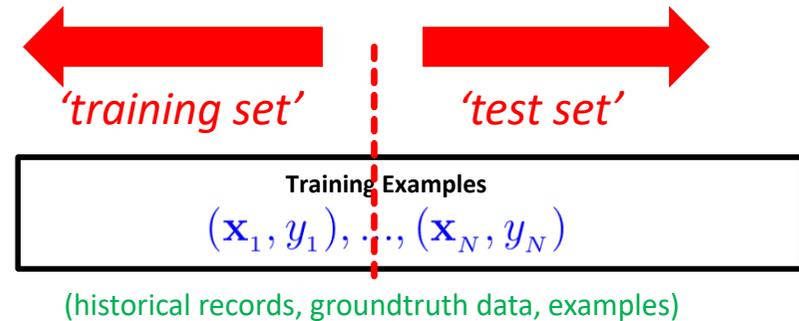
- The two general approaches to prevent overfitting are (1) regularization and (2) validation

Terminologies & Different Dataset Elements

- **Target Function** $f : X \rightarrow Y$
 - Ideal function that ‘explains’ the data we want to learn
- **Labelled Dataset (samples)**
 - ‘in-sample’ data given to us: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- **Learning vs. Memorizing**
 - The goal is to create a system that works well ‘out of sample’
 - In other words we want to classify ‘future data’ (out of sample) correct
- **Dataset Part One: Training set**
 - Used for training a machine learning algorithms
 - Result after using a training set: a trained system
- **Dataset Part Two: Test set**
 - Used for testing whether the trained system might work well
 - Result after using a test set: accuracy of the trained model

Model Evaluation – Training and Testing Phases

- Different Phases in Learning (cf. day one remote sensing)
 - **Training** phase is a hypothesis search
 - **Testing** phase checks if we are on right track (once the hypothesis clear)
- Work on **‘training examples’**
 - Create **two disjoint datasets**
 - One used **for training only** (aka training set)
 - Another **used for testing only** (aka test set)
 - Exact separation is **rule of thumb per use case** (e.g. 10 % training, 90% test)
 - Practice: If you get a dataset take immediately test data away (**‘throw it into the corner and forget about it during modelling’**)
 - Reasoning: Once we learned from training data it has an **‘optimistic bias’**



Learning Approaches – Supervised Learning – Formalization

- Each observation of the predictor measurement(s) has an associated response measurement:
 - Input $\mathbf{x} = x_1, \dots, x_d$
 - Output $y_i, i = 1, \dots, n$
 - Data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- Goal: Fit a model that relates the response to the predictors
 - **Prediction:** Aims of accurately predicting the response for future observations
 - **Inference:** Aims to better understanding the relationship between the response and the predictors

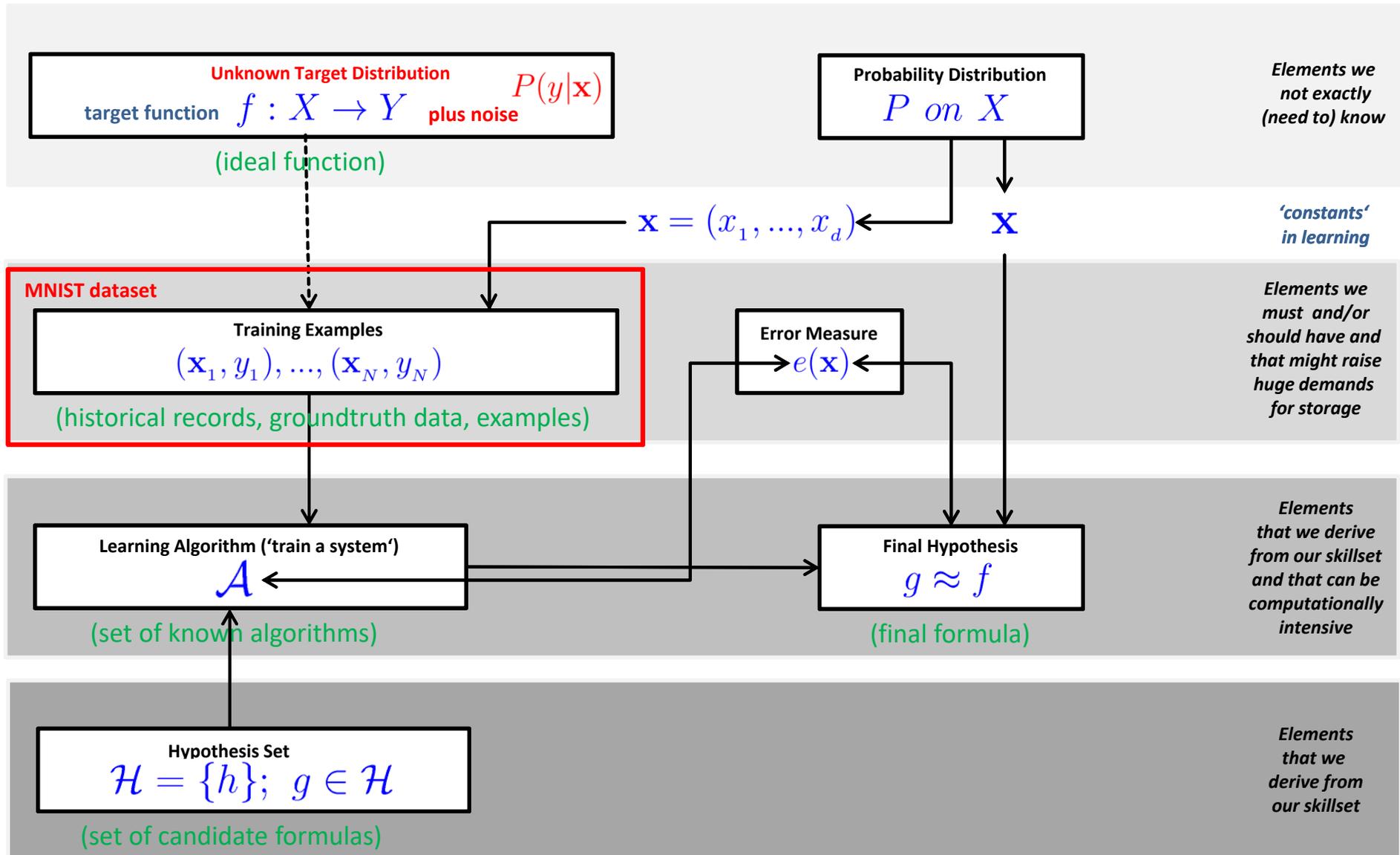
Training Examples
 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

(historical records, groundtruth data, examples)

- Supervised learning approaches fits a model that related the response to the predictors
- Supervised learning approaches are used in classification algorithms such as SVMs
- Supervised learning works with data = [input, correct output]

[9] *An Introduction to Statistical Learning*

Supervised Learning – Training Examples



Handwritten Character Recognition MNIST Dataset

- Metadata

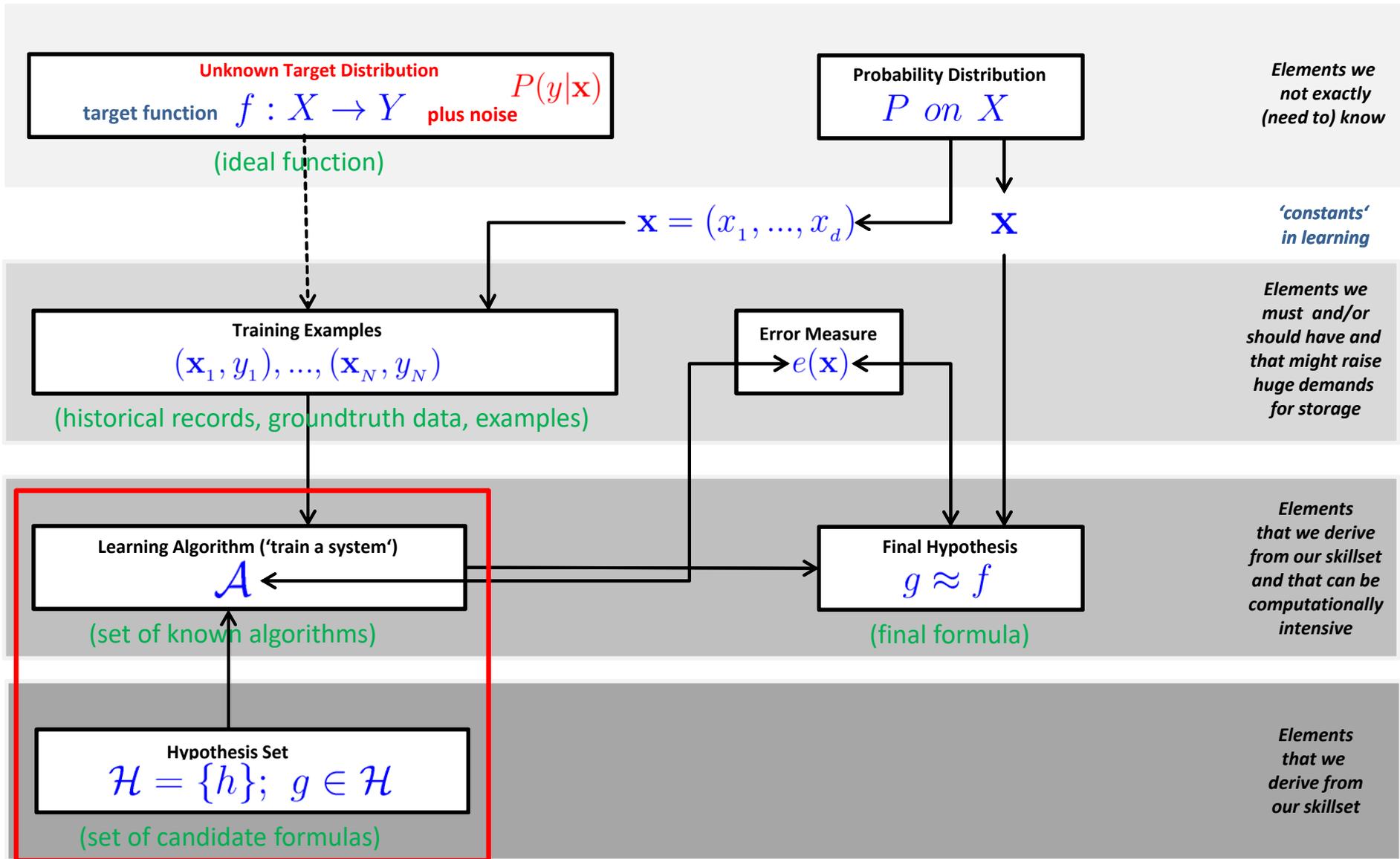
- Subset of a larger dataset from US National Institute of Standards (NIST)
- Handwritten digits including corresponding labels with values 0 to 9
- All digits have been size-normalized to 28 * 28 pixels and are centered in a fixed-size image for direct processing
- Not very challenging dataset, but good for experiments / tutorials

- Dataset Samples

- Labelled data (10 classes)
- Two separate files for training and test
- 60000 training samples (~47 MB)
- 10000 test samples (~7.8 MB)



Supervised Learning – Many Hypothesis to Choose



Different Models – Understanding the Hypothesis Set

Hypothesis Set

$$\mathcal{H} = \{h\}; g \in \mathcal{H}$$

$$\mathcal{H} = \{h_1, \dots, h_m\};$$

(all candidate functions derived from models and their parameters)

- Choosing from various model approaches h_1, \dots, h_m is a different hypothesis
- Additionally a change in model parameters of h_1, \dots, h_m means a different hypothesis too

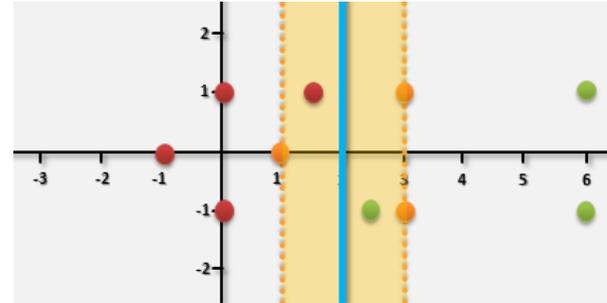
‘select one function’ that best approximates

Final Hypothesis

$$g \approx f$$

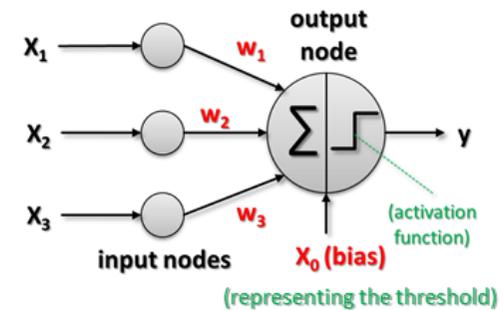


h_1



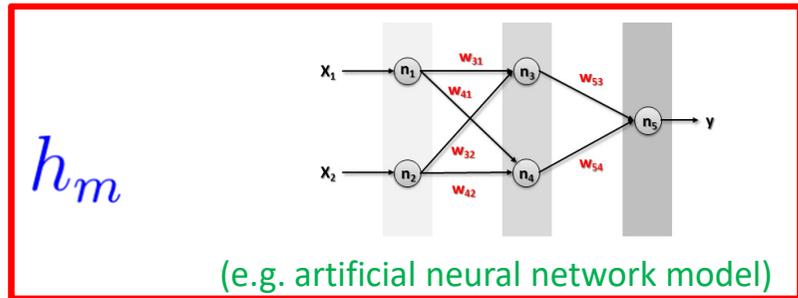
(e.g. support vector machine model)

h_2



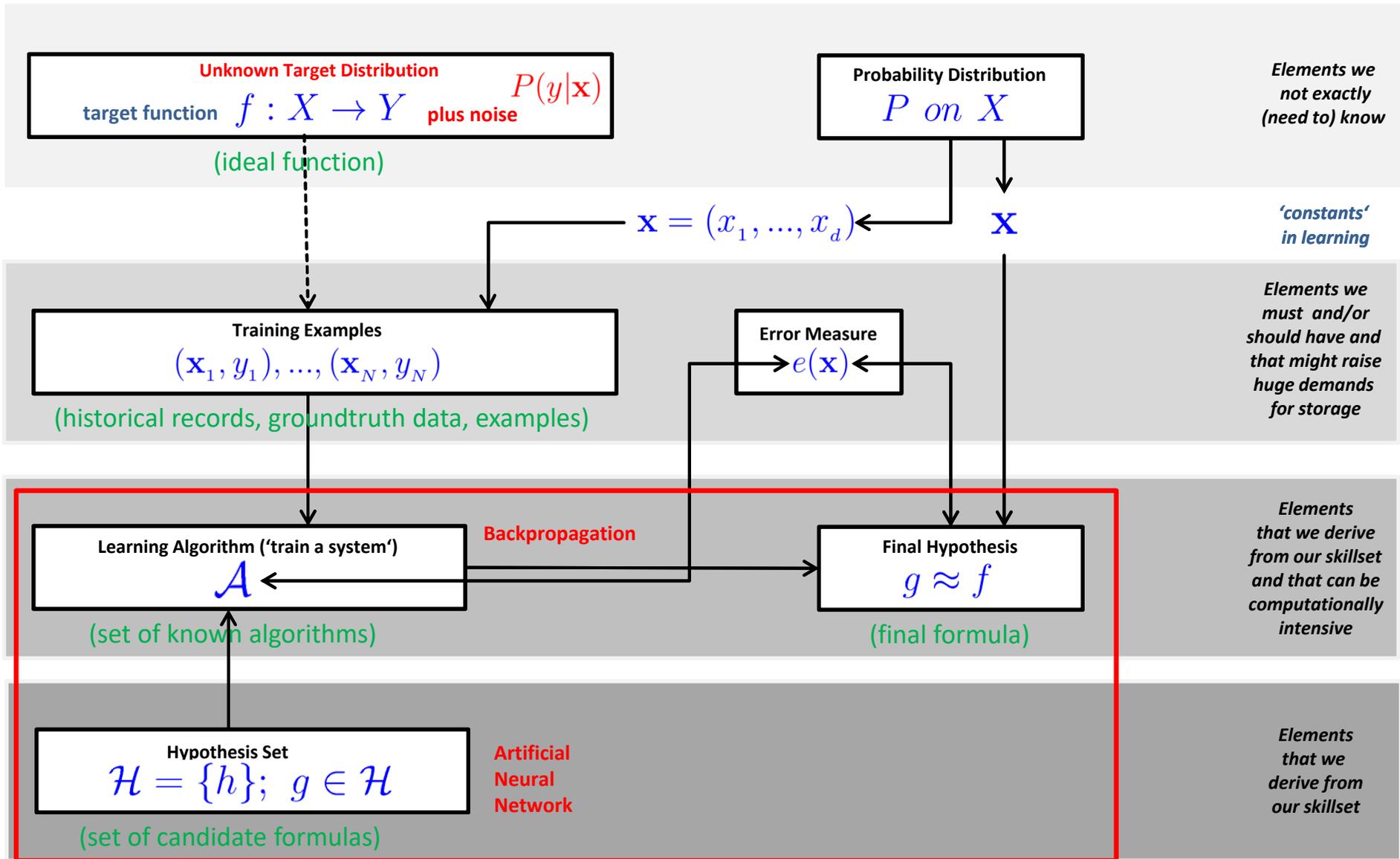
(e.g. linear perceptron model)

h_m



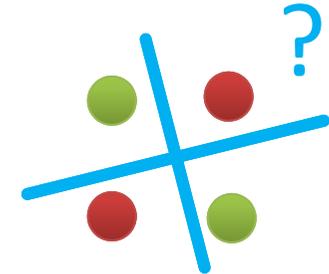
(e.g. artificial neural network model)

Supervised Learning – Training Examples



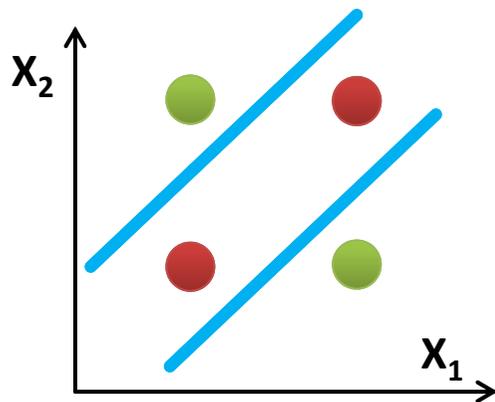
Artificial Neural Network (ANN) – Revisited

- Simple perceptrons fail: 'not linearly seperable'



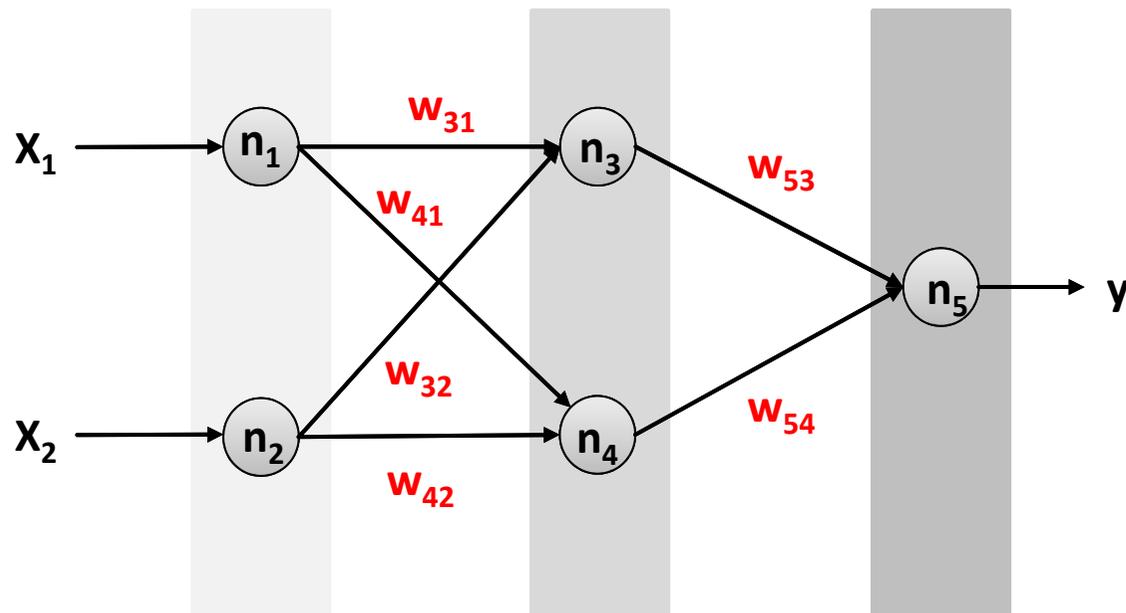
x_1	x_2	y
0	0	-1
1	0	1
0	1	1
1	1	-1

Labelled Data Table



Decision Boundary

(Idea: instances can be classified using two lines at once to model XOR)



Two-Layer, feed-forward Artificial Neural Network topology

High-level Tools – Keras

- Keras is a high-level deep learning library implemented in Python that works on top of existing other rather low-level deep learning frameworks like Tensorflow, CNTK, or Theano
- The key idea behind the Keras tool is to enable faster experimentation with deep networks
- Created deep learning models run seamlessly on CPU and GPU via low-level frameworks

```
keras.layers.Dense(units,  
                    activation=None,  
                    use_bias=True,  
                    kernel_initializer='glorot_uniform',  
                    bias_initializer='zeros',  
                    kernel_regularizer=None,  
                    bias_regularizer=None,  
                    activity_regularizer=None,  
                    kernel_constraint=None,  
                    bias_constraint=None)
```

```
keras.optimizers.SGD(lr=0.01,  
                     momentum=0.0,  
                     decay=0.0,  
                     nesterov=False)
```

- Tool Keras supports inherently the creation of artificial neural networks using Dense layers and optimizers (e.g. SGD)
- Includes regularization (e.g. weight decay) or momentum



Keras

[3] *Keras Python Deep Learning Library*

ANN – MNIST Dataset – Create ANN Blueprint

✓ Data Preprocessing done (i.e. data normalization, reshape, etc.)

1. Define a neural network topology

- Which layers are required?
- Think about input layer need to match the data – what data we had?
- Maybe hidden layers?
- Think Dense layer – Keras?
- Think about final Activation as Softmax (cf. Day One) → output probability

2. Compile the model → model representation for Tensorflow et al.

- Think about what loss function you want to use in your problem?
- What is your optimizer strategy, e.g. SGD (cf. Day One)

3. Fit the model → the model learning takes place

- How long you want to train (e.g. NB_EPOCHS)
- How much samples are involved (e.g. BATCH_SIZE)

Exercises – Create a Simple ANN Model – One Dense



MNIST Dataset – Model Parameters & Data Normalization

```
import numpy as np
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD
from keras.utils import np_utils
```

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
```

```
# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# X_train is 60000 rows of 28x28 values --> reshaped in 60000 x 784
RESHAPED = 784
X_train = X_train.reshape(60000, RESHAPED)
X_test = X_test.reshape(10000, RESHAPED)
X_train = X_train.astype('float32')
X_test = X_test.astype('float32')

# normalize
X_train /= 255
X_test /= 255
```

```
# output number of samples
print(X_train.shape[0], 'train samples')
print(X_test.shape[0], 'test samples')
```

- **NB_CLASSES: 10 Class Problem**
- **NB_EPOCH: number of times the model is exposed to the overall training set – at each iteration the optimizer adjusts the weights so that the objective function is minimized**
- **BATCH_SIZE: number of training instances taken into account before the optimizer performs a weight update to the model**
- **OPTIMIZER: Stochastic Gradient Descent ('SGD') – only one training sample/iteration**

- **Data load shuffled between training and testing set in files**
- **Data preparation, e.g. X_train is 60000 samples / rows of 28 x 28 pixel values that are reshaped in 60000 x 784 including type specification (i.e. float32)**
- **Data normalization: divide by 255 – the max intensity value to obtain values in range [0,1]**

MNIST Dataset – A Multi Output Perceptron Model

- The Sequential() Keras model is a linear pipeline (aka 'a stack') of various neural network layers including Activation functions of different types (e.g. softmax)

- Dense() represents a fully connected layer used in ANNs that means that each neuron in a layer is connected to all neurons located in the previous layer

- The non-linear activation function 'softmax' is a generalization of the sigmoid function – it squashes an n-dimensional vector of arbitrary real values into a n-dimensional vector of real values in the range of 0 and 1 – here it aggregates 10 answers provided by the Dense layer with 10 neurons

```
# convert class label vectors using one hot encoding
Y_train = np_utils.to_categorical(y_train, NB_CLASSES)
Y_test = np_utils.to_categorical(y_test, NB_CLASSES)

# model Keras sequential
model = Sequential()

# add fully connected layer - input with output
model.add(Dense(NB_CLASSES, input_shape=(RESHAPED,)))

# add activation function layer to get class probabilities
model.add(Activation('softmax'))

# printout a summary of the model to understand model complexity
model.summary()

# specify loss, optimizer and metric
model.compile(loss='categorical_crossentropy', optimizer=OPTIMIZER, metrics=['accuracy'])

# model training
history = model.fit(X_train, Y_train, batch_size=BATCH_SIZE, epochs=NB_EPOCH, verbose=VERBOSE)

# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

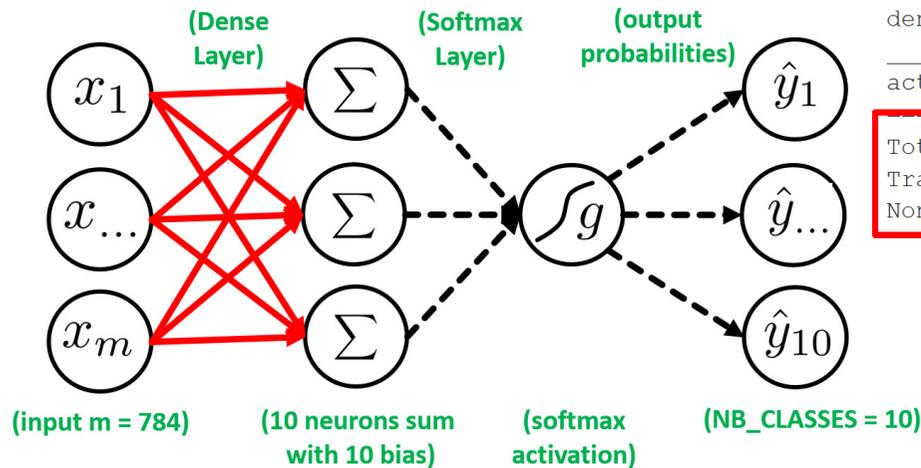
- Loss function is a multi-class logarithmic loss: target is $t_{i,j}$ and prediction is $p_{i,j}$

$$L_i = -\sum_j t_{i,j} \log(p_{i,j})$$

- Train the model ('fit')

MNIST Dataset & Model Summary & Parameters

- Activation Function Softmax
 - Softmax enables probabilities for 10 classes



Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 10)	7850
activation_1 (Activation)	(None, 10)	0

Total params: 7,850
 Trainable params: 7,850
 Non-trainable params: 0

(parameters = 784 * 10 + 10 bias = 7850)

■ Relevant for validation: Choosing a model with different layers is a model selection that directly also influences the number of parameters (e.g. add Dense layer from Keras means new weights)



```
# printout a summary of the model to understand model complexity
model.summary()
```

Model Evaluation – Testing Phase & Confusion Matrix

- Model is fixed
 - Model is just used with the testset
 - Parameters are set
- Evaluation of model performance
 - Counts of test records that are incorrectly predicted
 - Counts of test records that are correctly predicted
 - E.g. create **confusion matrix** for a two class problem

Counting per sample		Predicted Class	
		Class = 1	Class = 0
Actual Class	Class = 1	f_{11}	f_{10}
	Class = 0	f_{01}	f_{00}

(serves as a basis for further performance metrics usually used)

Model Evaluation – Testing Phase & Performance Metrics

Counting per sample		Predicted Class	
		Class = 1	Class = 0
Actual Class	Class = 1	f_{11}	f_{10}
	Class = 0	f_{01}	f_{00}

(100% accuracy in learning often points to problems using machine learning methods in practice)

- Accuracy (usually in %)

$$Accuracy = \frac{\text{number of correct predictions}}{\text{total number of predictions}}$$

- Error rate

$$Error\ rate = \frac{\text{number of wrong predictions}}{\text{total number of predictions}}$$

Exercises – Evaluate Multi Output Perceptron Model



MNIST Dataset – A Multi Output Perceptron Model – Output

```
Epoch 7/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4419 - acc: 0.8838
Epoch 8/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4271 - acc: 0.8866
Epoch 9/20
60000/60000 [=====] - 2s 25us/step - loss: 0.4151 - acc: 0.8888
Epoch 10/20
60000/60000 [=====] - 2s 26us/step - loss: 0.4052 - acc: 0.8910
Epoch 11/20
60000/60000 [=====] - 2s 26us/step - loss: 0.3968 - acc: 0.8924
Epoch 12/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3896 - acc: 0.8944
Epoch 13/20
60000/60000 [=====] - 2s 26us/step - loss: 0.3832 - acc: 0.8956
Epoch 14/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3777 - acc: 0.8969
Epoch 15/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3727 - acc: 0.8982
Epoch 16/20
60000/60000 [=====] - 1s 24us/step - loss: 0.3682 - acc: 0.8989
Epoch 17/20
60000/60000 [=====] - 1s 25us/step - loss: 0.3641 - acc: 0.9001
Epoch 18/20
60000/60000 [=====] - 1s 25us/step - loss: 0.3604 - acc: 0.9007
Epoch 19/20
60000/60000 [=====] - 2s 25us/step - loss: 0.3570 - acc: 0.9016
Epoch 20/20
60000/60000 [=====] - 1s 24us/step - loss: 0.3538 - acc: 0.9023
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 41us/step
Test score: 0.33423959468007086
Test accuracy: 0.9101
```

✓ **Multi Output Perceptron:
~91,01% (20 Epochs)**

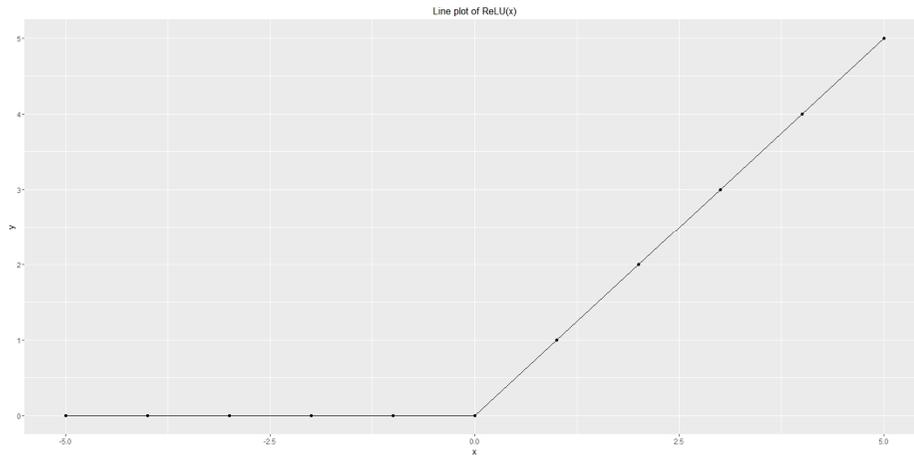
ANN – MNIST Dataset – Extend ANN Blueprint

- ✓ Data Preprocessing done (i.e. data normalization, reshape, etc.)
- ✓ Initial ANN topology existing
- ✓ Initial setup of model works (create, compile, fit)

- **Extend the neural network topology**
 - Which layers are required?
 - Think about input layer need to match the data – what data we had?
 - Maybe hidden layers?
 - How many hidden layers?
 - What activation function for which layer (e.g. maybe ReLU)?
 - Think Dense layer – Keras?
 - Think about final Activation as Softmax (cf. Day One) → output probability

Selected Activation Functions

■ Rectified Linear Unit



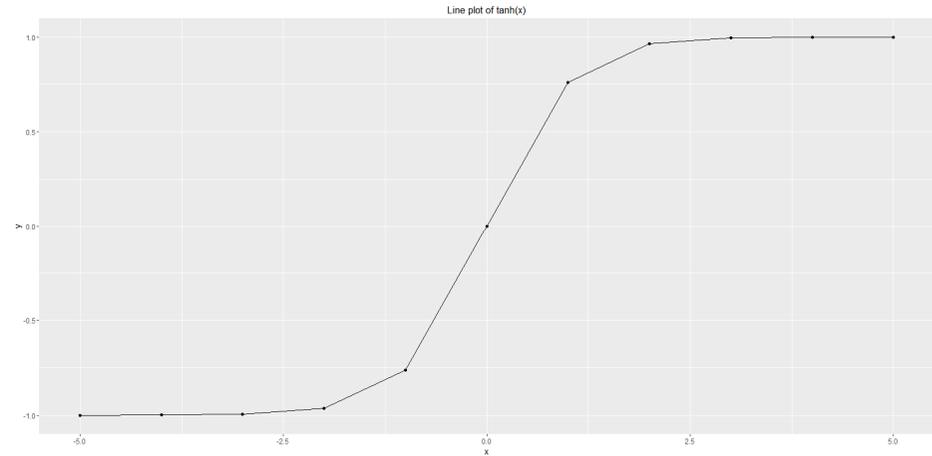
[10] *big-data.tips*,
'Relu Neural Network'

	x	y
1	-5	0
2	-4	0
3	-3	0
4	-2	0
5	-1	0
6	0	0
7	1	1
8	2	2
9	3	3
10	4	4
11	5	5



```
model.add(Dense(N_HIDDEN))  
model.add(Activation('relu'))
```

■ Tanh



[11] *big-data.tips*,
'tanh'

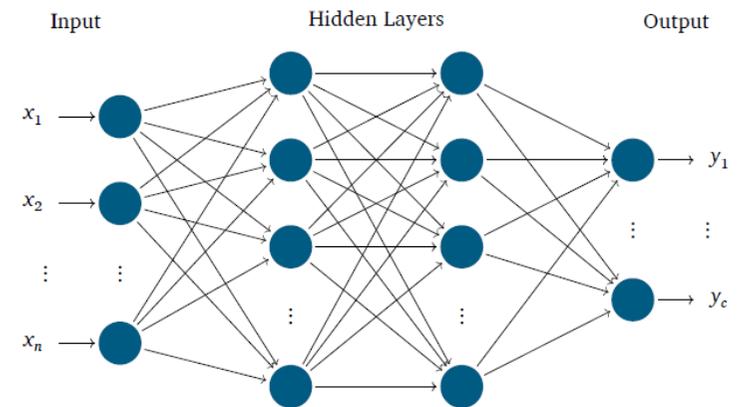
	x	y
1	-5	-0.9999092
2	-4	-0.9993293
3	-3	-0.9950548
4	-2	-0.9640276
5	-1	-0.7615942
6	0	0.0000000
7	1	0.7615942
8	2	0.9640276
9	3	0.9950548
10	4	0.9993293
11	5	0.9999092



```
model.add(Dense(N_HIDDEN))  
model.add(Activation('tanh'))
```

Exercises – Add Two Hidden Layers

✓ Multi Output Perceptron: ~91,01% (20 Epochs)



ANN – MNIST Dataset – Add Two Hidden Layers

- All parameter value remain the same as before
- We add N_HIDDEN as parameter in order to set 128 neurons in one hidden layer – this number is a hyperparameter that is not directly defined and needs to be find with parameter search

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
N_HIDDEN = 128 # number of neurons in one hidden layer
```

```
# model Keras sequential
model = Sequential()
```

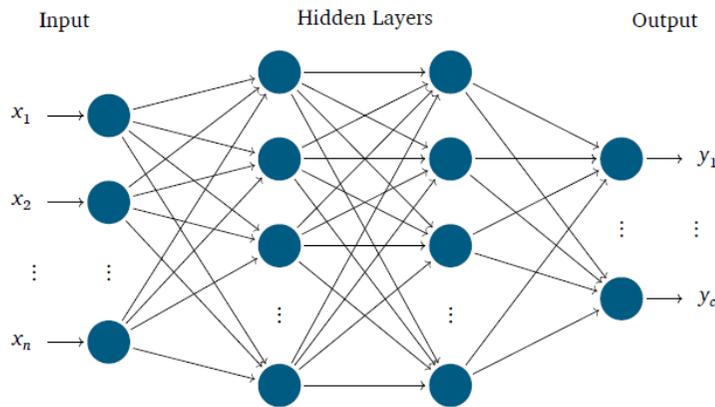
```
# modeling step
# 2 hidden layers each N_HIDDEN neurons
model.add(Dense(N_HIDDEN, input_shape=(RESHAPED,)))
model.add(Activation('relu'))
model.add(Dense(N_HIDDEN))
model.add(Activation('relu'))
model.add(Dense(NB_CLASSES))
```

```
# add activation function layer to get class probabilities
model.add(Activation('softmax'))
```

- The non-linear Activation function 'relu' represents a so-called Rectified Linear Unit (ReLU) that only recently became very popular because it generates good experimental results in ANNs and more recent deep learning models – it just returns 0 for negative values and grows linearly for only positive values
- A hidden layer in an ANN can be represented by a fully connected Dense layer in Keras by just specifying the number of hidden neurons in the hidden layer

MNIST Dataset & Model Summary & Parameters

- Added two Hidden Layers
 - Each hidden layers has 128 neurons



Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 128)	100480
activation_1 (Activation)	(None, 128)	0
dense_2 (Dense)	(None, 128)	16512
activation_2 (Activation)	(None, 128)	0
dense_3 (Dense)	(None, 10)	1290
activation_3 (Activation)	(None, 10)	0

Total params: 118,282
Trainable params: 118,282
Non-trainable params: 0

■ Relevant for validation: Choosing a model with different layers is a model selection that directly also influences the number of parameters (e.g. add Dense layer from Keras means new weights)



```
# printout a summary of the model to understand model complexity  
model.summary()
```

ANN 2 Hidden – MNIST Dataset – Output

```
Epoch 7/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2743 - acc: 0.9223
Epoch 8/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2601 - acc: 0.9266
Epoch 9/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2477 - acc: 0.9301
Epoch 10/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2365 - acc: 0.9329
Epoch 11/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2264 - acc: 0.9356
Epoch 12/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2175 - acc: 0.9386
Epoch 13/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2092 - acc: 0.9412
Epoch 14/20
60000/60000 [=====] - 1s 18us/step - loss: 0.2013 - acc: 0.9432
Epoch 15/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1942 - acc: 0.9454
Epoch 16/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1876 - acc: 0.9472
Epoch 17/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1813 - acc: 0.9487
Epoch 18/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1754 - acc: 0.9502
Epoch 19/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1700 - acc: 0.9522
Epoch 20/20
60000/60000 [=====] - 1s 18us/step - loss: 0.1647 - acc: 0.9536
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 33us/step
Test score: 0.16286438911408185
Test accuracy: 0.9514
```

- ✓ **Multi Output Perceptron:
~91,01% (20 Epochs)**
- ✓ **ANN 2 Hidden Layers:
~95,14 % (20 Epochs)**

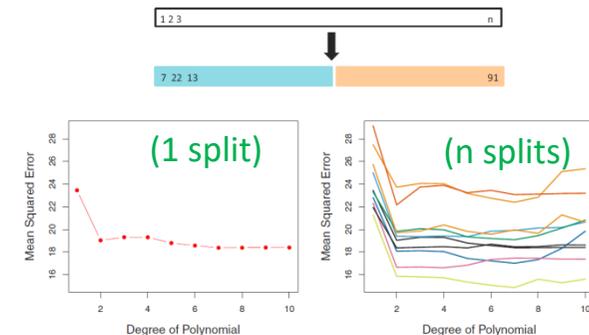
Validation & Model Selection – Terminology

- The ‘Validation technique’ should be used in all machine learning or data mining approaches
- Model assessment is the process of evaluating a models performance
- Model selection is the process of selecting the proper level of flexibility for a model

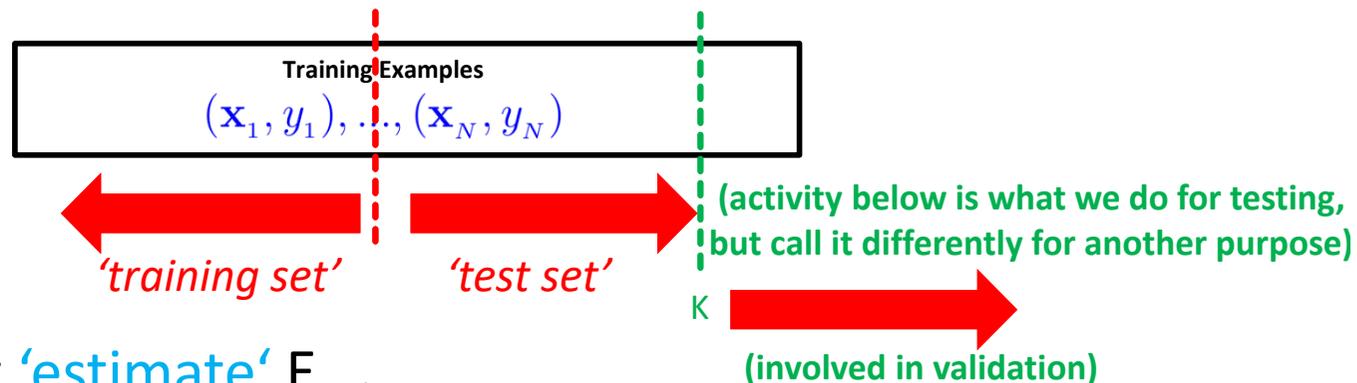
modified from [9] ‘An Introduction to Statistical Learning’

- ‘Training error’
 - Calculated when learning from data (i.e. dedicated training set)
- ‘Test error’
 - Average error resulting from using the model with ‘new/unseen data’
 - ‘new/unseen data’ was **not used in training** (i.e. dedicated test set)
 - In many practical situations, a dedicated test set is not really available
- ‘Validation Set’
 - Split data into training & validation set
- ‘Variance’ & ‘Variability’
 - Result in **different random splits** (right)

(split creates a two subsets of comparable size)

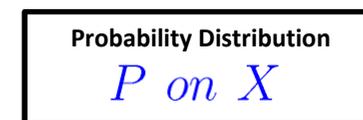


Validation Technique – Pick one point & Estimate E_{out}



- Understanding ‘estimate’ E_{out}
 - On one out-of-sample point (\mathbf{x}, y) the error is $e(h(\mathbf{x}), y)$
 - E.g. use squared error: $e(h(\mathbf{x}), f(\mathbf{x})) = (h(\mathbf{x}) - f(\mathbf{x}))^2$
 $e(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2$
 - Use this quantity as estimate for E_{out} (poor estimate)
 - Term ‘expected value’ to formalize (probability theory)

(Taking into account the theory of Lecture 1 with probability distribution on X etc.)



(aka ‘random variable’)

$$\mathbf{x} = (x_1, \dots, x_d)$$

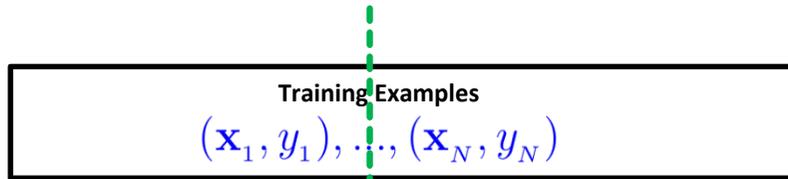
$$\mathbb{E}[e(h(\mathbf{x}), y)] = E_{out}(h) \text{ (aka the long-run average value of repetitions of the experiment)}$$

(one point as unbiased estimate of E_{out} that can have a high variance leads to bad generalization)

Validation Technique – Validation Set

- Validation set consists of data that has been not used in training to estimate true out-of-sample
- Rule of thumb from practice is to take 20% (1/5) for validation of the learning model

- Solution for **high variance** in expected values $\mathbb{E}[e(h(\mathbf{x}), y)] = E_{out}(h)$
 - Take a **‘whole set’** instead of just one point (\mathbf{x}, y) for validation



(we need points not used in training to estimate the out-of-sample performance)

(involved in training+test) K (involved in validation)

(we do the same approach with the testing set, but here different purpose)

- Idea: K data points for validation

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_K, y_K)$ (validation set)

$$E_{val}(h) = \frac{1}{K} \sum_{k=1}^K e(h(\mathbf{x})_k, y_k) \text{ (validation error)}$$

- Expected value to **‘measure’** the out-of-sample error

(expected values averaged over set)

- ‘Reliable estimate’** if K is large

$$\mathbb{E}[E_{val}(h)] = \frac{1}{K} \sum_{k=1}^K \mathbb{E}[e(h(\mathbf{x})_k, y_k)] = E_{out}$$

(on rarely used validation set, otherwise data gets contaminated)

(this gives a much better (lower) variance than on a single point given K is large)

Validation Technique – Model Selection Process

- Model selection is choosing (a) different types of models or (b) parameter values inside models
- Model selection takes advantage of the validation error in order to decide → ‘pick the best’

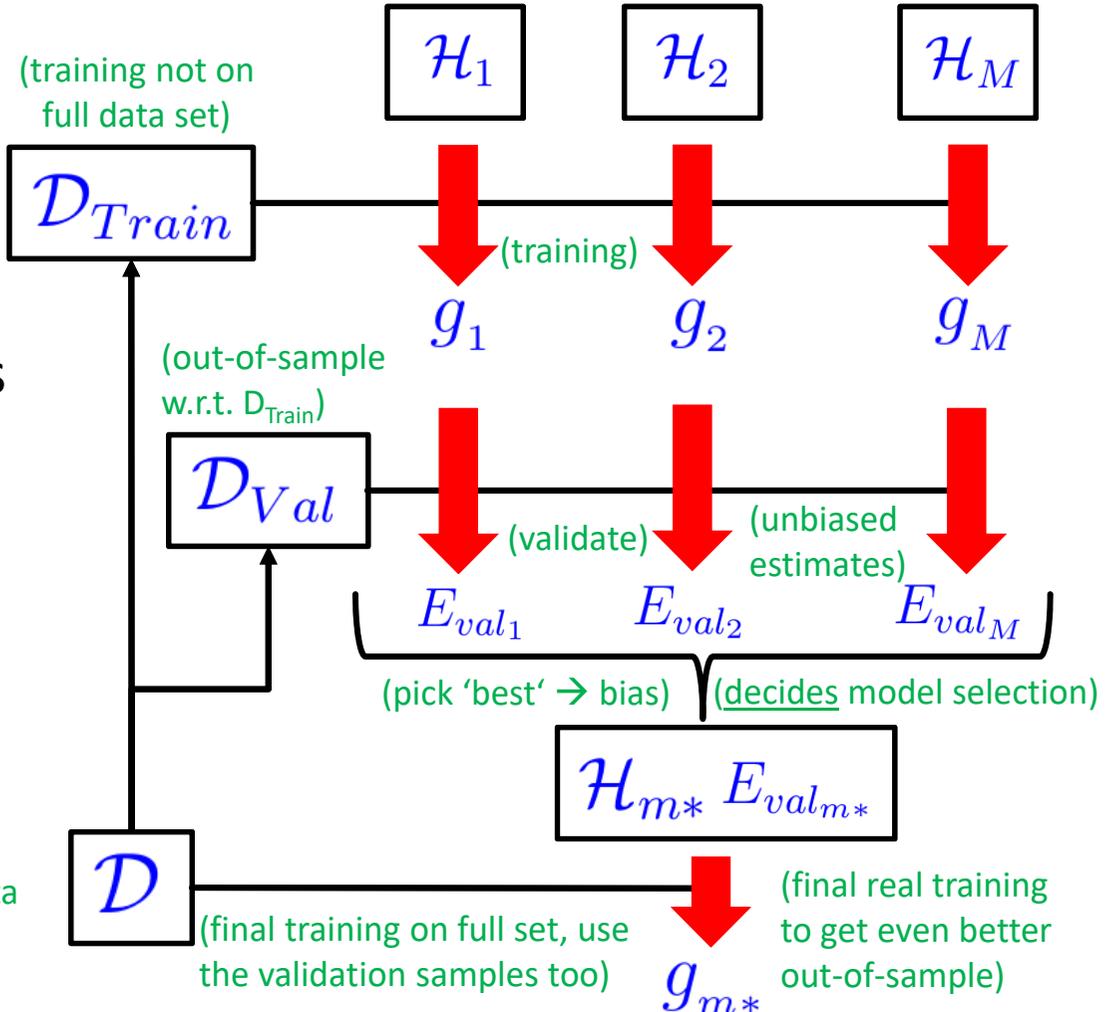
Hypothesis Set
 $\mathcal{H} = \{h\}; g \in \mathcal{H}$

(set of candidate formulas across models)

- Many different models
Use validation error to perform select decisions
- Careful consideration:
 - ‘Picked means decided’ hypothesis has already bias (→ contamination)
 - Using \mathcal{D}_{Val} M times

Final Hypothesis
 $g_{m^*} \approx f$

(test this on unseen data good, but depends on availability in practice)



Exercises – Add Validation – Table & Groups

- ✓ Multi Output Perceptron: ~91,01% (20 Epochs)
- ✓ ANN 2 Hidden Layers: ~95,14% (20 Epochs) – overfit?



VAL_SPLIT	Accuracy Groups
0.0	
0.1	
0.2	
0.3	
0.4	
0.5	

ANN 2 Hidden 1/5 Validation – MNIST Dataset

- If there is enough data available one rule of thumb is to take 1/5 (0.2) 20% of the datasets for validation only
- Validation data is used to perform model selection (i.e. parameter / topology decisions)

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
N_HIDDEN = 128 # number of neurons in one hidden layer
VAL_SPLIT = 0.2 # 1/5 for validation rule of thumb
```

- The validation split parameter enables an easy validation approach during the model training (aka fit)
- Expectations should be a higher accuracy for unseen data since training data is less biased when using validation for model decisions (check statistical learning theory)
- **VALIDATION_SPLIT**: Float between 0 and 1
- Fraction of the training data to be used as validation data
- The model fit process will set apart this fraction of the training data and will not train on it
- Instead it will evaluate the loss and any model metrics on the validation data at the end of each epoch.

```
# model training
history = model.fit(X_train, Y_train, batch_size=BATCH_SIZE, epochs=NB_EPOCH, verbose=VERBOSE, validation_split = VAL_SPLIT)
```

Train on 48000 samples, validate on 12000 samples

ANN 2 Hidden – 1/5 Validation – MNIST Dataset – Output

```
Epoch 7/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2967 - acc: 0.9148 - val_loss: 0.2759 - val_acc: 0.9212  
Epoch 8/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2825 - acc: 0.9187 - val_loss: 0.2636 - val_acc: 0.9248  
Epoch 9/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2702 - acc: 0.9222 - val_loss: 0.2550 - val_acc: 0.9272  
Epoch 10/20  
48000/48000 [=====] - 1s 17us/step - loss: 0.2593 - acc: 0.9259 - val_loss: 0.2461 - val_acc: 0.9311  
Epoch 11/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2494 - acc: 0.9283 - val_loss: 0.2367 - val_acc: 0.9338  
Epoch 12/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2403 - acc: 0.9309 - val_loss: 0.2304 - val_acc: 0.9348  
Epoch 13/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2319 - acc: 0.9334 - val_loss: 0.2228 - val_acc: 0.9392  
Epoch 14/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2242 - acc: 0.9358 - val_loss: 0.2172 - val_acc: 0.9397  
Epoch 15/20  
48000/48000 [=====] - 1s 17us/step - loss: 0.2172 - acc: 0.9381 - val_loss: 0.2105 - val_acc: 0.9418  
Epoch 16/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2103 - acc: 0.9394 - val_loss: 0.2059 - val_acc: 0.9431  
Epoch 17/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.2040 - acc: 0.9417 - val_loss: 0.2007 - val_acc: 0.9447  
Epoch 18/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.1982 - acc: 0.9432 - val_loss: 0.1949 - val_acc: 0.9473  
Epoch 19/20  
48000/48000 [=====] - 1s 18us/step - loss: 0.1926 - acc: 0.9447 - val_loss: 0.1920 - val_acc: 0.9472  
Epoch 20/20  
48000/48000 [=====] - 1s 17us/step - loss: 0.1876 - acc: 0.9464 - val_loss: 0.1866 - val_acc: 0.9499
```

```
# model evaluation  
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)  
print("Test score:", score[0])  
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 21us/step  
Test score: 0.18584023508876563  
Test accuracy: 0.9462
```

Machine Learning Challenges – Problem of Overfitting

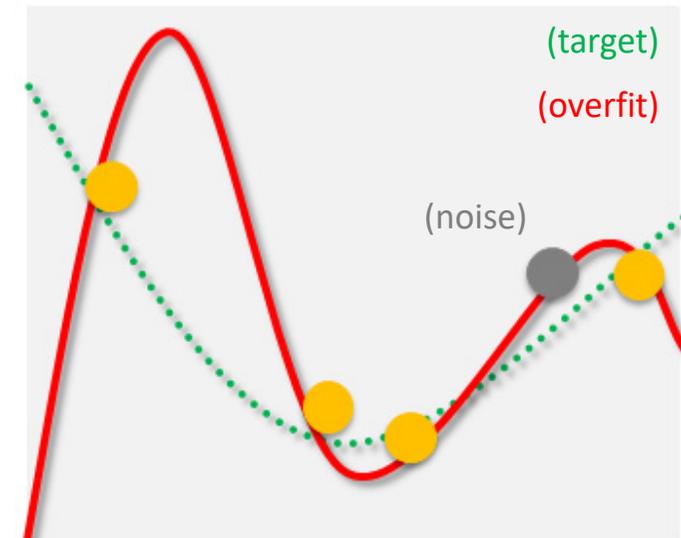
- Overfitting refers to fit the data too well – more than is warranted – thus may misguide the learning
- Overfitting is not just ‘bad generalization’ - e.g. the VC dimension covers noiseless & noise targets
- Theory of Regularization are approaches against overfitting and prevent it using different methods

- Key problem: noise in the target function leads to overfitting

- Effect: ‘noisy target function’ and its noise misguides the fit in learning
- There is always ‘some noise’ in the data
- Consequence: poor target function (‘distribution’) approximation

- Example: Target functions is second order polynomial (i.e. parabola)

- Using a higher-order polynomial fit
- Perfect fit: low $E_{in}(g)$, but large $E_{out}(g)$



(but simple polynomial works good enough)

(‘over’: here meant as 4th order, a 3rd order would be better, 2nd best)

Problem of Overfitting – Clarifying Terms

- A good model must have low training error (E_{in}) and low generalization error (E_{out})
- Model overfitting is if a model fits the data too well (E_{in}) with a poorer generalization error (E_{out}) than another model with a higher training error (E_{in})

- **Overfitting & Errors**

- $E_{in}(g)$ goes **down**

- $E_{out}(g)$ goes **up**

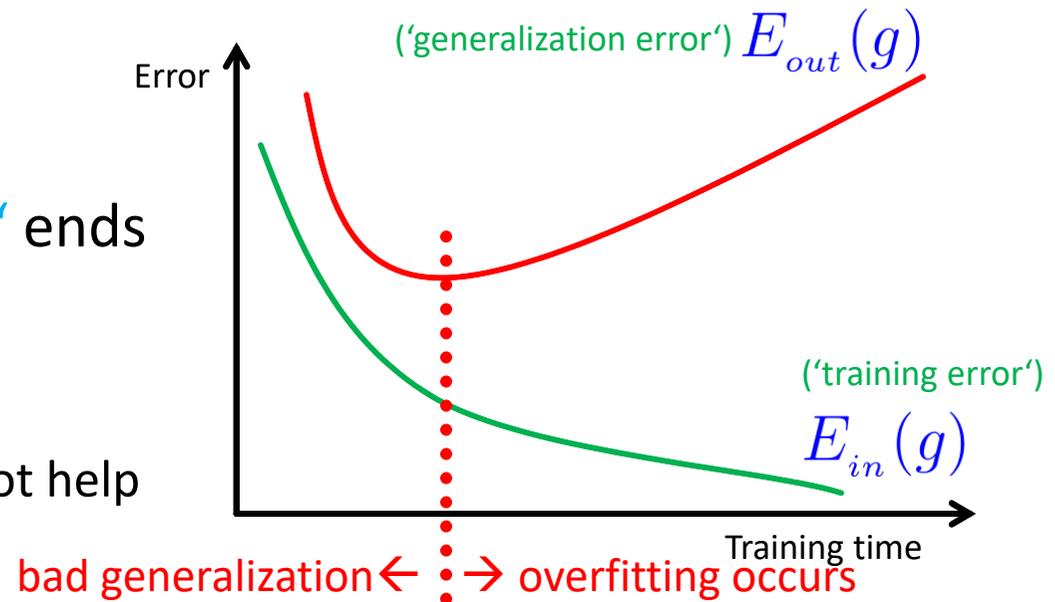
- **'Bad generalization area' ends**

- Good to reduce $E_{in}(g)$

- **'Overfitting area' starts**

- Reducing $E_{in}(g)$ does not help

- Reason **'fitting the noise'**



- The two general approaches to prevent overfitting are (1) regularization and (2) validation

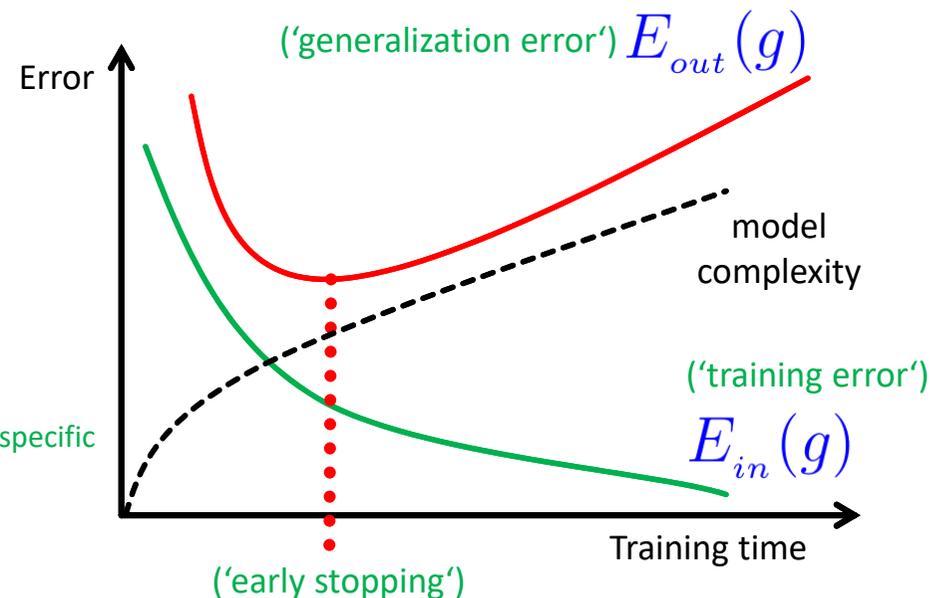
Problem of Overfitting – Model Relationships

- Review ‘overfitting situations’
 - When comparing ‘various models’ and related to ‘model complexity’
 - Different models are used, e.g. 2nd and 4th order polynomial
 - Same model is used with e.g. two different instances (e.g. two neural networks but with different parameters)

- Intuitive solution

- Detect when it happens
- ‘Early stopping regularization term’ to stop the training
- Early stopping method (later)

(‘model complexity measure: the VC analysis was independent of a specific target function – bound for all target functions’)

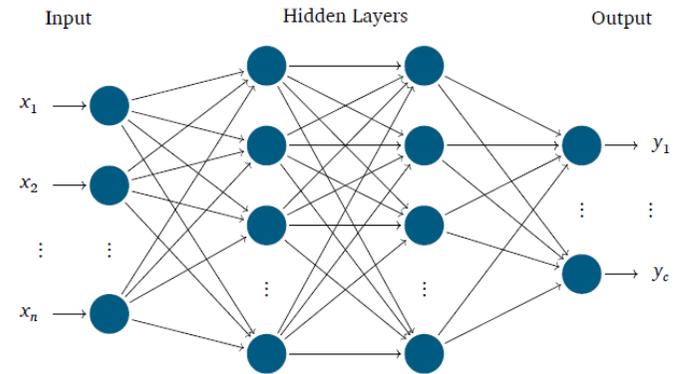


▪ ‘Early stopping’ approach is part of the theory of regularization, but based on validation methods

Problem of Overfitting – ANN Model Example

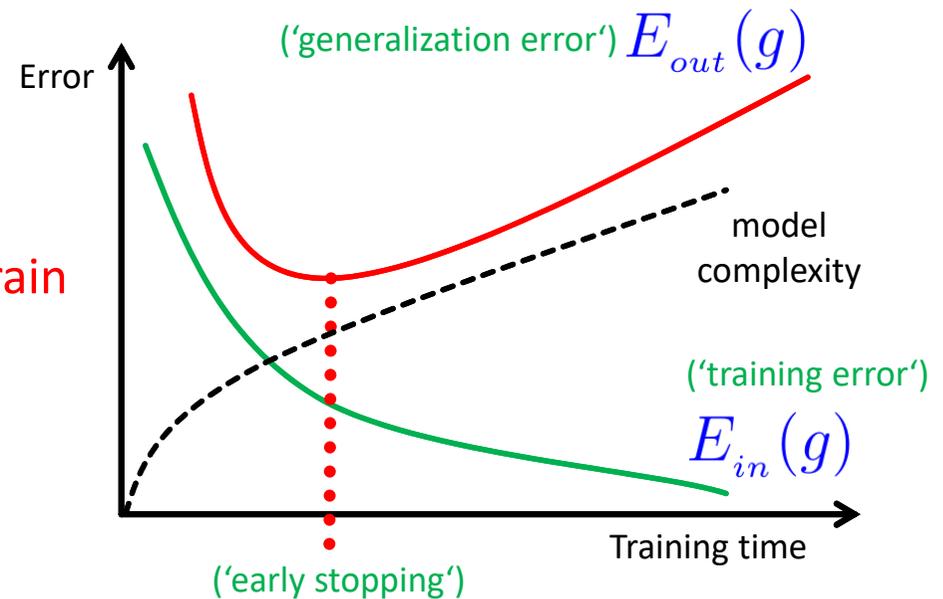
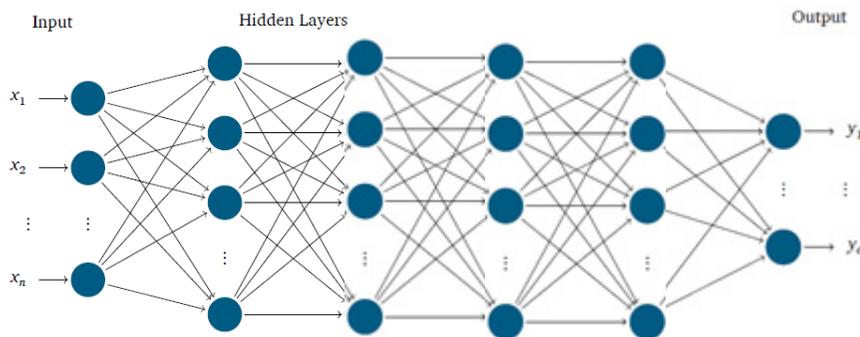
- Two Hidden Layers

- Good accuracy and works well
- Model complexity seem to match the application & data

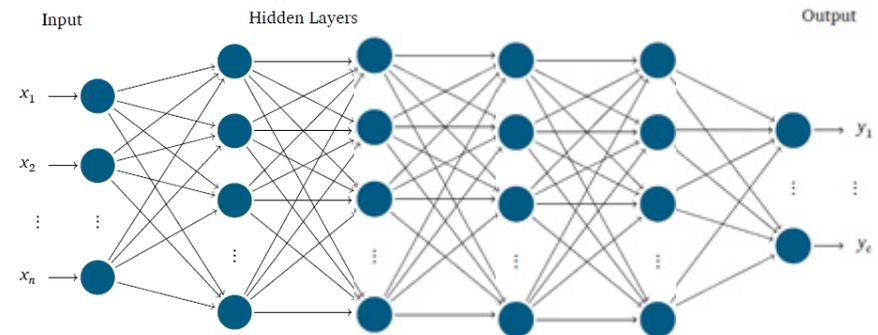
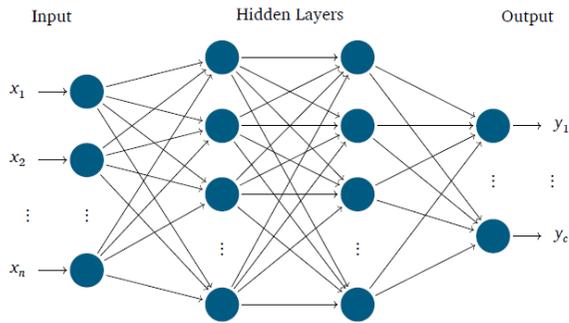


- Four Hidden Layers

- Accuracy goes down
- $E_{in}(g)$ goes down
- $E_{out}(g)$ goes up
- Significantly more weights to train
- Higher model complexity

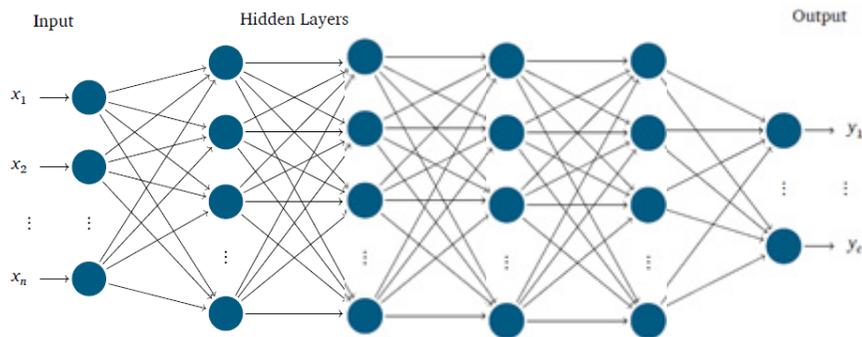


Exercises - Add more Hidden Layers – Accuracy?



MNIST Dataset & Model Summary & Parameters

- Four Hidden Layers
 - Each hidden layers has 128 neurons



Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 128)	100480
activation_1 (Activation)	(None, 128)	0
dense_2 (Dense)	(None, 128)	16512
activation_2 (Activation)	(None, 128)	0
dense_3 (Dense)	(None, 128)	16512
activation_3 (Activation)	(None, 128)	0
dense_4 (Dense)	(None, 128)	16512
activation_4 (Activation)	(None, 128)	0
dense_5 (Dense)	(None, 10)	1290
activation_5 (Activation)	(None, 10)	0

```
Total params: 151,306  
Trainable params: 151,306  
Non-trainable params: 0
```



```
# printout a summary of the model to understand model complexity  
model.summary()
```

Exercises - Add more Hidden Layers – 4 Hidden Layers

```
Epoch 7/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.2614 - acc: 0.9237 - val_loss: 0.2364 - val_acc: 0.9323  
Epoch 8/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.2431 - acc: 0.9290 - val_loss: 0.2243 - val_acc: 0.9347  
Epoch 9/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.2270 - acc: 0.9339 - val_loss: 0.2158 - val_acc: 0.9377  
Epoch 10/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.2130 - acc: 0.9385 - val_loss: 0.1995 - val_acc: 0.9427  
Epoch 11/20  
48000/48000 [=====] - 1s 23us/step - loss: 0.2001 - acc: 0.9425 - val_loss: 0.1908 - val_acc: 0.9451  
Epoch 12/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1888 - acc: 0.9445 - val_loss: 0.1866 - val_acc: 0.9464  
Epoch 13/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1783 - acc: 0.9479 - val_loss: 0.1750 - val_acc: 0.9497  
Epoch 14/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1701 - acc: 0.9507 - val_loss: 0.1675 - val_acc: 0.9529  
Epoch 15/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1615 - acc: 0.9533 - val_loss: 0.1631 - val_acc: 0.9537  
Epoch 16/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1539 - acc: 0.9555 - val_loss: 0.1553 - val_acc: 0.9555  
Epoch 17/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1469 - acc: 0.9575 - val_loss: 0.1536 - val_acc: 0.9558  
Epoch 18/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1405 - acc: 0.9590 - val_loss: 0.1505 - val_acc: 0.9560  
Epoch 19/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1351 - acc: 0.9609 - val_loss: 0.1456 - val_acc: 0.9574  
Epoch 20/20  
48000/48000 [=====] - 1s 24us/step - loss: 0.1295 - acc: 0.9625 - val_loss: 0.1398 - val_acc: 0.9600
```

```
# model evaluation  
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)  
print("Test score:", score[0])  
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 33us/step  
Test score: 0.13893915132246912  
Test accuracy: 0.9571
```

- Training accuracy should still be above the test accuracy – otherwise overfitting starts!

Exercises - Add more Hidden Layers – 6 Hidden Layers

```

Epoch 7/20
48000/48000 [=====] - 1s 28us/step - loss: 0.2567 - acc: 0.9231 - val_loss: 0.2370 - val_acc: 0.9311
Epoch 8/20
48000/48000 [=====] - 1s 28us/step - loss: 0.2333 - acc: 0.9312 - val_loss: 0.2229 - val_acc: 0.9342
Epoch 9/20
48000/48000 [=====] - 1s 28us/step - loss: 0.2141 - acc: 0.9372 - val_loss: 0.1979 - val_acc: 0.9429
Epoch 10/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1963 - acc: 0.9415 - val_loss: 0.1860 - val_acc: 0.9461
Epoch 11/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1812 - acc: 0.9470 - val_loss: 0.1779 - val_acc: 0.9487
Epoch 12/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1693 - acc: 0.9496 - val_loss: 0.1717 - val_acc: 0.9504
Epoch 13/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1580 - acc: 0.9540 - val_loss: 0.1651 - val_acc: 0.9543
Epoch 14/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1477 - acc: 0.9573 - val_loss: 0.1535 - val_acc: 0.9552
Epoch 15/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1381 - acc: 0.9594 - val_loss: 0.1461 - val_acc: 0.9577
Epoch 16/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1309 - acc: 0.9616 - val_loss: 0.1427 - val_acc: 0.9582
Epoch 17/20
48000/48000 [=====] - 1s 28us/step - loss: 0.1240 - acc: 0.9630 - val_loss: 0.1495 - val_acc: 0.9573
Epoch 18/20
48000/48000 [=====] - 1s 27us/step - loss: 0.1170 - acc: 0.9663 - val_loss: 0.1447 - val_acc: 0.9563
Epoch 19/20
48000/48000 [=====] - 1s 27us/step - loss: 0.1114 - acc: 0.9674 - val_loss: 0.1391 - val_acc: 0.9587
Epoch 20/20
48000/48000 [=====] - 1s 27us/step - loss: 0.1053 - acc: 0.9696 - val_loss: 0.1355 - val_acc: 0.9601

```

```

# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])

```

```

10000/10000 [=====] - 0s 34us/step
Test score: 0.13102742895036937
Test accuracy: 0.9614

```

▪ Training accuracy should still be above the test accuracy – otherwise overfitting starts!

AUDIENCE QUESTION

Why it does not make sense to add more and more layers?



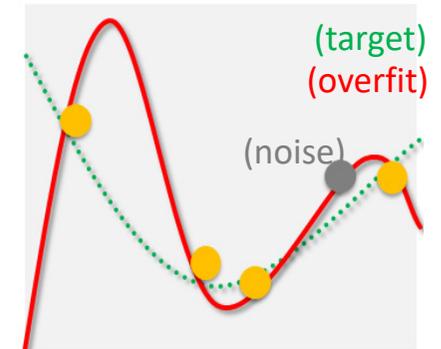
Problem of Overfitting – Noise Term Revisited

- ‘(Noisy) Target function’ is not a (deterministic) function
 - Getting with ‘same x in’ the ‘same y out’ is not always given in practice
 - Idea: Use a ‘target distribution’ instead of ‘target function’

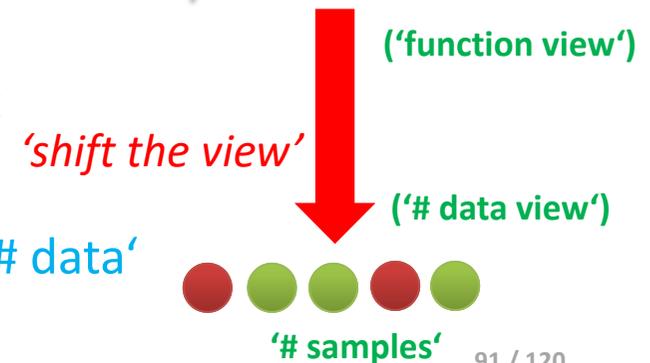
Unknown Target Distribution $P(y|x)$
target function $f : X \rightarrow Y$ plus noise
(ideal function)

■ Fitting some noise in the data is the basic reason for overfitting and harms the learning process

■ Big datasets tend to have more noise in the data so the overfitting problem might occur even more intense



- ‘Different types of some noise’ in data
 - Key to understand overfitting & preventing it
 - ‘Shift of view’: refinement of noise term
 - Learning from data: ‘matching properties of # data’



Problem of Overfitting – Stochastic Noise

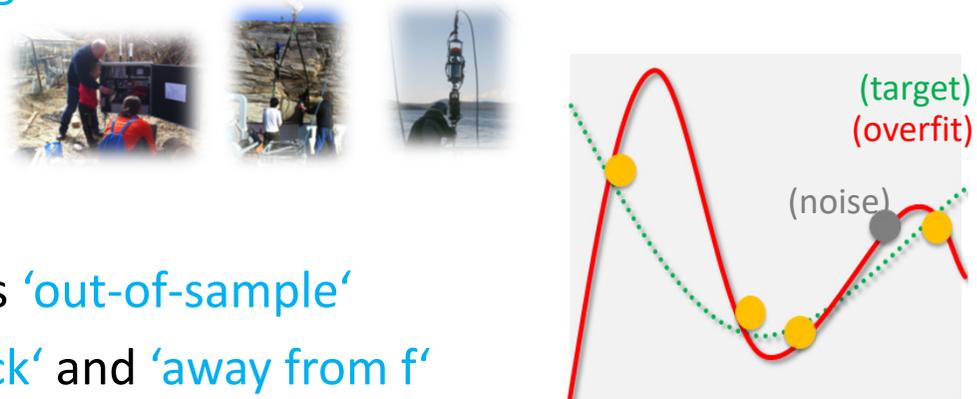
- Stochastic noise is a part ‘on top of’ each learnable function
 - Noise in the data that can not be captured and thus not modelled by f
 - Random noise : aka ‘non-deterministic noise’
 - Conventional understanding established early in this course
 - Finding a ‘non-existing pattern in noise not feasible in learning’

$$\text{target function } f : X \rightarrow Y \text{ plus noise } P(y|x)$$

(ideal function)

Practice Example

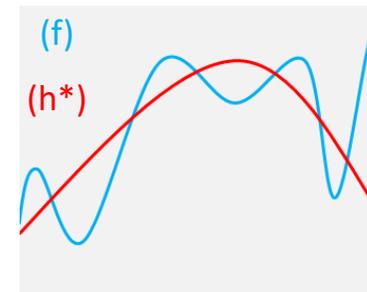
- Random fluctuations and/or measurement errors in data
- Fitting a pattern that not exists ‘out-of-sample’
- Puts learning progress ‘off-track’ and ‘away from f ’



■ Stochastic noise here means noise that can't be captured, because it's just pure 'noise as is' (nothing to look for) – aka no pattern in the data to understand or to learn from

Problem of Overfitting – Deterministic Noise

- Part of target function f that H can not capture: $f(\mathbf{x}) - h^*(\mathbf{x})$
 - Hypothesis set H is limited so best h^* can not fully approximate f
 - h^* approximates f , but fails to pick certain parts of the target f
 - ‘Behaves like noise’, existing even if data is ‘stochastic noiseless’
- Different ‘type of noise’ than stochastic noise
 - Deterministic noise depends on \mathcal{H} (determines how much more can be captured by h^*)
 - E.g. same f , and more sophisticated \mathcal{H} : noise is smaller ^{h^*} (stochastic noise remains the same, nothing can capture it)
 - Fixed for a given \mathbf{x} , clearly measurable (stochastic noise may vary for values of \mathbf{x})



(learning deterministic noise is outside the ability to learn for a given h^*)

■ Deterministic noise here means noise that can't be captured, because it is a limited model (out of the league of this particular model), e.g. ‘learning with a toddler statistical learning theory’

Problem of Overfitting – Impacts on Learning

- The higher the degree of the polynomial (cf. model complexity), the more degrees of freedom are existing and thus the more capacity exists to overfit the training data

- Understanding **deterministic noise & target complexity**
 - Increasing target complexity **increases deterministic noise** (at some level)
 - Increasing the number of data N **decreases the deterministic noise**
- **Finite N case:** \mathcal{H} tries to fit the noise
 - Fitting the noise straightforward (e.g. Perceptron Learning Algorithm)
 - **Stochastic (in data)** and **deterministic (simple model)** noise will be part of it
- **Two ‘solution methods’** for avoiding overfitting
 - **Regularization:** ‘Putting the brakes in learning’, e.g. early stopping (more theoretical, hence ‘theory of regularization’)
 - **Validation:** ‘Checking the bottom line’, e.g. other hints for out-of-sample (more practical, methods on data that provides ‘hints’)

High-level Tools – Keras – Regularization Techniques

- Keras is a high-level deep learning library implemented in Python that works on top of existing other rather low-level deep learning frameworks like Tensorflow, CNTK, or Theano
- The key idea behind the Keras tool is to enable faster experimentation with deep networks
- Created deep learning models run seamlessly on CPU and GPU via low-level frameworks

```
keras.layers.Dropout(rate,  
                     noise_shape=None,  
                     seed=None)
```

- Dropout is randomly setting a fraction of input units to 0 at each update during training time, which helps prevent overfitting (using parameter rate)

```
from keras import regularizers  
model.add(Dense(64, input_dim=64,  
               kernel_regularizer=regularizers.l2(0.01),  
               activity_regularizer=regularizers.l1(0.01)))
```

- L2 regularizers allow to apply penalties on layer parameter or layer activity during optimization itself – therefore the penalties are incorporated in the loss function during optimization



Keras

[3] *Keras Python Deep Learning Library*

Exercises – Underfitting & Add Dropout Regularizer

- Run with 20 Epochs first (not trained enough); then 200 Epochs
 - Training accuracy should be above the test accuracy – otherwise ‘underfitting’

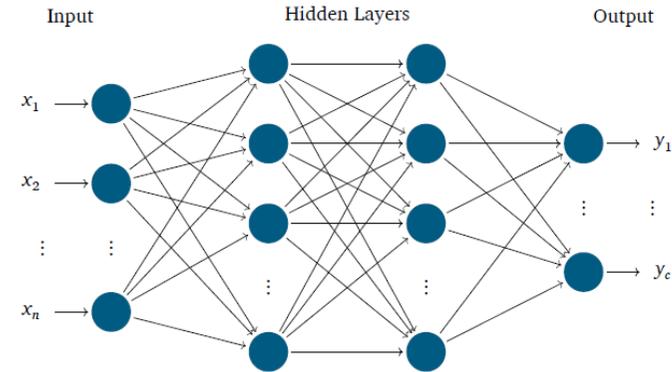


VAL_SPLIT	Dropout	Accuracy Groups
0.0	0.10	
0.1	0.20	
0.2	0.25	
0.3	0.30	
0.4	0.40	

ANN – MNIST Dataset – Add Weight Dropout Regularizer

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
N_HIDDEN = 128 # number of neurons in one hidden layer
VAL_SPLIT = 0.2 # 1/5 for validation rule of thumb
DROPOUT = 0.3 # regularization
```

```
# modeling step
# 2 hidden layers each N_HIDDEN neurons
model.add(Dense(N_HIDDEN, input_shape=(RESHAPED,)))
model.add(Activation('relu'))
model.add(Dropout(DROPOUT))
model.add(Dense(N_HIDDEN))
model.add(Activation('relu'))
model.add(Dropout(DROPOUT))
model.add(Dense(NB_CLASSES))
```



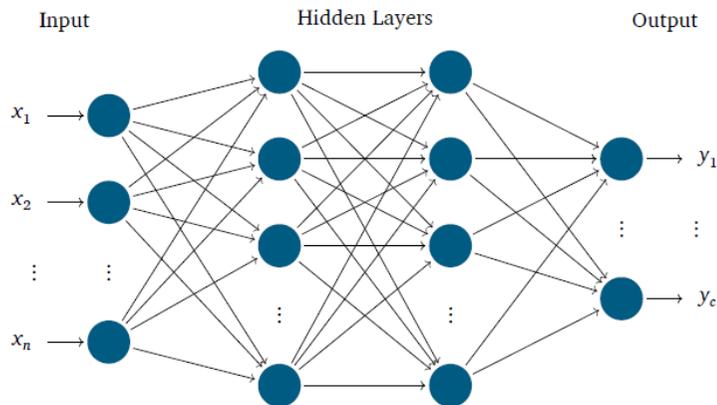
- A Dropout() regularizer randomly drops with its dropout probability some of the values propagated inside the Dense network hidden layers improving accuracy again
- Our standard model is already modified in the python script but needs to set the DROPOUT rate
- A Dropout() regularizer randomly drops with its dropout probability some of the values propagated inside the Dense network hidden layers improving accuracy again



```
model.add(Activation('relu'))
model.add(Dropout(DROPOUT))
```

MNIST Dataset & Model Summary & Parameters

- Only two Hidden Layers but with Dropout
 - Each hidden layers has 128 neurons



Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 128)	100480
activation_1 (Activation)	(None, 128)	0
dropout_1 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 128)	16512
activation_2 (Activation)	(None, 128)	0
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 10)	1290
activation_3 (Activation)	(None, 10)	0
Total params: 118,282		
Trainable params: 118,282		
Non-trainable params: 0		



```
# printout a summary of the model to understand model complexity  
model.summary()
```

ANN – MNIST – DROPOUT (20 Epochs)

```
Epoch 7/20
48000/48000 [=====] - 1s 22us/step - loss: 0.4616 - acc: 0.8628 - val_loss: 0.3048 - val_acc: 0.9127
Epoch 8/20
48000/48000 [=====] - 1s 22us/step - loss: 0.4386 - acc: 0.8688 - val_loss: 0.2896 - val_acc: 0.9172
Epoch 9/20
48000/48000 [=====] - 1s 22us/step - loss: 0.4181 - acc: 0.8762 - val_loss: 0.2776 - val_acc: 0.9198
Epoch 10/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3990 - acc: 0.8838 - val_loss: 0.2657 - val_acc: 0.9234
Epoch 11/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3819 - acc: 0.8876 - val_loss: 0.2551 - val_acc: 0.9258
Epoch 12/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3688 - acc: 0.8920 - val_loss: 0.2465 - val_acc: 0.9283
Epoch 13/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3571 - acc: 0.8943 - val_loss: 0.2388 - val_acc: 0.9299
Epoch 14/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3466 - acc: 0.8991 - val_loss: 0.2319 - val_acc: 0.9323
Epoch 15/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3359 - acc: 0.9015 - val_loss: 0.2261 - val_acc: 0.9339
Epoch 16/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3244 - acc: 0.9055 - val_loss: 0.2180 - val_acc: 0.9352
Epoch 17/20
48000/48000 [=====] - 1s 22us/step - loss: 0.3142 - acc: 0.9085 - val_loss: 0.2122 - val_acc: 0.9375
Epoch 18/20
48000/48000 [=====] - 1s 21us/step - loss: 0.3103 - acc: 0.9095 - val_loss: 0.2076 - val_acc: 0.9390
Epoch 19/20
48000/48000 [=====] - 1s 21us/step - loss: 0.3019 - acc: 0.9118 - val_loss: 0.2018 - val_acc: 0.9409
Epoch 20/20
48000/48000 [=====] - 1s 21us/step - loss: 0.2931 - acc: 0.9132 - val_loss: 0.1974 - val_acc: 0.9419
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])

10000/10000 [=====] - 0s 29us/step
Test score: 0.19944561417847873
Test accuracy: 0.9404
```

- Regularization effect not yet because too little training time (i.e. other regularization ,early stopping' here)

ANN – MNIST – DROPOUT (200 Epochs)

```
Epoch 187/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0780 - acc: 0.9755 - val_loss: 0.0810 - val_acc: 0.9764
Epoch 188/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0795 - acc: 0.9753 - val_loss: 0.0799 - val_acc: 0.9765
Epoch 189/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0774 - acc: 0.9763 - val_loss: 0.0802 - val_acc: 0.9763
Epoch 190/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0773 - acc: 0.9770 - val_loss: 0.0799 - val_acc: 0.9758
Epoch 191/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0746 - acc: 0.9771 - val_loss: 0.0804 - val_acc: 0.9762
Epoch 192/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0761 - acc: 0.9771 - val_loss: 0.0805 - val_acc: 0.9762
Epoch 193/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0750 - acc: 0.9772 - val_loss: 0.0800 - val_acc: 0.9763
Epoch 194/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0753 - acc: 0.9766 - val_loss: 0.0804 - val_acc: 0.9767
Epoch 195/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0748 - acc: 0.9768 - val_loss: 0.0799 - val_acc: 0.9767
Epoch 196/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0755 - acc: 0.9767 - val_loss: 0.0795 - val_acc: 0.9765
Epoch 197/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0740 - acc: 0.9771 - val_loss: 0.0799 - val_acc: 0.9767
Epoch 198/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0744 - acc: 0.9769 - val_loss: 0.0792 - val_acc: 0.9772
Epoch 199/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0759 - acc: 0.9769 - val_loss: 0.0794 - val_acc: 0.9767
Epoch 200/200
48000/48000 [=====] - 1s 21us/step - loss: 0.0730 - acc: 0.9778 - val_loss: 0.0794 - val_acc: 0.9771
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])

10000/10000 [=====] - 0s 27us/step
Test score: 0.07506137332450598
Test accuracy: 0.9775
```

- **Regularization effect visible by long training time using dropouts and achieving highest accuracy**
- **Note: Convolutional Neural Networks more!**

ANN – MNIST – w/o DROPOUT (200 Epochs)

```
Epoch 187/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0173 - acc: 0.9973 - val_loss: 0.0888 - val_acc: 0.9753
Epoch 188/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0170 - acc: 0.9975 - val_loss: 0.0896 - val_acc: 0.9742
Epoch 189/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0169 - acc: 0.9975 - val_loss: 0.0888 - val_acc: 0.9750
Epoch 190/200
48000/48000 [=====] - 1s 19us/step - loss: 0.0168 - acc: 0.9973 - val_loss: 0.0880 - val_acc: 0.9752
Epoch 191/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0165 - acc: 0.9977 - val_loss: 0.0884 - val_acc: 0.9747
Epoch 192/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0164 - acc: 0.9976 - val_loss: 0.0887 - val_acc: 0.9751
Epoch 193/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0162 - acc: 0.9976 - val_loss: 0.0888 - val_acc: 0.9747
Epoch 194/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0160 - acc: 0.9977 - val_loss: 0.0891 - val_acc: 0.9752
Epoch 195/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0159 - acc: 0.9977 - val_loss: 0.0889 - val_acc: 0.9752
Epoch 196/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0157 - acc: 0.9979 - val_loss: 0.0886 - val_acc: 0.9752
Epoch 197/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0155 - acc: 0.9980 - val_loss: 0.0890 - val_acc: 0.9748
Epoch 198/200
48000/48000 [=====] - 1s 19us/step - loss: 0.0153 - acc: 0.9980 - val_loss: 0.0893 - val_acc: 0.9747
Epoch 199/200
48000/48000 [=====] - 1s 19us/step - loss: 0.0152 - acc: 0.9980 - val_loss: 0.0892 - val_acc: 0.9746
Epoch 200/200
48000/48000 [=====] - 1s 20us/step - loss: 0.0151 - acc: 0.9980 - val_loss: 0.0894 - val_acc: 0.9749
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 27us/step
Test score: 0.07599342362476745
Test accuracy: 0.9764
```

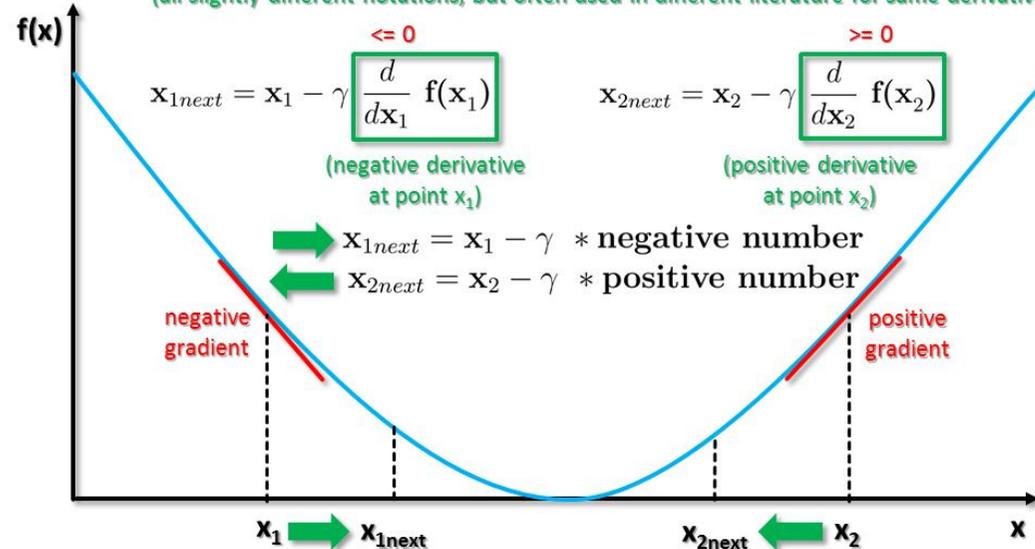
- No regularization method by long training time for comparison – slight drop in accuracy since simple dataset

MNIST Dataset & SGD Method – Revisited

- Gradient Descent (GD) uses all the training samples available for a step within a iteration
- Stochastic Gradient Descent (SGD) converges faster: only one training samples used per iteration

$$b = a - \gamma \nabla f(a) \quad b = a - \gamma \frac{\partial}{\partial a} f(a) \quad b = a - \gamma \frac{d}{da} f(a)$$

(all slightly different notations, but often used in different literature for same derivative term)



```
from keras.optimizers import SGD
```

```
OPTIMIZER = SGD() # optimization technique
```

[4] Big Data Tips,
Gradient Descent

MNIST Dataset & RMSprop & Adam Optimization Methods

- RMSProp is an advanced optimization technique that in many cases enable earlier convergence
- Adam includes a concept of momentum (i.e. velocity) in addition to the acceleration of SGD

```
Epoch 7/20
48000/48000 [=====] - 1s 25us/step - loss: 0.1127 - acc: 0.9668 - val_loss: 0.1014 - val_acc: 0.9723
Epoch 8/20
48000/48000 [=====] - 1s 25us/step - loss: 0.1051 - acc: 0.9690 - val_loss: 0.0984 - val_acc: 0.9735
Epoch 9/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0970 - acc: 0.9706 - val_loss: 0.0996 - val_acc: 0.9747
Epoch 10/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0949 - acc: 0.9716 - val_loss: 0.0958 - val_acc: 0.9754
Epoch 11/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0880 - acc: 0.9734 - val_loss: 0.0945 - val_acc: 0.9763
Epoch 12/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0873 - acc: 0.9745 - val_loss: 0.0957 - val_acc: 0.9761
Epoch 13/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0842 - acc: 0.9745 - val_loss: 0.0952 - val_acc: 0.9757
Epoch 14/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0804 - acc: 0.9763 - val_loss: 0.1002 - val_acc: 0.9767
Epoch 15/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0788 - acc: 0.9771 - val_loss: 0.0991 - val_acc: 0.9772
Epoch 16/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0756 - acc: 0.9772 - val_loss: 0.0988 - val_acc: 0.9761
Epoch 17/20
48000/48000 [=====] - 1s 25us/step - loss: 0.0758 - acc: 0.9776 - val_loss: 0.1033 - val_acc: 0.9753
Epoch 18/20
48000/48000 [=====] - 1s 26us/step - loss: 0.0755 - acc: 0.9781 - val_loss: 0.0996 - val_acc: 0.9773
Epoch 19/20
48000/48000 [=====] - 1s 26us/step - loss: 0.0725 - acc: 0.9784 - val_loss: 0.1055 - val_acc: 0.9764
Epoch 20/20
48000/48000 [=====] - 1s 26us/step - loss: 0.0712 - acc: 0.9791 - val_loss: 0.1014 - val_acc: 0.9778
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [=====] - 0s 33us/step
Test score: 0.09596708530617616
Test accuracy: 0.9779
```



```
from keras.optimizers import RMSprop
```

```
OPTIMIZER = RMSprop() # optimization technique
```

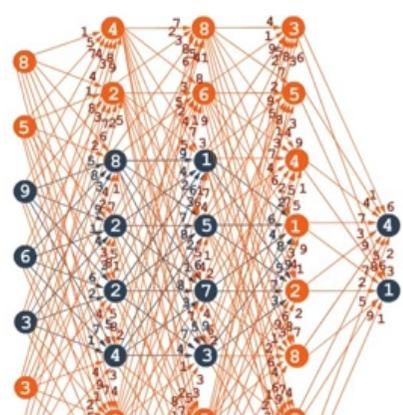
Exercises – Underfitting & Change to Adam

- Run with 20 Epochs With Adam Optimizer



[Video] Overfitting in Deep Neural Networks

Causes and Outcomes



will assign weights to features that are not needed and will add unnecessary complexity

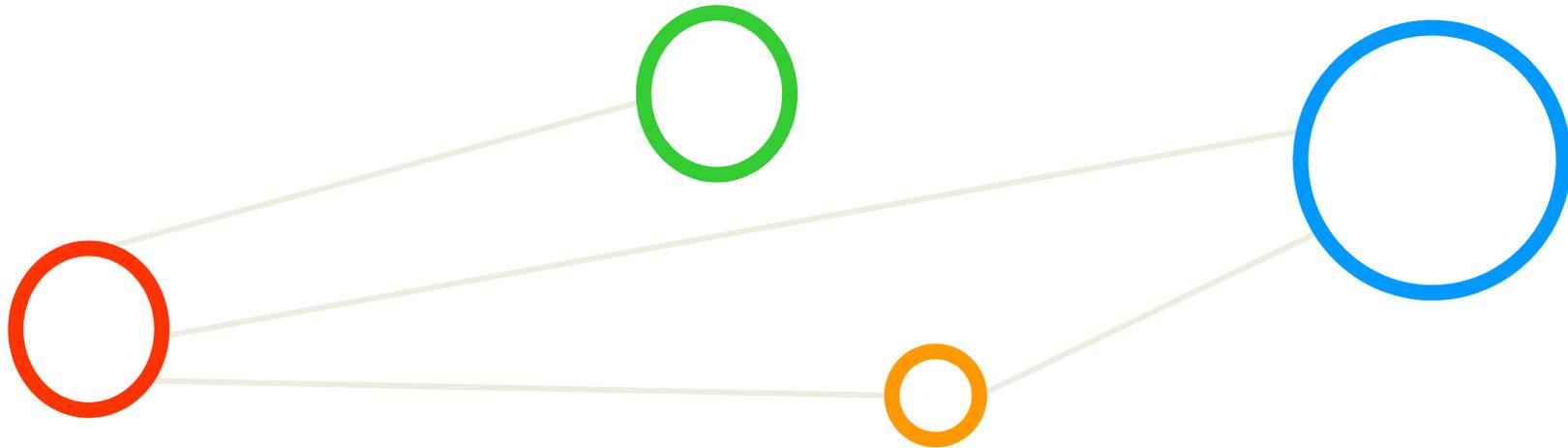
TODSI

3:42 / 5:38

CC

[12] *How good is your fit?, YouTube*

Appendix A: CRISP-DM Process

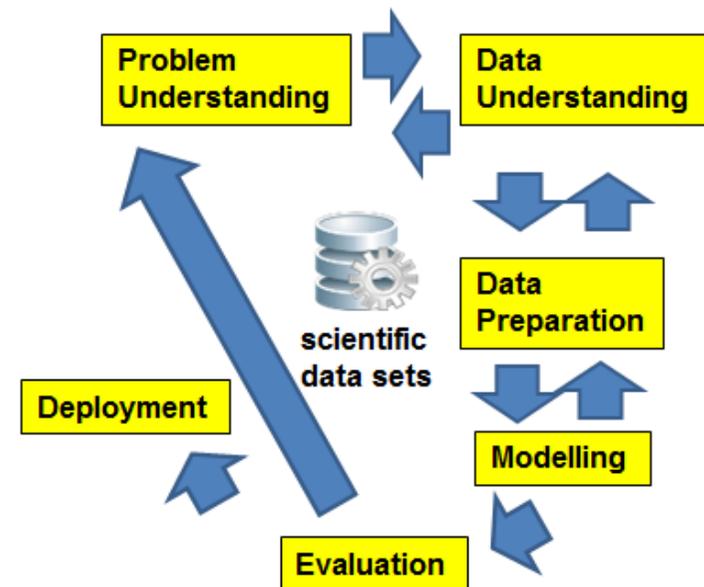


Summary: Systematic Process

- Systematic data analysis guided by a ‘standard process’
 - Cross-Industry Standard Process for Data Mining (CRISP-DM)

- A data mining project is guided by these six phases:
 - (1) Problem Understanding;
 - (2) Data Understanding;
 - (3) Data Preparation;
 - (4) Modeling;
 - (5) Evaluation;
 - (6) Deployment

- Lessons Learned from Practice
 - Go back and forth between the different six phases



[7] C. Shearer, CRISP-DM model, Journal Data Warehousing, 5:13

1 – Problem (Business) Understanding

- The Business Understanding phase consists of four distinct tasks: (A) Determine Business Objectives; (B) Situation Assessment; (C) Determine Data Mining Goal; (D) Produce Project Plan

[6] CRISP-DM User Guide

- **Task A – Determine Business Objectives**
 - Background, Business Objectives, Business Success Criteria
- **Task B – Situation Assessment**
 - Inventory of Resources, Requirements, Assumptions, and Constraints
 - Risks and Contingencies, Terminology, Costs & Benefits
- **Task C – Determine Data Mining Goal**
 - Data Mining Goals and Success Criteria
- **Task D – Produce Project Plan**
 - Project Plan
 - Initial Assessment of Tools & Techniques

2 – Data Understanding

- The Data Understanding phase consists of four distinct tasks:
(A) Collect Initial Data; (B) Describe Data; (C) Explore Data; (D) Verify Data Quality

[6] CRISP-DM User Guide

- Task A – Collect Initial Data
 - Initial Data Collection Report
- Task B – Describe Data
 - Data Description Report
- Task C – Explore Data
 - Data Exploration Report
- Task D – Verify Data Quality
 - Data Quality Report

3 – Data Preparation

- The Data Preparation phase consists of six distinct tasks: (A) Data Set; (B) Select Data; (C) Clean Data; (D) Construct Data; (E) Integrate Data; (F) Format Data

[6] CRISP-DM User Guide

- Task A – Data Set
 - Data set description
- Task B – Select Data
 - Rationale for inclusion / exclusion
- Task C – Clean Data
 - Data cleaning report
- Task D – Construct Data
 - Derived attributes, generated records
- Task E – Integrate Data
 - Merged data
- Task F – Format Data
 - Reformatted data

4 – Modeling

- The Data Preparation phase consists of four distinct tasks: (A) Select Modeling Technique; (B) Generate Test Design; (C) Build Model; (D) Assess Model;

[6] CRISP-DM User Guide

- **Task A – Select Modeling Technique**
 - Modeling assumption, modeling technique
- **Task B – Generate Test Design**
 - Test design
- **Task C – Build Model**
 - Parameter settings, models, model description
- **Task D – Assess Model**
 - Model assessment, revised parameter settings

5 – Evaluation

- The Data Preparation phase consists of three distinct tasks: (A) Evaluate Results; (B) Review Process; (C) Determine Next Steps

[6] CRISP-DM User Guide

- **Task A – Evaluate Results**
 - Assessment of data mining results w.r.t. business success criteria
 - List approved models
- **Task B – Review Process**
 - Review of Process
- **Task C – Determine Next Steps**
 - List of possible actions, decision

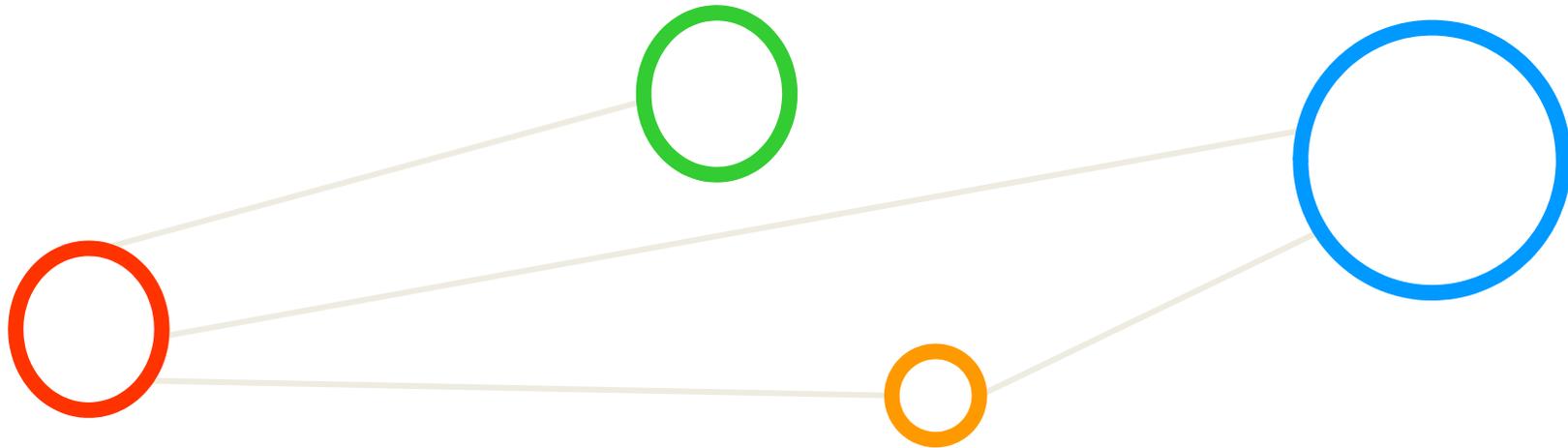
6 – Deployment

- The Data Preparation phase consists of three distinct tasks: (A) Plan Deployment; (B) Plan Monitoring and Maintenance; (C) Produce Final Report; (D) Review Project

[6] CRISP-DM User Guide

- **Task A – Plan Deployment**
 - Establish a deployment plan
- **Task B – Plan Monitoring and Maintenance**
 - Create a monitoring and maintenance plan
- **Task C – Product Final Report**
 - Create final report and provide final presentation
- **Task D – Review Project**
 - Document experience, provide documentation

Appendix B – SSH Commands JURECA



Appendix B – SSH Commands JURECA

- `salloc --gres=gpu:4 --partition=gpus --nodes=1 --account=training1911 --time=00:30:00 --reservation=intro-dl-wed`
- `module --force purge;`
`module use /usr/local/software/jureca/OtherStages`
`module load Stages/Devel-2018b GCCcore/.7.3.0`
`module load TensorFlow/1.12.0-GPU-Python-3.6.6`
`module load Keras/2.2.4-GPU-Python-3.6.6`
- `srun python PYTHONSCRIPTNAME`

Appendix B – JURECA Login Screen

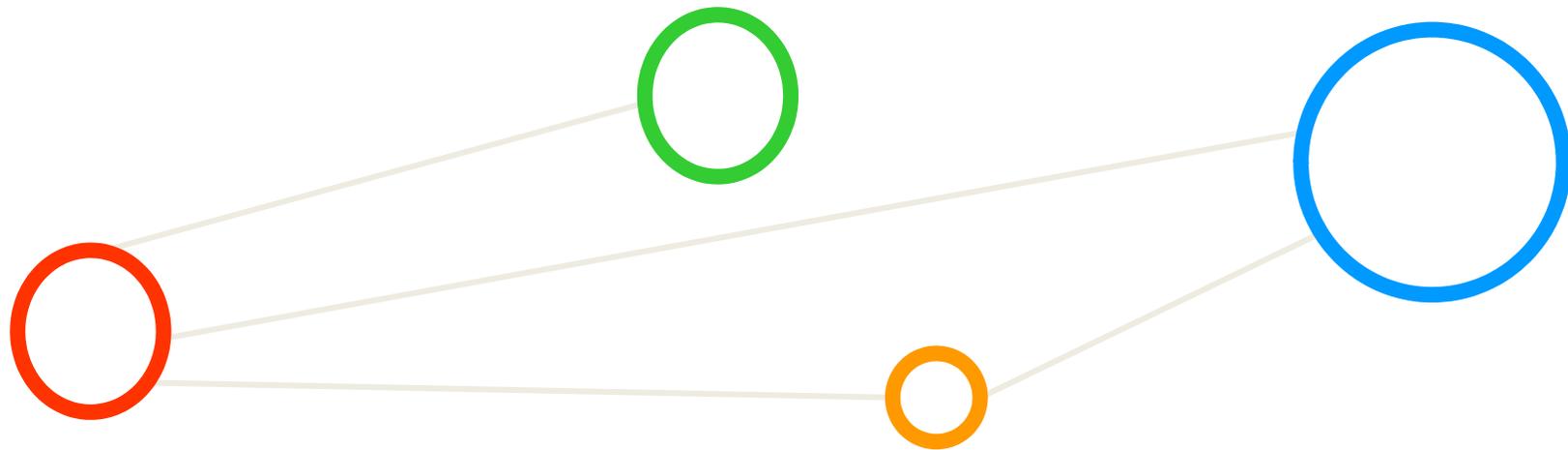
```
2. jureca.fz-juelich.de (riedel1)
• MobaXterm 11.0 •
(SSH client, X-server and networking tools)

> SSH session to riedell@jureca.fz-juelich.de
• SSH compression : ✓
• SSH-browser      : ✓
• X11-forwarding  : ✓ (remote display is forwarded through SSH)
• DISPLAY         : ✓ (automatically set on remote server)

> For more info, ctrl+click on help or visit our website

Last login: Wed Feb 27 09:16:51 2019 from pool-230-44-zam2168.wlan.kfa-juelich.de
*****
*                               Welcome to JURECA                               *
*                                                                           *
* Information about the system, latest changes, user documentation and FAQs: *
*                               http://www.fz-juelich.de/ias/jsc/jureca        *
*****
*                               ### Known Issues ###                         *
*                                                                           *
* An up-to-date list of known issues on the system is maintained at         *
*                               http://www.fz-juelich.de/ias/jsc/jureca-known-issues *
*                                                                           *
*                               2019-03-11 *
*****
*                               ### New software modules ###                 *
*                                                                           *
* Since 14th of May 2019, a new default modules stage (2019a) is available *
* with the latest versions of software and compilers.                       *
*                                                                           *
* Users who like to keep the old stage (2018b) are advised to perform the   *
* following adaptation in their job scripts:                                 *
*   1. Cluster: ml use /usr/local/software/jureca/OtherStages                *
*      Booster: ml use /usr/local/software/jurecabooster/OtherStages         *
*   2. Both:   ml Stages/2018b                                               *
*                                                                           *
*                               2019-05-14 *
*****
[riedell@jrl05 ~]$
```

Lecture Bibliography



Lecture Bibliography (1)

- [1] Lego Bricks Images,
Online: https://tucsonbotanical.org/wp-content/uploads/2015/08/You-build-it_image.jpg
- [2] F. Rosenblatt, 'The Perceptron--a perceiving and recognizing automaton',
Report 85-460-1, Cornell Aeronautical Laboratory, 1957
- [3] Keras Python Deep Learning Library,
Online: <https://keras.io/>
- [4] Big Data Tips, 'Gradient Descent',
Online: <http://www.big-data.tips/gradient-descent>
- [5] YouTube Video, 'Neural Networks, A Simple Explanation',
Online: http://www.youtube.com/watch?v=gck_5x2KsLA
- [6] Pete Chapman, 'CRISP-DM User Guide', 1999,
Online: <http://lyle.smu.edu/~mhd/8331f03/crisp.pdf>
- [7] C. Shearer, CRISP-DM model, Journal Data Warehousing, 5:13
- [8] Morris Riedel, 'Introduction to Machine Learning Algorithms', Invited YouTube Lecture, six lectures
University of Ghent, 2017,
Online: <https://www.youtube.com/watch?v=KgiuUZ3WeP8&list=PLrmNhuZo9sgbcWtMGN0i6G9HEvh08JG0J>
- [9] An Introduction to Statistical Learning with Applications in R,
Online: <http://www-bcf.usc.edu/~gareth/ISL/index.html>
- [10] www.big-data.tips, 'Relu Neural Network'
Online: <http://www.big-data.tips/relu-neural-network>

Lecture Bibliography (2)

- [11] www.big-data.tips, 'tanh',
Online: <http://www.big-data.tips/tanh>
- [12] YouTube Video, 'How good is your fit? - Ep. 21 (Deep Learning SIMPLIFIED)',
Online: <https://www.youtube.com/watch?v=cJA5IHIL30>

