

Introduction to Data Analysis and Machine Learning in Physics:

3. Machine Learning Basics

Martino Borsato, Jörg Marks, Klaus Reygers

Studierendentage, 11-14 April 2023

Books / Material

- Books on Machine Learning

- Ian Goodfellow and Yoshua Bengio and Aaron Courville, MIT Press Book
Deep Learning
free online <http://www.deeplearningbook.org/>
- Aurelien Geron, Hands-On Machine Learning with Scikit-Learn and TensorFlow
- Francois Chollet, Deep Learning with Python
- Deep Learning for Physics Research, Martin Erdmann, Jonas Glombitza, Gregor Kasieczka, Uwe Klemradt (available via Kindle)

- Papers

- A high-bias, low-variance introduction to Machine Learning for physicists
<https://arxiv.org/abs/1803.08823>
- Machine learning and the physical sciences
<https://arxiv.org/abs/1903.10563>

- Credits: The talk contains material of a tutorial at the 3rd IML workshop 2019 at CERN by Yannik Rath and the Train-the-Trainer Workshop "Deep Learning Basics" in February 2023

Introduction

- ChatGPT's opinion on the relation of Artificial Intelligence (AI) Machine Learning (ML) and Deep Learning

AI (Artificial Intelligence) is a broad field that encompasses a range of techniques and approaches for developing intelligent systems. Machine Learning (ML) is a subset of AI that focuses on teaching machines to learn from data and make predictions or decisions without being explicitly programmed. Deep Learning (DL) is a subset of ML that uses artificial neural networks with many layers to learn increasingly abstract features from data, and has achieved significant breakthroughs in areas such as computer vision, natural language processing, and robotics. Therefore, DL is a specific type of ML, and ML is a specific type of AI.

Question: Where are you in this hierarchy?

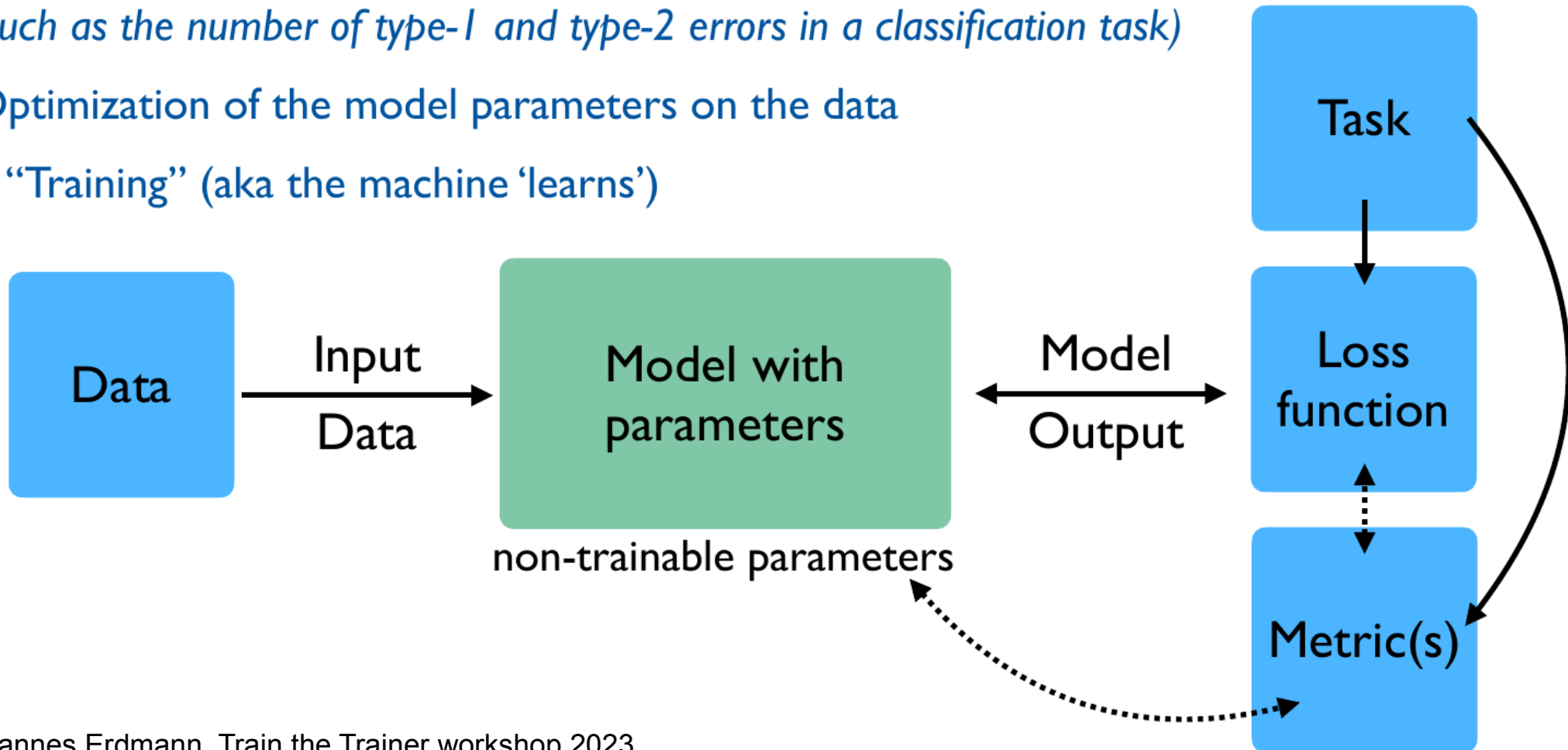
As an AI language model, I am built using various ML techniques, such as deep learning and natural language processing. Therefore, I am a subset of AI, specifically falling under the category of ML-based AI.

- Stable Diffusion is a latent text-to-image diffusion model capable of generating photo-realistic images given any text input



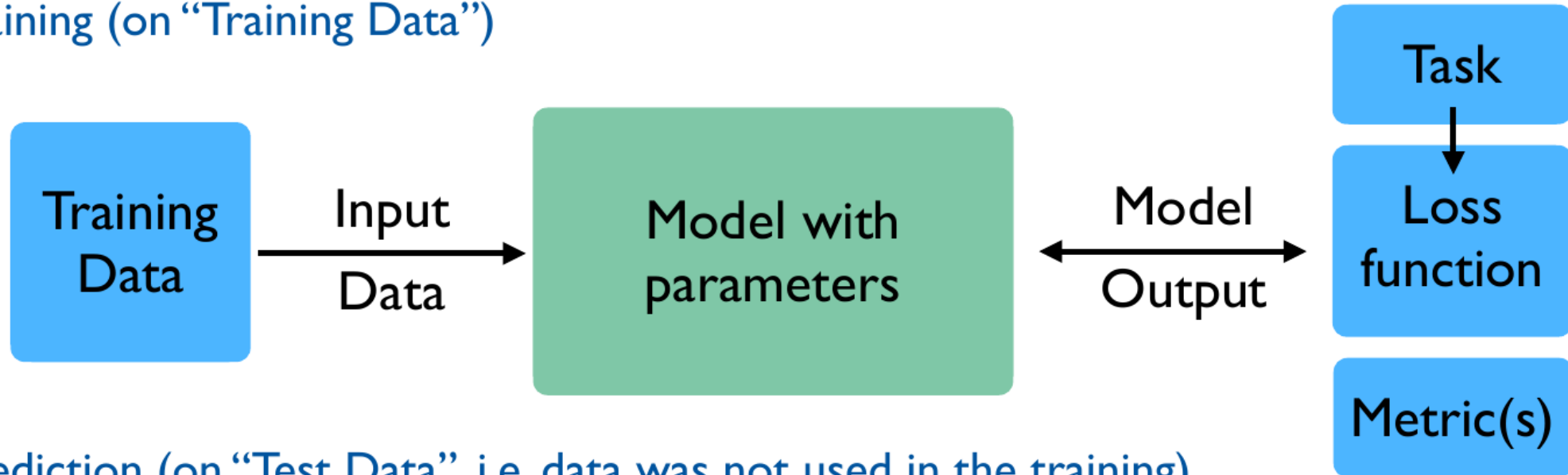
Introduction Machine Learning

- Data + Task + “a mathematical formulation of what means ‘better’ ” *
(such as the number of type-1 and type-2 errors in a classification task)
- Optimization of the model parameters on the data
= “Training” (aka the machine ‘learns’)



Introduction Machine Learning

- Training (on “Training Data”)



- Prediction (on “Test Data”, i.e. data was not used in the training)



Categories of Machine Learning

- Machine Learning types are closely related to data available for the task

- Supervised Learning

Here the model is trained on a labeled dataset, each data has an associated target variable or label. The goal of the model is to learn a mapping between the input features and the output label → make accurate predictions on new, unseen data.

- Unsupervised Learning

The model is trained on unlabeled data. The goal of the model is to find patterns and structure in the data, such as clusters, associations, or dependencies.

- Reinforcement Learning

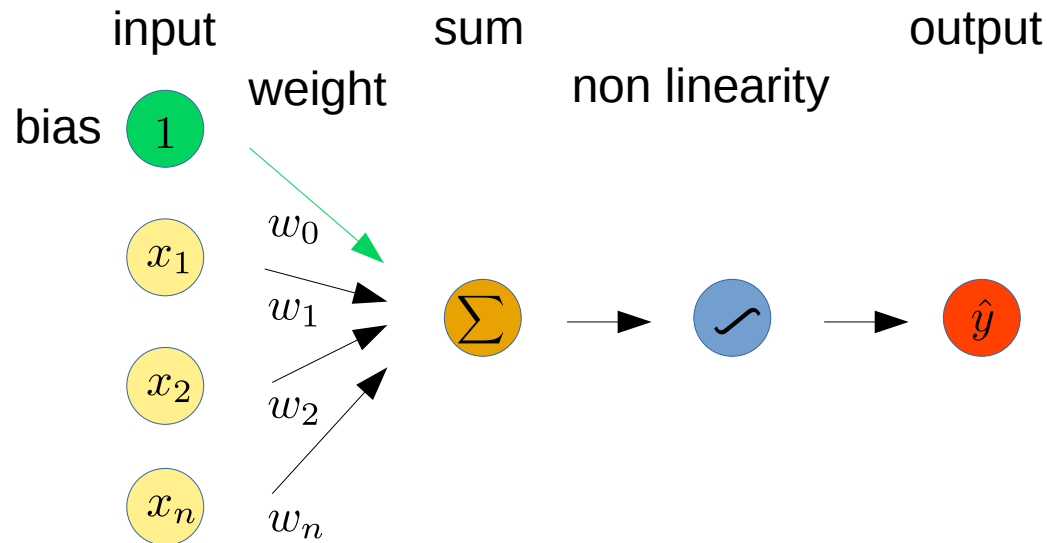
The model learns to make decisions based on feedback it receives. The goal is to maximize a reward signal (a score or a profit)

- Transfer Learning

A technique in machine learning where a pre-trained model is used as a starting point for a new task, rather than training a new model from scratch. Typically a pre-trained model has already learned useful features from a large dataset, that is then used for a new task with only a small amount of additional training data.

Multilayer Perceptron

- Forward propagating perceptron



$$\hat{y} = \theta \cdot (w_0 + \vec{x}^T \vec{W})$$

The bias term shifts the activation independent of the input.

Idea: minimize a cost function to determine the weights and bias of the perceptron by comparing the predicted output and the true output.

→ obtain a relation between the multidimensional input parameter space and some output

The input of the perceptron is passed through a non linear activation function, which determines the output of the perceptron. There is only a contribution to the perceptron if a threshold value is reached.

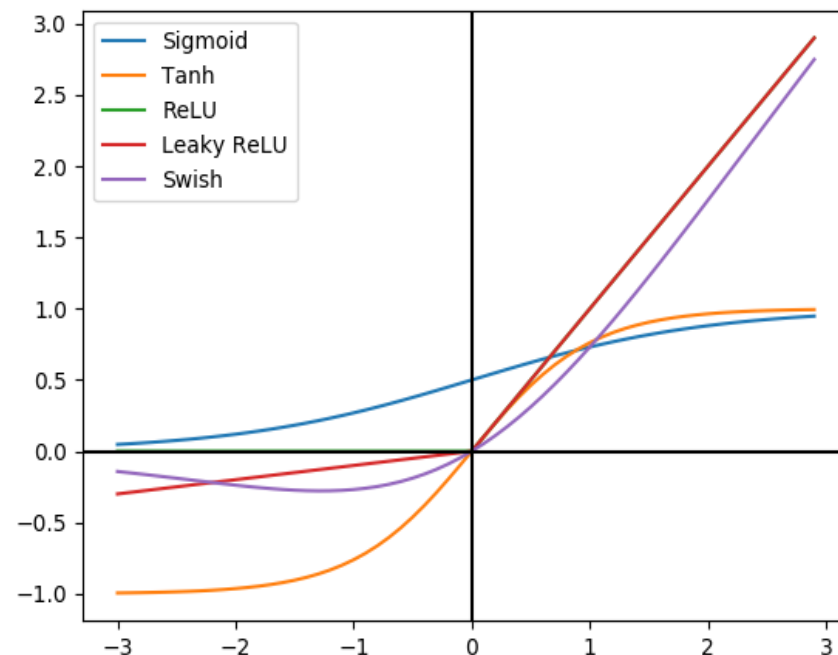
Multilayer Perceptron

- Activation function

- Activation functions are applied to each node of a neural network
- Introduce non linearities into the network → allows to approximate complex shapes

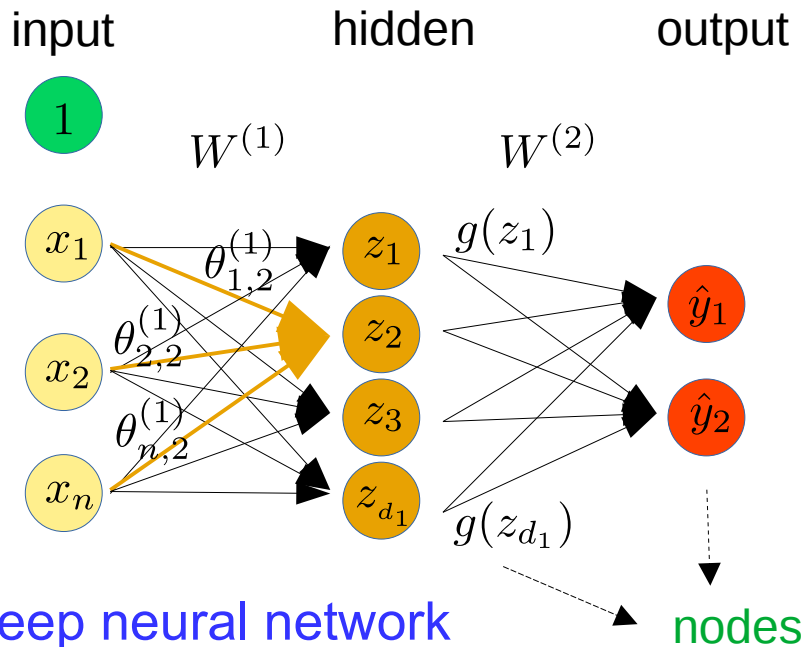
activation.py

ACTIVATION FUNCTION	EQUATION	RANGE
Linear Function	$f(x) = x$	$(-\infty, \infty)$
Step Function	$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$\{0, 1\}$
Sigmoid Function	$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$	$(0, 1)$
Hyperbolic Tanjant Function	$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$(-1, 1)$
ReLU	$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$
Leaky ReLU	$f(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$
Swish Function	$f(x) = 2x\sigma(\beta x) = \begin{cases} \beta = 0 & \text{for } f(x) = x \\ \beta \rightarrow \infty & \text{for } f(x) = 2\max(0, x) \end{cases}$	$(-\infty, \infty)$



Deep Neural Networks (DNN)

- Single layer neural network



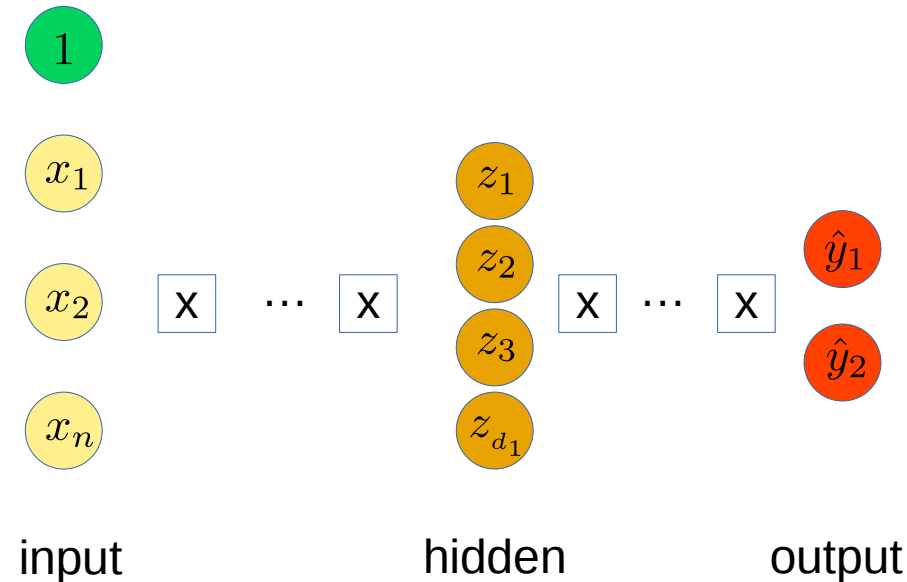
2nd element hidden layer 1 :

$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^n x_j w_{j,2}^{(1)}$$

ith output :

$$\hat{y}_i = g(w_{0,2}^{(2)} + \sum_{j=1}^{d_1} z_j w_{i,j}^{(2)})$$

- Deep neural network

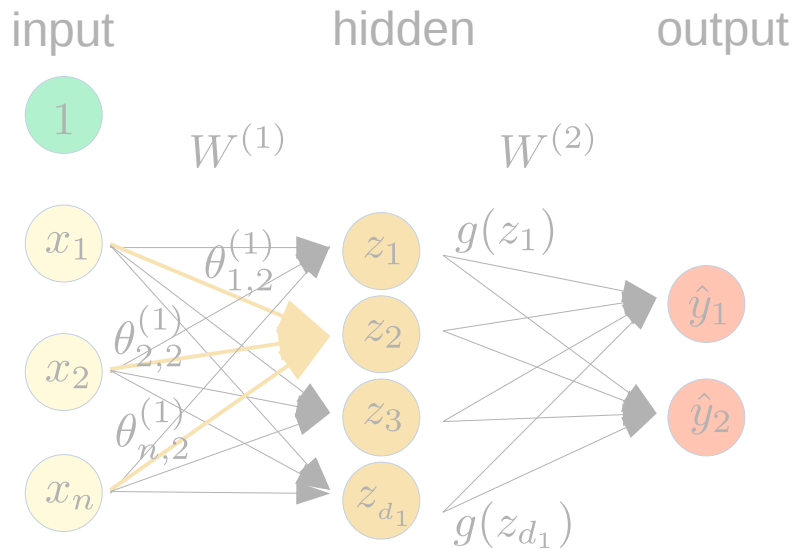


ith element hidden layer k :

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{i,j}^{(k)}$$

Deep Neural Networks (DNN)

- Single layer neural network



- Have to apply activation functions on nodes in each hidden layer and the final output layer

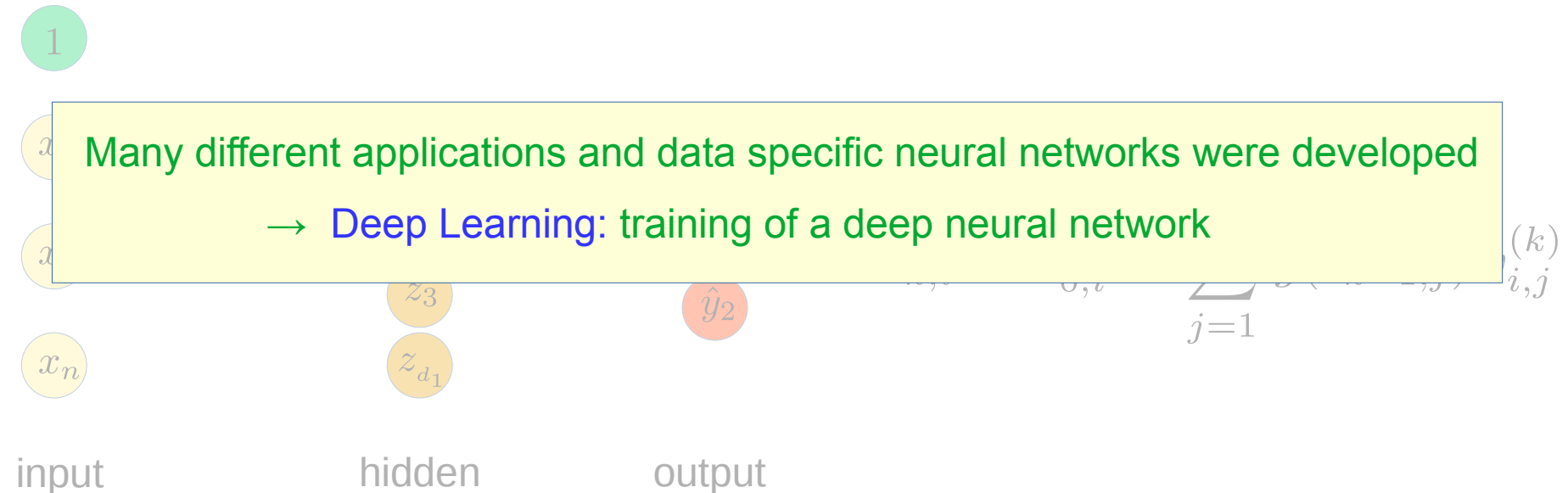
2nd element hidden layer 1 :

$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^n x_j w_{j,2}^{(1)}$$

$$z_i \rightarrow \theta(z_i)$$

Non-linear activation functions are really the key

- Deep neural network



Deep Neural Networks (DNN)

- Deep learning

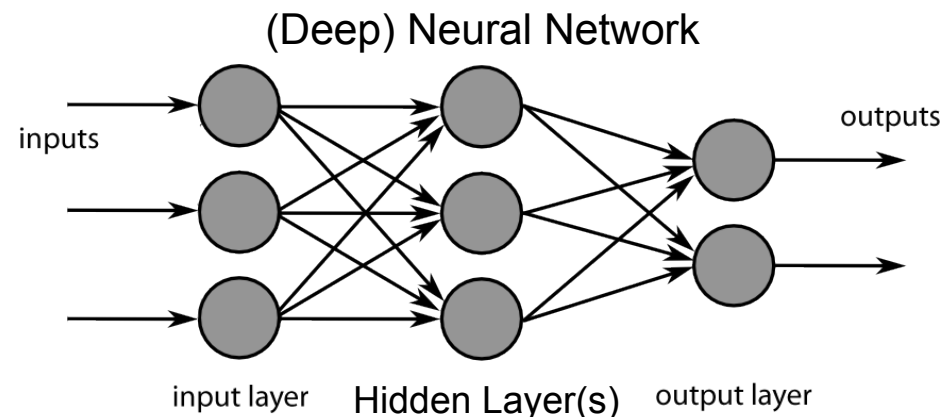
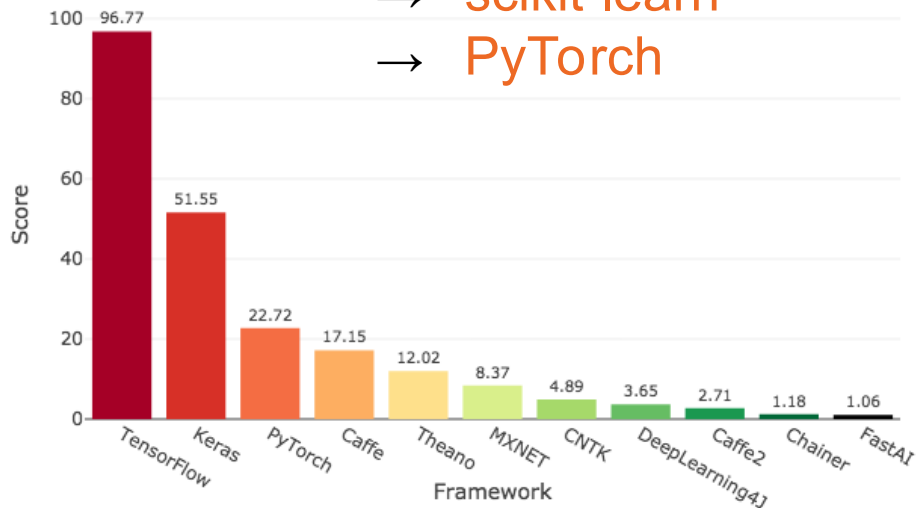
- Part of a broader family of machine learning methods based on artificial neural networks that use multiple layers to progressively extract higher level features from raw input

- Deep neural network

- Network with an input layer, at least a hidden layer and an output layer
- Each layer performs specific types of sorting and ordering in a process that some refer to as “feature hierarchy”
- Deal with unlabeled or unstructured data
- Algorithms are called **deep** if the input data is passed through a series of hidden layers with nonlinearities (**nonlinear transformations**) before it becomes output.

- Most Deep Learning frameworks (user interface) are based on Python

- TensorFlow and Keras are one of the most popular frameworks
- scikit-learn
- PyTorch



Scikit learn



- Scikit-learn is an open-source machine learning library for Python <https://scikit-learn.org/stable/>

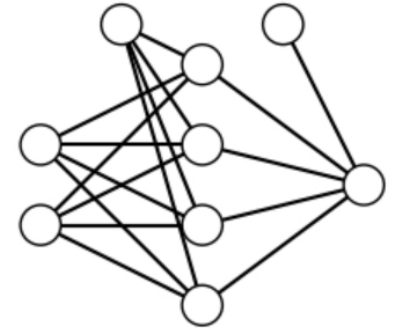
Current version 1.2.2

- Various classification, regression and clustering algorithms
Many algorithms are included: k-nearest neighbors, multi-layer, perceptrons, support vector machines, random forests, gradient boosting, k-means
- Dimensionality reduction with PCA, feature selection, non-negative matrix factorization
- Model selection with comparing, validating and choosing parameters and models.
- Feature extraction and normalization
- Examples from various areas are available via the web site

Example Python Code of a Neural Network

- A simple feed forward neural network with one hidden layer is programmed in pure python without using machine learning libraries

[03_ml_basics_simple_neural_network.ipynb](#)

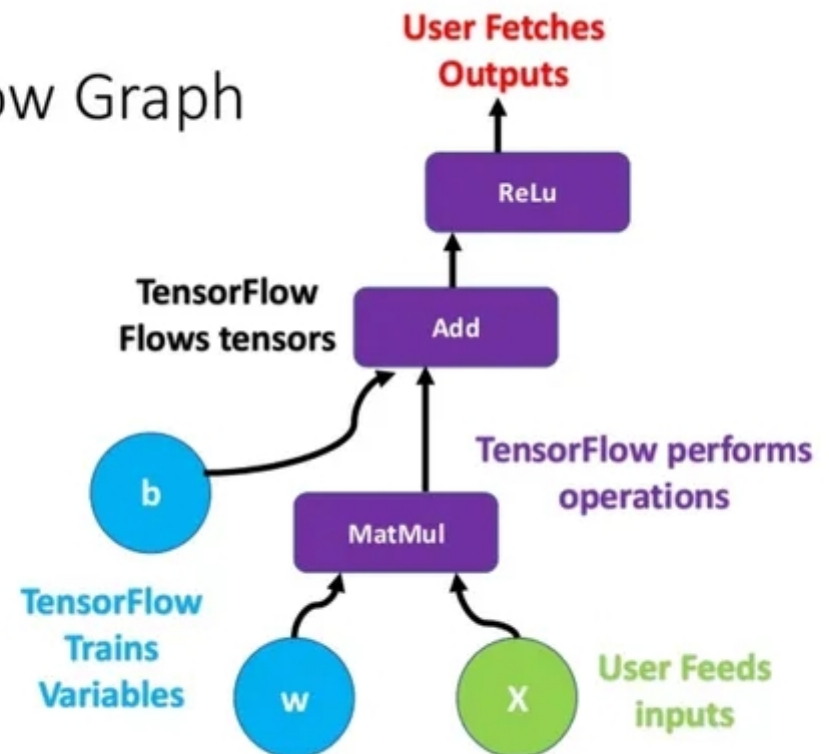


Introduction to Tensorflow

TensorFlow (TF), developed by the Google Brain team, is an open source software library for numerical computation using data flow graphs

- Tensors are multidimensional data, e.g a rank-0 tensor is a scalar, a rank-1 tensor is a vector, a rank-2 tensor is a matrix,
- Tensors represent in deep learning algorithms input data, model parameters and intermediate activations of the model during training and inference. They can be created and manipulated using TensorFlow operations, and are passed between operations
- TF implements standard mathematical operations on tensors, as well as many operations specialized for machine learning
- TF runs on CPU, GPU and TPU
- TF handles datasets, splitting data in test and training samples and compiles models
- TF implements automatic differentiation (autodiff) to compute loss functions
- Processes training loops and keeps track of weights

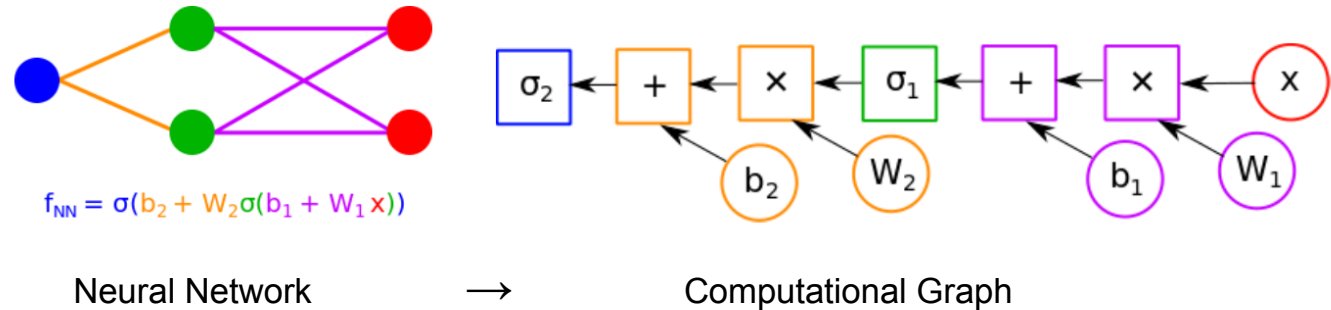
TensorFlow Graph



Computational Graphs

- Computational graphs are the basic concept of TensorFlow (TF)

Graph: S. Wunsch, IML workshop, 2019



- Nodes in the graph represent mathematical operations
 - Edges are represented by multidimensional data arrays (tensors) which communicate between the nodes
 - The nodes take 1 or more input tensors, does the computation and produces 1 or more output tensor which are then passed on to the next nodes
 - TensorFlow can optimize the graph by reordering and merging operations, eliminating redundant computations, and parallelizing operations
 - Computational graph can be modified dynamically during runtime, allowing to build models that can adapt to different inputs and perform different computations depending on the input
 - Computational graph can be visualized using tools like TensorBoard, to understand the structure of your model
- Eager execution \rightarrow Calculations are performed on demand (as other python code)

Introduction to Tensorflow

- **Tensors**

for details, see <https://www.tensorflow.org/guide/tensor>

```
rank_0_tensor = tf.constant(4) # int32 tensor
rank_1_tensor = tf.constant([2.0, 3.0, 4.0]) # float32 tensor
rank_3_tensor = tf.constant([ # int32 tensor
    [[0, 1, 2, 3, 4], # shape (3,2,5)
     [5, 6, 7, 8, 9]],
    [[10, 11, 12, 13, 14],
     [15, 16, 17, 18, 19]],
    [[20, 21, 22, 23, 24],
     [25, 26, 27, 28, 29]]])
```

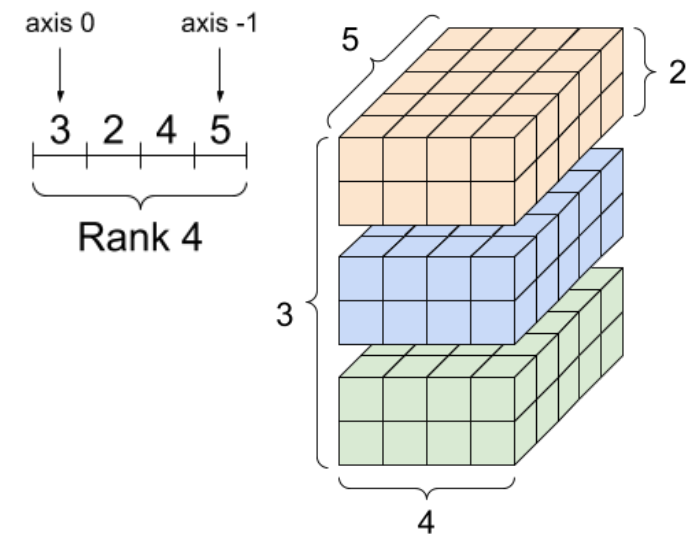
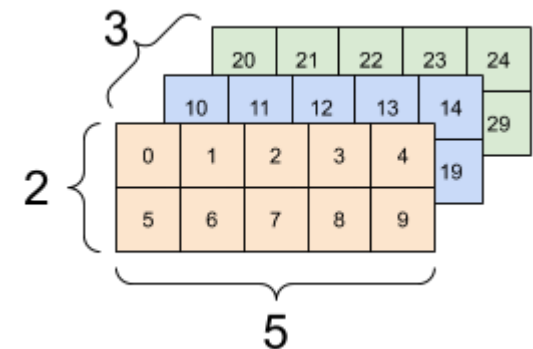
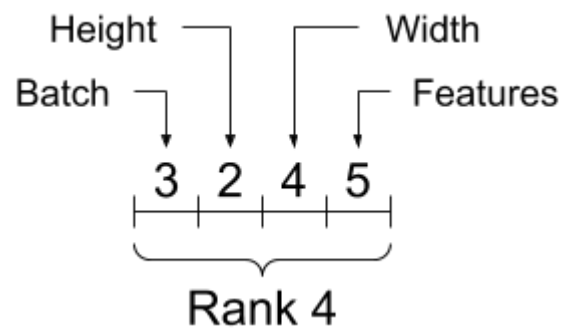
- can contain also complex and string
- can be created with arbitrary dimensions

```
rank_4_tensor = tf.zeros([3, 2, 4, 5])
```

- axis are referred to by their indices

important:

keep in mind the meaning



Introduction to Tensorflow

- **Tensor operations**

- **Export and import tensors to numpy arrays**

```
tensor = tf.constant([[1,2,3], [4,5,6], [7,8,9]]) # rank 3 tensor
arr = tensor.numpy() # numpy array

arr = np.array([[1,2,3], [4,5,6], [7,8,9]]) # numpy array
tensor = tf.constant(arr) # imported tensor
```

- **Accessing parts of tensors via numpy**

```
rank_2_tensor[1, 1].numpy() # access single tensor elements
rank_2_tensor[-1, :].numpy() # last row
rank_2_tensor[1:, :].numpy() # skip first row
rank_2_tensor[:, 1].numpy() # 2nd column
```

- **Reshaping of tensors, a new tensor is created pointing to the same data**

```
rank_3_tensor = tf.reshape(rank_3_tensor, [-1])
-1 means, shape it to whatever fits → in this case a 1 D tensor is created
tf.Tensor([ 0  1  2  3  4..... 26 27 28 29], shape=(30,), dtype=int32)
tf.reshape(rank_3_tensor, [3*2, 5]) # reshaping to (3x2)x5
tf.reshape(rank_3_tensor, [3, -1]) # reshaping to 3x(2x5)
Axis are swaped using tf.transpose
```

- **Tensors are rectangular structures, but there is `tf.RaggedTensor` (a non rectangular structure) and `SparseTensor` (only few elements are nonzero)**

Introduction to Tensorflow

- Tensor operations

- Basic math on tensors

```
tf.add(a,b) and tf.multiply(a,b)      # elementwise operation
tf.matmul(a,b)                        # matrix multiplication
tf.reduce_sum(a)                      # Sum all elements in the tensor
tf.square(a)                          # Square all elements
tf.divide(b, a)                       # Element-wise division
tf.negative(a)                        # Element-wise negation
tf.abs(tf.constant([-1, 2, -3, 4], dtype=tf.float32))
                                         # Element-wise absolute value
```

```
c = tf.constant([[4.0, 5.0], [10.0, 1.0]])
tf.reduce_max(c)                       # Find the largest value
tf.math.argmax(c)                     # Find the index of the largest value
tf.nn.softmax(c)                      # Compute softmax from tf.nn module
```

```
rank_3_tensor.dtype                   # type of the elements    dtype:'int32'
rank_3_tensor.ndim                    # number of axes          3
rank_3_tensor.shape                   # shape of tensor        (3, 2, 5)
rank_3_tensor.shape[0]                # elements along axis 0 of tensor
rank_3_tensor.shape[-1]               # elem along the last axis of tensor
tf.size(rank_3_tensor).numpy()        # total number of elements
```

Introduction to Tensorflow

- Tensor operations

- Broadcasting is a feature in TensorFlow that allows operations to be performed on **tensors of different shapes**, so that there is no need to manually reshape tensors to match each other.

Broadcasting works by replicating tensor elements along one or more dimensions to match the shape of the other tensor(s) in the operation.

Broadcasting rule:

- If the two tensors have the same rank, they are compatible if their shapes are equal in each dimension.
- if the two tensors have different ranks, the smaller tensor is broadcast to match the shape of the larger tensor. The smaller tensor is pre-padded with 1s on the left until it has the same rank as the larger tensor.
- If the two tensors have different shapes and cannot be broadcasted according to rules 1 and 2, an error is raised.

- Example of broadcasting in TensorFlow:

[03_ml_basics_tf_broadcasting.ipynb](#)

```
import tensorflow as tf
a=tf.constant([[1,2,3], [4,5,6]])      # Define 2 tensors with different shapes
b=tf.constant([10, 20, 30])
c = a * b                               # Perform element-wise multiplication using broadcasting
print(c)                                [[ 10  40  90]
                                       [ 40 100 180]],shape=(2, 3),dtype=int32)
```

Introduction to Tensorflow

- Variables in TF

- Variables are created and tracked via the `tf.Variable` class, they can be modified after initialization

```
tensor = tf.constant([[1.0, 2.0], [3.0, 4.0]])
variable = tf.Variable(tensor)           # same type as the initialization
variable.numpy()                        # convert to numpy
tf.convert_to_tensor(variable)          # convert to a tensor
tf.math.argmax(variable)                # index of highest value
```

```
b = tf.Variable([1, 2, 3])              # creates a variable tensor
b.assign_add([1, 1, 1])                 # and then modify its values
a = tf.Variable([2.0, 3.0])            # create variable a
b = tf.Variable(a)                     # Create b based on value of a
a.assign([5, 6])                       # assign a new value to a
a.assign_add([2,3]).numpy()            # add element wise to a
a.assign_sub([2,3]).numpy()            # subtract element wise from a
```

```
a = tf.Variable([[1.0, 2.0, 3.0], [4.0, 5.0, 6.0]])
b = tf.Variable([[1.0, 2.0, 3.0]])
k = a*b                                # element-wise multiply
```

```
a = tf.Variable([[1.0, 2.0, 3.0], [4.0, 5.0, 6.0]])
b = tf.constant([[1.0, 2.0], [3.0, 4.0], [5.0, 6.0]])
c = tf.matmul(a, b)                   # new tensor
```

Introduction to Tensorflow

- Gradients and differentiation in TF

To differentiate automatically, TF needs to remember what operations happen in which order during the forward pass. In backward pass the list of operations is done in reverse order. `GradientTape.gradient(target, sources)` is used to calculate the gradient of some target (often a loss) relative to some source.

- Example with a scalar

```
x = tf.Variable(3.0) # variable = 3
with tf.GradientTape() as tape:
    y = x**2
dy_dx = tape.gradient(y, x) # dy = 2x * dx
dy_dx.numpy() # convert to a number = 6.0
```

- Example with tensors

[03_ml_basics_tf_differentiate.ipynb](#)

```
w = tf.Variable(tf.random.normal((3, 2)), name='w')
b = tf.Variable(tf.zeros(2, dtype=tf.float32), name='b')
x = [[1., 2., 3.]]
with tf.GradientTape(persistent=True) as tape:
    y = x @ w + b # tape can be used multiple times to calc grad
    loss = tf.reduce_mean(y**2) # get the mean value
[d1_dw, d1_db] = tape.gradient(loss, [w, b]) # gradients to both vars
```

Machine Learning Datasets

- Machine Learning datasets in TF

- TensorFlow provides several built-in datasets for machine learning tasks from various areas and different topics, e.g. 3d - 2d image detection, classification and generation, categorical solutions, abstractive text summarization, anomaly detection and language modeling

detailed infos in <https://www.tensorflow.org/datasets/catalog/overview>

- cached locally in `~/tensorflow_datasets` in one directory per dataset

```
import tensorflow as tf
(x_train, y_train), (x_test, y_test) = tf.keras.datasets.mnist.load_data()
```

- **Example:** import horses_or_humans dataset, select horses and preprocessing

```
import tensorflow_datasets as tfds
dataset, info = tfds.load('horses_or_humans', with_info=True)
horse_ds = dataset['train'].filter(lambda x: x['label'] == 0) # horses
horse_examples = horse_ds.take(5) # 5 horses from the dataset

train_dataset, valid_dataset = tfds.load('horses_or_humans',
                                          split=['train', 'test'], as_supervised=True)

def preprocess(image, label):
    image = tf.cast(image, tf.float32)
    image = tf.image.resize(image, (IMG_SIZE, IMG_SIZE)) #resize to a fixed size
    image = image / 255.0 #rescale the pixel values to be between 0 and 1
    label = tf.one_hot(label, NUM_CLASSES) #assign labels
    return image, label
```

Machine Learning Datasets

- Input data usually has to be pre processed

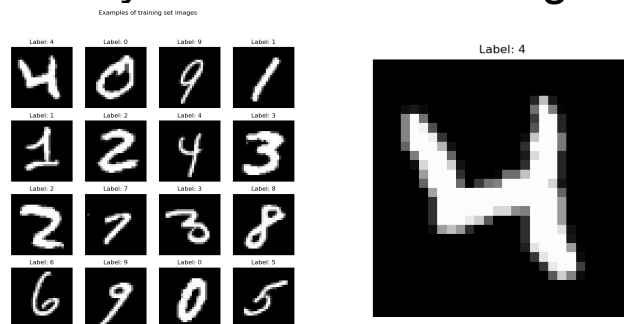
- Needed for numerical stability
- Want to have the mean around zero and the variance order of one
 - gaussian like distributed data → subtract mean and divide by standard deviation
 - data peaking at zero → apply logarithm and then subtract mean and divide by standard deviation

- Example datasets in the MNIST database (Modified National Institute of Standards and Technology database) used for machine learning and practicing techniques

The data can be loaded directly in Tensorflow using e.g. `tf.keras.datasets.mnist.load_data()`

- MNIST handwriting
60000 images

[03_ml_basics_display_Ha](#)



28 x 28 array with greyscale values, 0 – 255
→ preprocessing: divide by 255

use pixel array and label

- Fashion-MNIST dataset
60000 images

[03_ml_basics_display_Clothing.ip](#)



28 x 28 array with greyscale values, 0 – 255 and labels

- Input data needs often reshaping

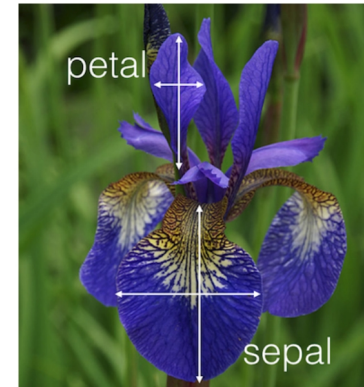
Machine Learning Datasets

- Example datasets for classification applications

- Iris dataset from Fisher (1936) with 3 classes, Iris Setosa, Iris Versicolour and Iris Virginica with 50 instances each. There are 4 attributes:

- sepal length [cm]
- sepal width [cm]
- petal length [cm]
- petal width [cm]

The application is for classification

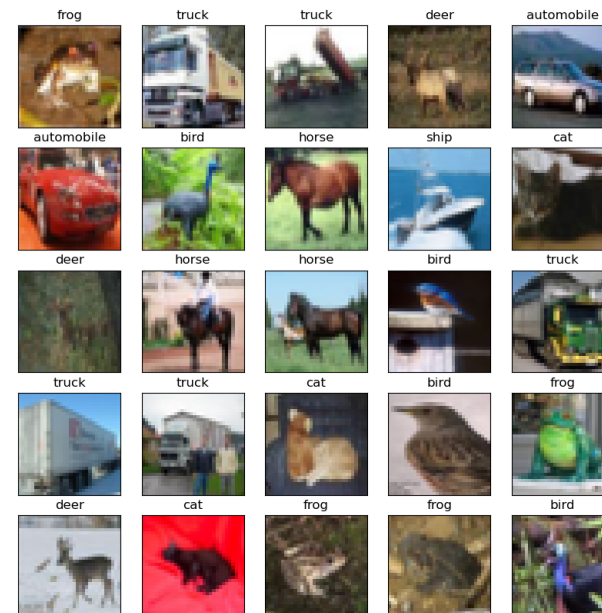


<https://archive.ics.uci.edu/ml/datasets/iris>

- Cifar-10 and Cifar-100 dataset of 32 x 32 color pictures in 10 classes (100 classes) with 60000 images per class and 10000 test images for Cifar-10.

The 100 classes in the CIFAR-100 are grouped also into 20 superclasses. There are 500 training images and 100 testing images per class

<https://www.cs.toronto.edu/~kriz/cifar.html>



Exercise 1

- Read/load the cifar10 dataset using `tf.keras.datasets`

There are the following classes

```
class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog',  
               'frog', 'horse', 'ship', 'truck']
```

Display the first 25 images

- They are color images, try to convert them to grey scale images by reducing the 3 colors (r,g,b) to one grey scale with the formula

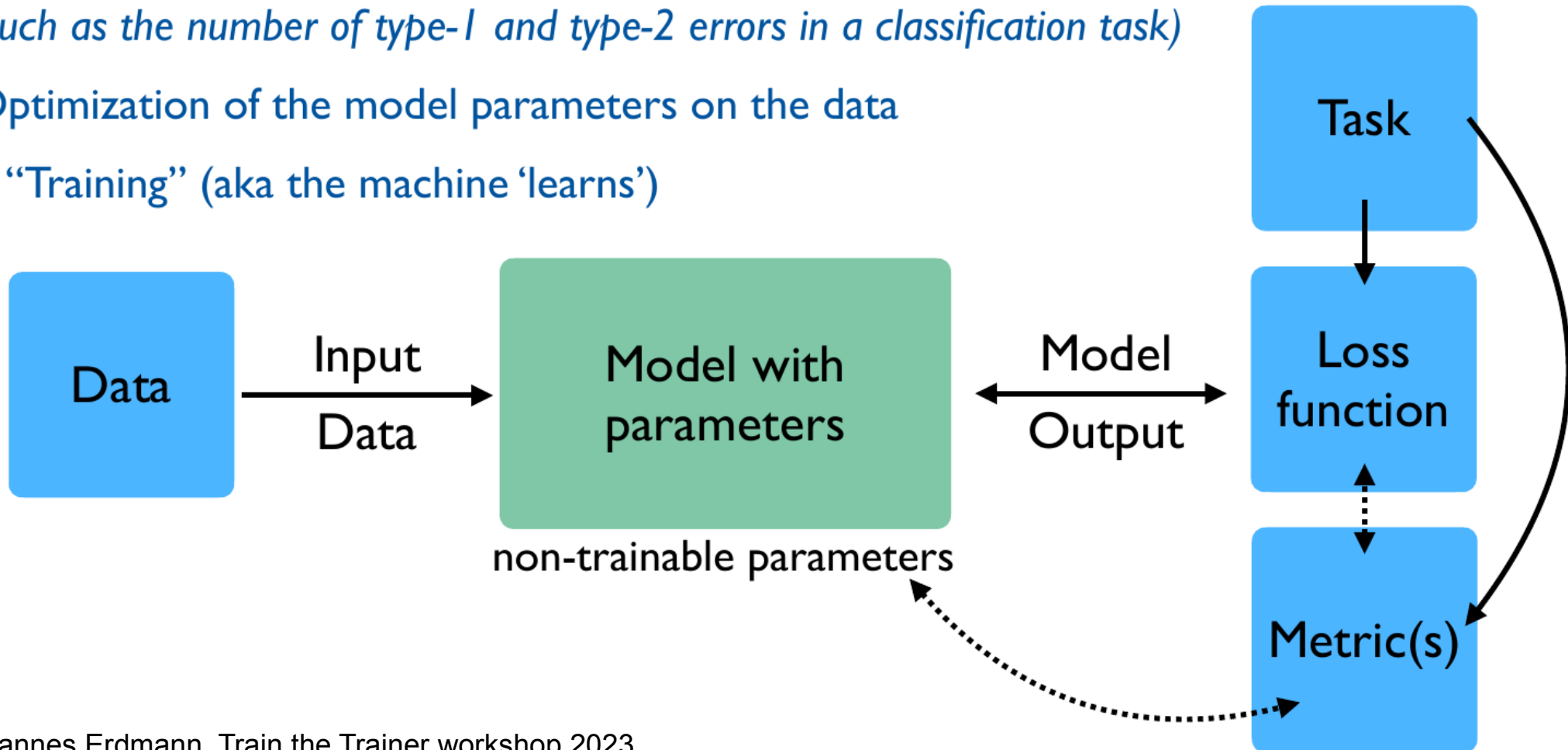
$$\text{gray} = 0.2989 * r + 0.5870 * g + 0.1140 * b$$

Normalize the data to the interval [0,1]

Display the first 25 greyscale images

Introduction Machine Learning

- Data + Task + “a mathematical formulation of what means ‘better’ ” *
(such as the number of type-1 and type-2 errors in a classification task)
- Optimization of the model parameters on the data
= “Training” (aka the machine ‘learns’)



DNN Training

- Metric

Machine learning datasets contain truth information to allow an independent quality measure of the neural network. Metrics are typically defined in terms of the neural network model's predictions and the true labels. They quantify how well the model is doing at its task.

- accuracy: number of correctly classified data points over the total number.
- precision: ratio of true positive (TP) predictions to the total number of positive (TP + false positive) predictions.
- recall: ratio of true positives to the sum of true positives and false negatives
- F1 score: $F1 \text{ score} = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$

- Splitting the dataset in 2 or 3 parts

- Training dataset – Do all the parameter determination (training) of the neural network.
- Validation dataset – Used to evaluate the model during training and to tune its hyperparameters (e.g. handled by tensorflow during minimization)
- Test dataset – Use it for evaluating the optimized parameters and the final performance of the model after it has been fully trained.

Splitting: Training:Validation:Test → 60%:20%:20%

DNN Minimization

- Find the network weights such that the loss function is minimal

$$W_{min} = \underset{w}{\operatorname{argmin}} J(W) = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

Repeat until convergence

- Initialize weights randomly
- Loop until convergence:

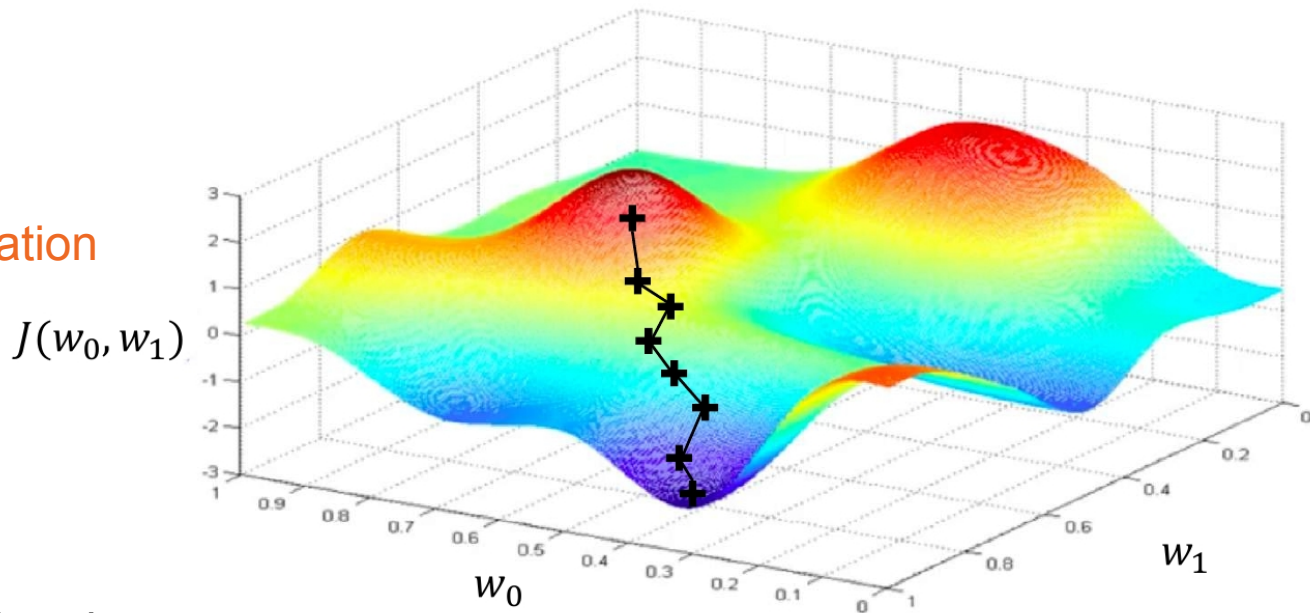
compute $\frac{\partial J(W)}{\partial W}$

backpropagation

update weights

$$W \leftarrow W - \eta \cdot \frac{\partial J(W)}{\partial W}$$

- return weights
- derivative calculation with chain rule



Introduction to Deep Neural Networks

- Find the network weights such that the loss function is minimal

$$W_{min} = \underset{w}{\operatorname{argmin}} J(W) = \operatorname{argmin} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

- Initialize weights randomly
- Loop until convergence:

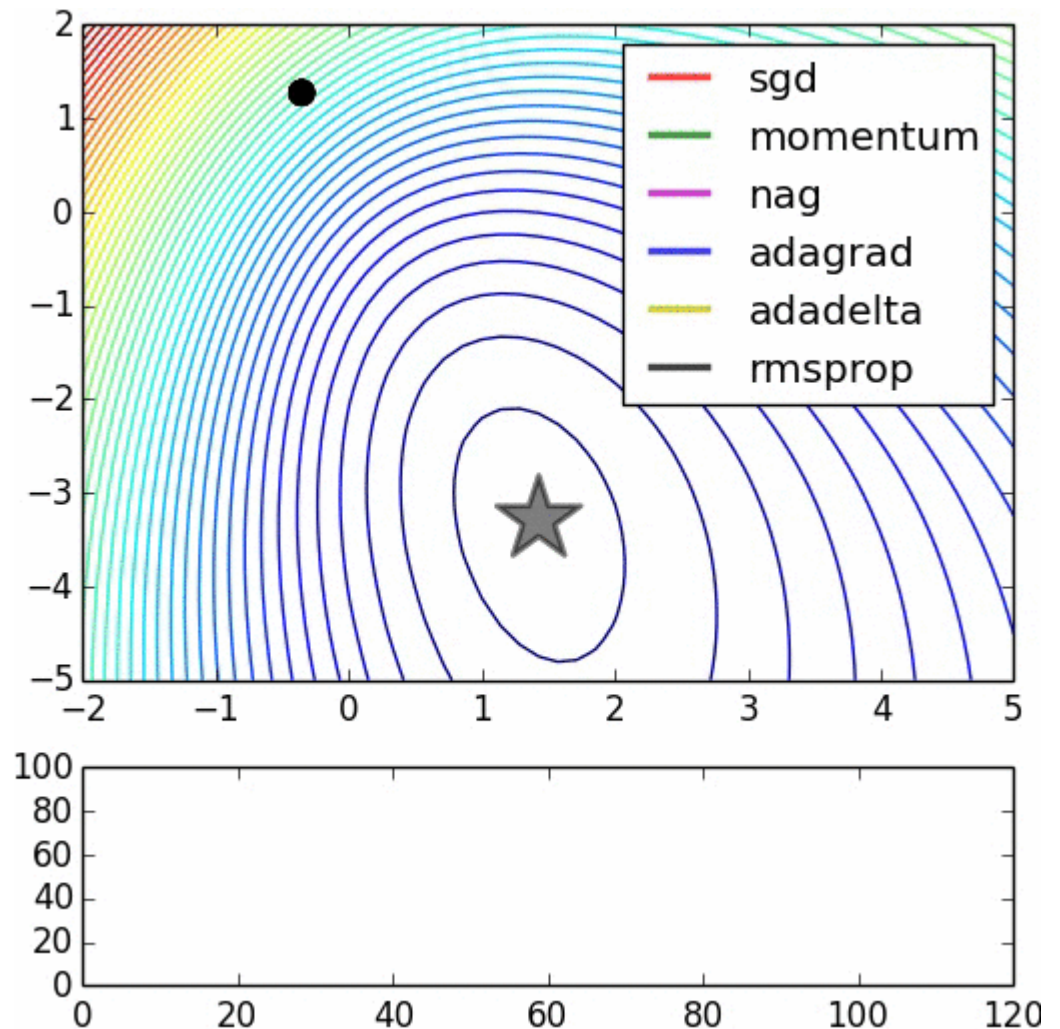
compute $\frac{\partial J(W)}{\partial W}$

backpropagation

update weights

$$W \leftarrow W - \eta \cdot \frac{\partial J(W)}{\partial W}$$

- return weights
- derivative calculation with chain rule



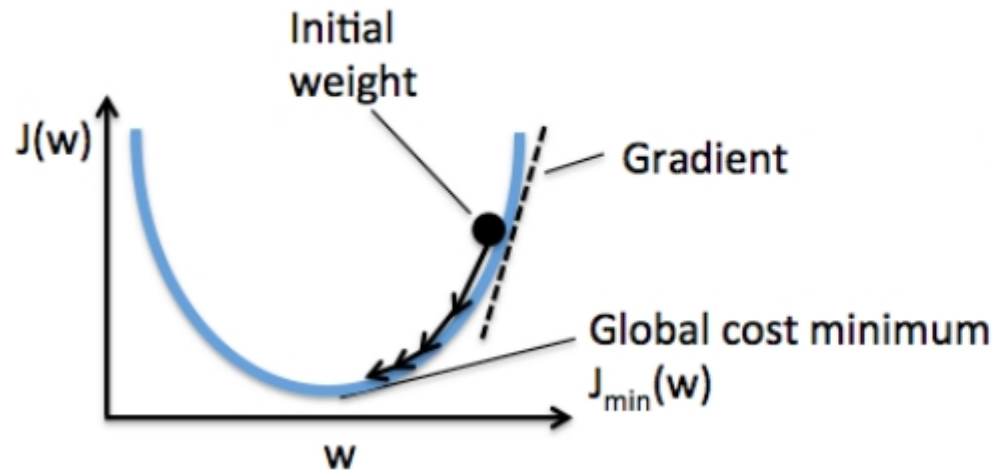
DNN Minimization

- Training strategies

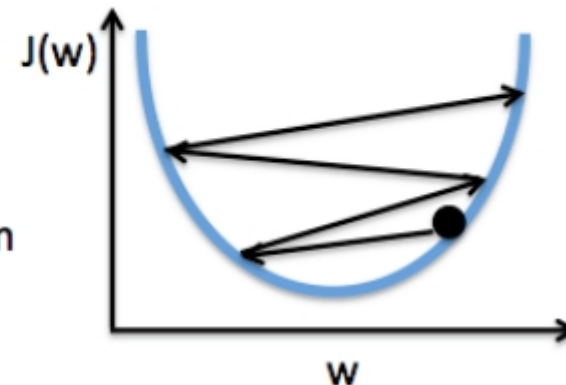
- Minimization: difficult to converge to a global minimum with large number of parameters.
- Learning rate

$$W \rightarrow W - \alpha \cdot \frac{\partial J(W)}{\partial W}$$

α is called learning rate



choosing the right learning rate is important for convergence



learning rate is too large

[https://raw.githubusercontent.com/rasbt/python-machine-learning-book/master/code/ch02/images/02_12.png, MIT License]

- Adaptive moment estimation (Adam) uses adaptive learning rates which reduces α during minimization for each parameter separately

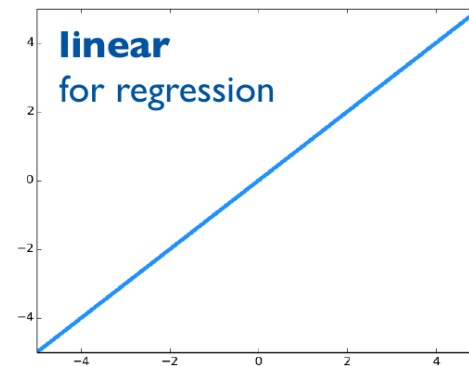
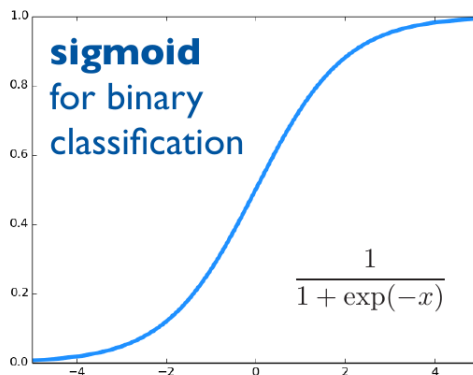
$$W \rightarrow W - (1 - \beta) \cdot \alpha \cdot \frac{\partial J(W)}{\partial W}$$

β change in W from previous step

DNN Minimization

- Output to the last node(s)

- The output function(s) of the neural network (nodes in the final layer) are determined by activation functions. Their choice depends on the problem to be solved.
- For binary classification problems the sigmoid function is used and interpreted as probability.
- For regression problems a linear activation function or no activation function is used.

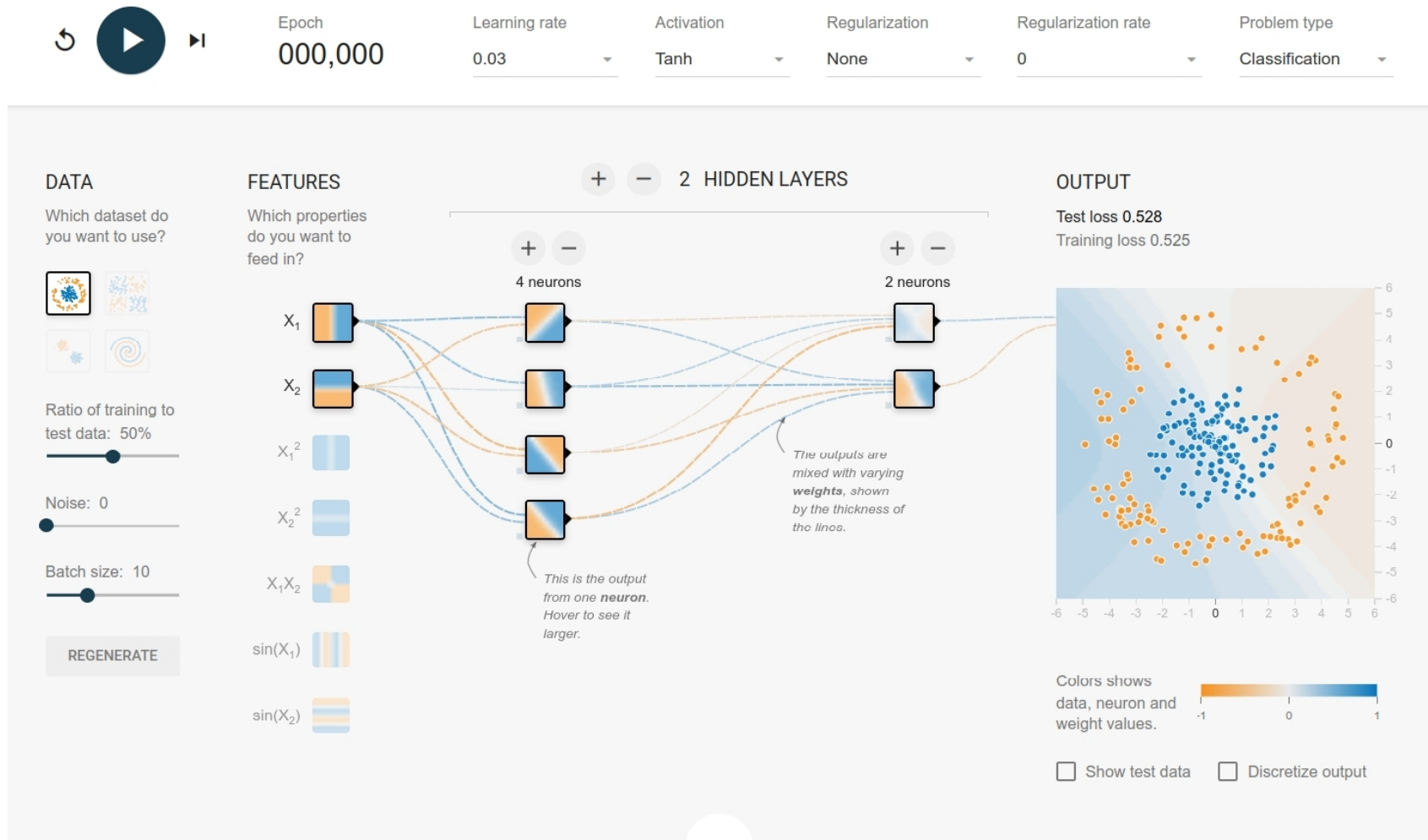


- For multi-class classification problems the softmax function is mostly in use and gives a probability distribution over the classes. Softmax assigns decimal probabilities to each class i of a multi class problem with its N multiclass output values:

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$$

Exercise 2

The following page illustrates a classification task which can be configured online via the web page <https://playground.tensorflow.org/>
- Try different activation functions and change the learning rate



DNN Minimization

- Minimizer algorithms: difficult to converge to global minimum with large number of parameters
 - various minimizers are available in the `tf.keras.optimizers` module of TensorFlow
 - Stochastic Gradient Descent (SGD)
gradient descent algorithm that updates the model parameters (weights and biases) in small batches based on the gradients computed on each batch in the direction of the negative gradient scaled by α
 - Momentum optimizer
updates the parameters by computing the gradient of the loss function with respect to the parameters and adjusting the parameters in the direction of the negative gradient taking into account previous updates to get the direction and magnitude
 - AdaGrad optimizer
adapt the learning rate for each parameter based on the historical gradient by keeping a running sum of the squared gradients for each parameter (accumulator) and divide the learning rate by the square root of the accumulator for each parameter, which adapts it.
 - RMSProp optimizer
In RMSprop (Root Mean Square Propagation) the learning rate is adaptively scaled for each weight parameter based on the root mean square (RMS) of the gradients. This reduces the learning rate for parameters with large and rapidly changing gradients
 - Adaptive Moment Estimation (Adam)
 - AdaDelta optimizer
 - Ftrl optimizer
 - Choice of optimizer often depends on the specifics of the problem
 - experiment with different optimizers

Examples - Deep Neural Networks

- We need start values for the network

- Initialize randomly, a range of values is needed, suitable values depend on the details of the network, like layer size and activation functions

- In general:

$$\text{var(input)} \approx \text{var(output)} \quad \text{with} \quad \text{var} \approx 2 / (N_{\text{input nodes}} + N_{\text{output nodes}})$$

draw from gaussian or uniform distributions within a range $\pm\sqrt{3\text{var}}$

- Usually input range differs largely

- transform to mean 0 and variance 1

- perform de-correlation of input data

- Simple example using TensorFlow (TF)

tf_intro.py

- Generate toy sample with 2 normalized gaussian distributions with mean (-1,-1) and (1,1)

- Each sample gets a label and then they are combined to a training set

- The data is plotted and colored according to their labels

- In TensorFlow 's a model is set up specifying the optimizer, loss function and metric

- Use AdamOptimizer to find the minimum and accuracy as metric

- The data and labels are given to the fit method, also the training iterations and the batch size. The model doesn't process the entire dataset at once, but rather in batches.

- Loss and accuracy are plotted as number of iterations

- Display classification results for sample points together with labeled data points

Examples - Deep Neural Networks

- Using the handwriting MNIST example with TensorFlow (TF)
 - Read in the dataset with a pre defined function in TF. There are 10 categories.
 - Each sample has a label and the 28 x 28 pixel grey scale data array.
 - The data is normalized and reshaped to a 1D array.
 - In TensorFlow 's a model is set up specifying a first dense hidden layer with 512 nodes and a ReLU activation function. Dense means all nodes are fully connected. A second layer also with 512 nodes and ReLU activation. The 3rd layer is the output layer with 10 nodes and Softmax is used as activation.
 - Use AdamOptimizer to find the minimum and accuracy as metric
 - The data and labels are given to the fit method, also the training iterations and the batch size. The model doesn't process the entire dataset at once, but rather in batches.
 - Loss and Accuracy are plotted as number of iterations
There is a difference between the training and the test data → Overtraining, the net learns features of the data which are not existing. To prevent this reduce the number of nodes or refine the model
 - Test accuracy = 0.98
 - The data is plotted and true and predicted label are written [MLP_MNIST_digits_tf.py](#)

Convolutional Neural Networks

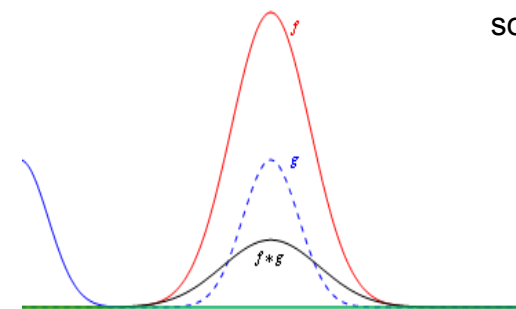
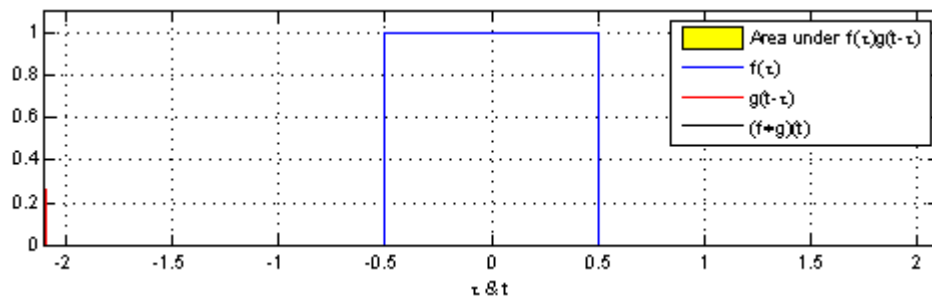
- A Convolutional Neural Network (CNN) learns features directly from data without specifying the features explicitly. CNNs are useful for pattern recognition in images and computer vision, but also for classification tasks in audio data and signal processing.
- Learning the features happens via convolution filter, which are activated by scanning over the pixel of an image. Edges or structure and color changes are transported to another filter layer. A Rectified Linear Unit (ReLU) is applied as activation function. Via a pooling technique the number of parameters is reduced from layer to layer.

- Convolution
$$(f * g)(x) = \int_{\mathbb{R}^n} f(\tau)g(x - \tau)d\tau$$

$f(x)$: input distribution

$g(x - \tau)$: kernel or filter function

The function value of the weight function f at the position τ tells us how strong the contribution of $g(x - \tau)$ in f at the position x is.



source Wikipedia

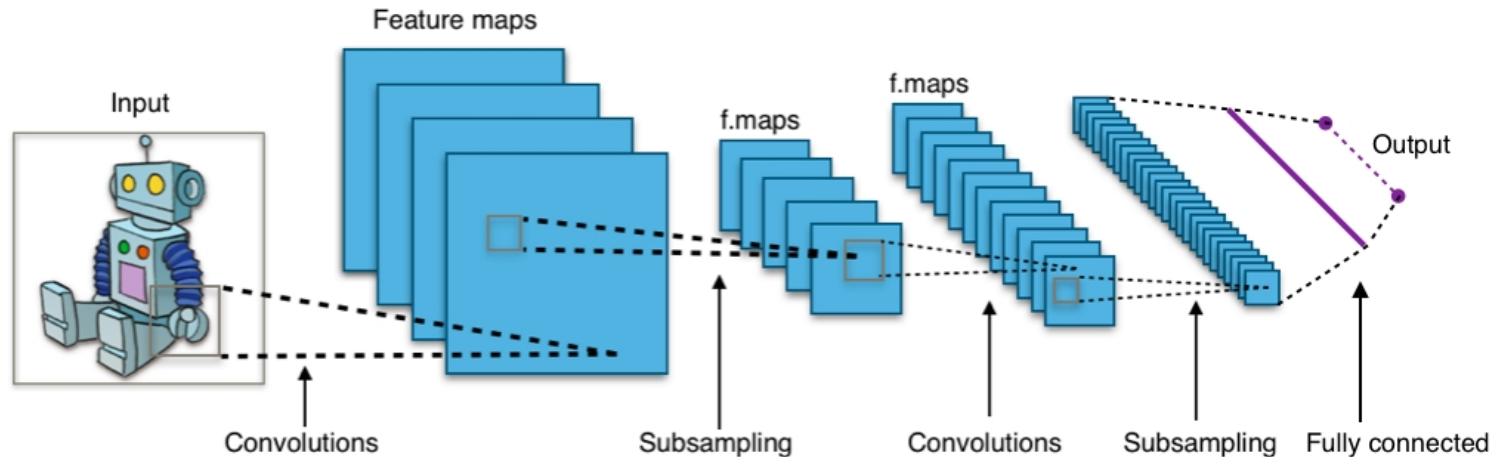
- Discrete convolution
$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$

Play with images and various kernel for image processing: <https://setosa.io/ev/image-kernels/>

Convolutional Neural Networks

More information can be found in the Stanford course <https://cs231n.github.io/convolutional-networks/>

- Structure of a typical CNN used in image classification



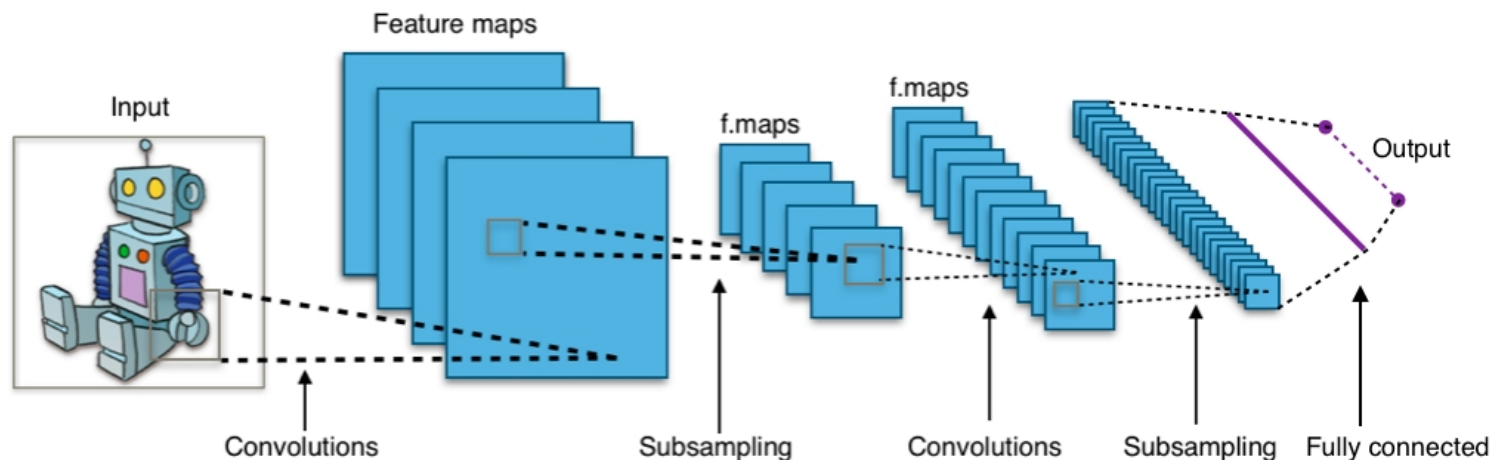
- Main idea is to extract particular localized features of data, eg. an image, using a convolution as filter mechanism. The input is a data structure [xPix,yPix,3] of raw pixel values with three color channels R,G,B. Each color channel is treated independently.
- There are 3 main building blocks:

1) Convolutional layer

- Obtain a weight matrix by computing the dot product of the weight matrix and a small sub region of the input data and scanning over the whole image. This results in a weight matrix which transports certain features of the image to further layers.
- To the weights activation functions like ReLU are applied. The convolution operation (weight matrix · sub input structure) behaves like a filter.
- The weight matrix is determined by a loss function.
- Multiple convolutional layers extract with increasing depth more and more complex features

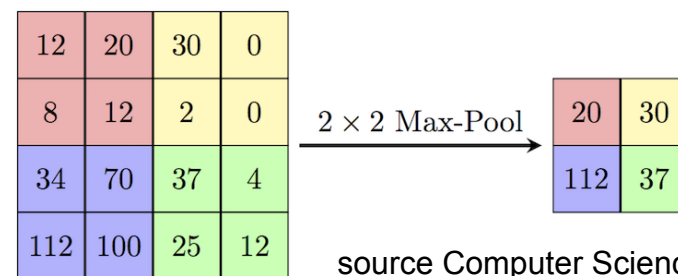
Convolutional Neural Networks

- Structure of a typical CNN used in image classification



II) Pooling layer

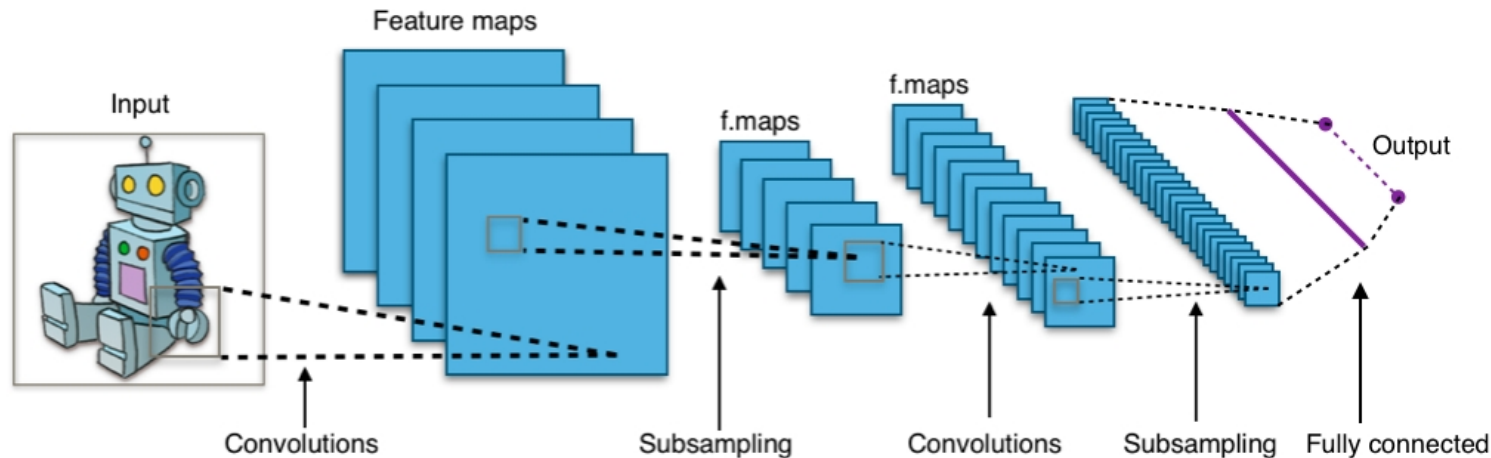
- Several neighbouring pixel are pooled together by averaging or by taking their maximum.
- Max Pooling algorithm: Take a 4x4 input (W), step (S) with a 2x2 filter (F) over the input, take the maximum entry
- Zero Padding: to the edges of the picture pixel are added (P) and filled with zeros



- Output size OP :
$$OP_{x,y} = \frac{W - F_{x,y} + 2P}{S_{x,y}} + 1$$

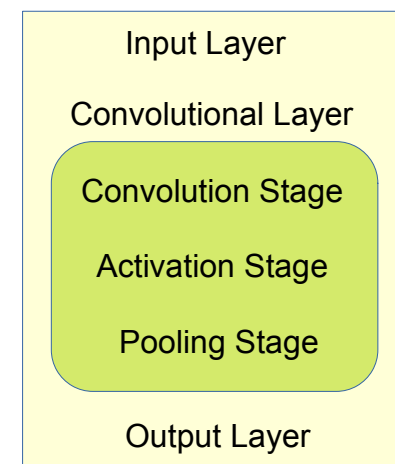
Convolutional Neural Networks

- Structure of a typical CNN used in image classification



III) Output layer

- The final stage is a fully connected layer which works as MLP to generate an output
- The very last layer is equal to the number of output classes. In order to obtain probabilities for each class SoftMax is applied
- A loss function determines all the weights



- As CNN example the MNIST handwriting dst is used which over comes overfitting

Examples - Deep Neural Networks

- Two extreme cases of training results

- If the model does not reflect the data content or the training is insufficient
→ bad network performance
- If the model allows for too much complexity it learns features of the training data sample
→ network can be applied to other samples (overtraining effect)
test overtraining in the example clothing dataset by including the method Dropout in TF or changing the number of nodes in the hidden layer

- Other classification examples with Tensorflow and Keras

- uses the Fashion MNIST dataset of Zalando, which contains 60,000 grey scale images in 10 categories each showing low resolution clothing pictures.
- sequential model with two dense layers → significant overtraining with test accuracy = 87.7%
`clothingSequential.py`
- sequential model with two dense layers adding Dropout and softmax activation
→ no overtraining with test accuracy = 88.2 %
`clothingSequentialDropout.py`
- CNN model with MaxPooling2D and 3x3 Filters and ReLu activation. The final output is 10 nodes with Softmax activation → overtraining with test accuracy = 90.6 %
`clothingCNN.py`
- CNN model where the input data is shaped to a 4D tensor of shape (samples, height, width, channels). MaxPooling2D and 3x3 Filters and ReLu activation and the Dropout feature is used → no overtraining with test accuracy = 91.4 %
`clothingCNN_4D.py`