

# 1 Deep Learning Theory

The following section covers the most important background theory for the experiments in Section ???. This includes an introduction to various types of neural networks, as well as an introduction to the optimization of such networks.

## 1.1 Feedforward Neural Networks

**Feedforward neural networks** are the most basic type of neural networks. The aim of a feedforward neural network is to approximate some function  $f^*$ , by defining a mapping  $\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta})$  and learning the parameters  $\boldsymbol{\theta}$ , that results in the best approximation of  $f^*$ . These models are called **feedforward** because there are no **feedback** connections in which the outputs of the model are fed back into itself. Instead, information flows through the function being evaluated from  $\mathbf{x}$ , through the intermediate computations used to define  $f$ , and finally to the output  $\mathbf{y}$ . Feedforward neural networks generally consists of multiple **layers**, arranged in a chain structure, with each layer being a function of the layer that preceded it [4].

### 1.1.1 Fully-connected Layers

The most simple type of layer found in a feedforward neural network is the **fully-connected layer**. The fully-connected layer usually consists of some learnable parameter matrix and learnable parameter vector, as well as a **activation function**, which is a non-linear function, that makes it possible for the network to model non-linearity. Three of the most common activation functions are the **Rectified Linear Unit (ReLU)**, **Sigmoid** and **Hyperbolic Tangent (Tanh)** activation functions, defined as

$$\text{ReLU}(x) = \max\{0, x\} \quad (1)$$

$$\text{Sigmoid}(x) = \sigma(x) = \frac{1}{1 + \exp(x)} \quad (2)$$

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \quad (3)$$

Given the learnable parameter matrix  $\mathbf{W}$ , learnable parameter vector  $\mathbf{b}$  and activation function  $g$ , the  $i$ 'th fully-connected layer is defined as

$$\mathbf{h}^{(i)} = \begin{cases} g^{(i)}(\mathbf{W}^{(i)\top} \mathbf{h}^{(i-1)} + \mathbf{b}^{(i)}) & \text{if } i > 1 \\ g^{(1)}(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) & \text{if } i = 1 \end{cases} \quad (4)$$

Thus, for a neural network with  $n$  layers, we have the mapping [4]

$$\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{h}^{(n)}. \quad (5)$$

### 1.1.2 Convolutional Layer

A **convolutional layer** is a specialized kind of feedforward layer, usually used in analysis of time-series or image data. If a network has at least one convolutional layer, it is called a **Convolutional neural network (CNN)** [4].

The convolutional layer consists of a set of **kernels**, each to be applied to the entire input vector, where each kernel is a learnable parameter matrix  $k \times k$  [7]. Each kernel is applied on the input to produce a **feature map**. The kernels are applied to the input by "sliding" over the input (where the step size is called **stride**). Each  $k \times k$  grid of the input is then used to

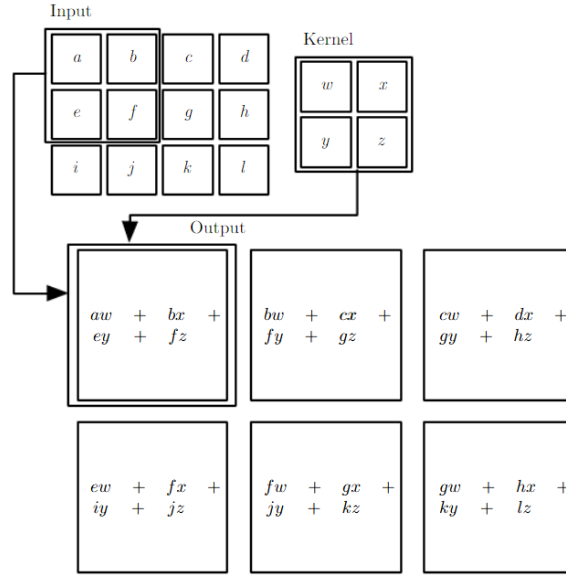


Figure 1: An example of applying a 2d kernel on an input [4].

compute the dot-product between the grid and each kernel, which is then placed in the corresponding feature map of each kernel, as visualized in Figure 1. [1]. To control the dimensions of the output, one might **pad** the sides with a constant value. Commonly, zero is used as the padding-value.

As seen in Figure 1, each kernel produces a linear combination of all pixel values in a neighbourhood defined by the size of the kernel. Thus, unlike a fully-connected layer, a convolutional layer captures the high correlation between a pixel and its neighbours. Further, by limiting the size of the kernel, the network will use much fewer parameters, than if it was replaced by a fully-connected layer [4].

## 1.2 Recurrent Neural Networks

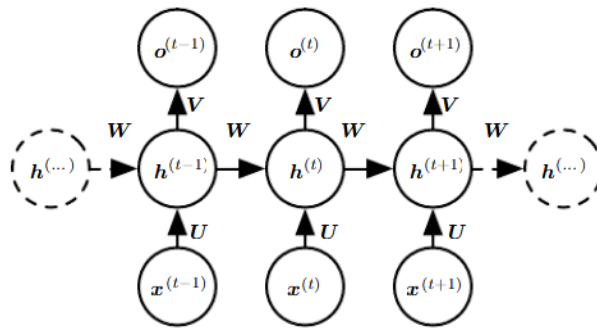


Figure 2: An illustration of an RNN [4].

**Recurrent neural networks (RNNs)** are a family of neural networks for processing sequential data. Figure 2 illustrates the general setup of such a network, which maps an input sequence of  $x$  values to a corresponding sequence of output  $o$  values. Generally, an RNN consists of three parts: (1) the input ( $x^{(i)}$ ), (2) the hidden state ( $h^{(i)}$ ), and (3) the output ( $o^{(i)}$ ). During inference, the model maps each input value to an output value in a sequential matter, where it first maps the first input value, then the second, then the third, and so forth. The network maps each input value to an output value by making use of the hidden state from the preceding step,

where the first hidden state has to be initialized. This kind of RNN is **unidirectional**, as it only considers the data in one sequence-direction. To utilize both directions of the sequence one can use two RNNs, where each RNN processes the data in opposite direction. This type of RNN is called being a **bidirectional** RNN (**bi-RNN**) [4].

### 1.2.1 Convolutional Long Short-Term Memory

One common recurrent neural network unit is the **convolutional long short-term memory (ConvLSTM)**, which is an adaptation of the standard **long short-term memory (LSTM)** for sequences of images. Both LSTMs work by potentially stacking multiple **cells** together, such that each LSTM cell is taking in outputs of the preceding LSTM cell as its input.

The idea of both LSTM cells is to create paths through time that have derivatives that neither vanish nor explode. This is done by introducing a memory cell  $C_t$ , which accumulates state information over a long duration. This cell is accessed, written and cleared by several controlling gates. The model learns during training, when to access, write and clear this memory cell. By using the memory cell and gates to control the flow of information, the gradient will be trapped in the cell and thus be prevented from vanishing too quickly [4, 8].

The ConvLSTM cell consists of three gates. The first gate is the input gate  $I_t$ , which controls whether or not to accumulate the information of a new input to the cell. This gate is defined as

$$I_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \circ C_{t-1} + B_i) \quad (6)$$

where  $X_t$  is the current input image,  $H_{t-1}$  is the hidden state of the previous time step,  $C_{t-1}$  is the cell output of the previous time step, and  $W_{xi}$ ,  $W_{hi}$ ,  $W_{ci}$  and  $B_i$  are learnable parameters

The second gate is the **forget gate** unit  $F_t$  (at time step  $t$ ), which controls whether or not to "forget" the status of the cell output of the previous time step. The forget gate is defined as

$$F_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \circ C_{t-1} + B_f) \quad (7)$$

where  $W_{xf}$ ,  $W_{hf}$ ,  $W_{cf}$  and  $B_f$  are learnable parameters.

The last gate is the **output gate** unit  $O_t$ , which controls whether or not the latest cell output will be propagated to the final state  $H_t$ . This cell is defined as

$$O_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \circ C_t + B_o) \quad (8)$$

where  $W_{xo}$ ,  $W_{ho}$ ,  $W_{co}$  and  $B_o$  are learnable parameters.

By combining the three gates, we get the following definition of the ConvLSTM Cell. Let  $X_1, X_2, \dots, X_t$  be a sequence of input images. Then, at each time step  $t$ , we compute

$$I_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \circ C_{t-1} + B_i) \quad (9)$$

$$F_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \circ C_{t-1} + B_f) \quad (10)$$

$$C_t = F_t \circ C_{t-1} + I_t \circ \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + B_c) \quad (11)$$

$$O_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \circ C_t + B_o) \quad (12)$$

$$H_t = O_t \circ \tanh(C_t) \quad (13)$$

For  $t = 1$ , both  $H_{t-1}$  and  $C_{t-1}$  have to be initialized [8].

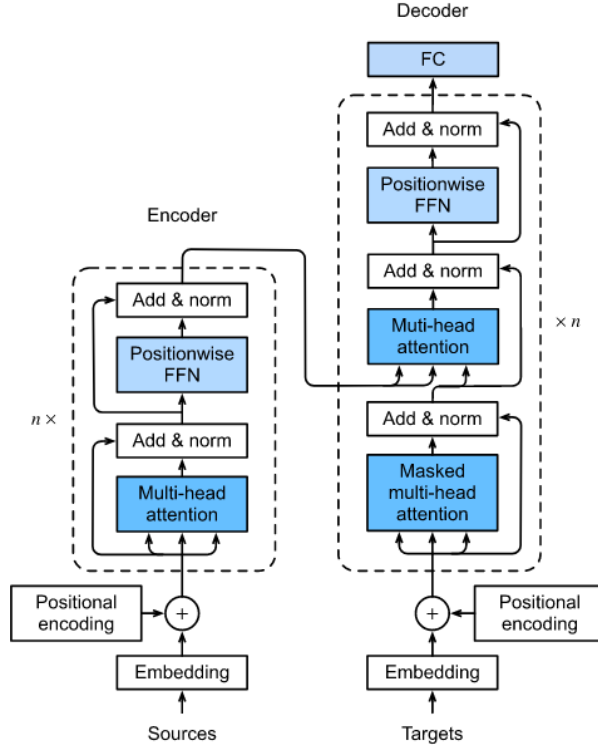


Figure 3: An illustration of the Transformer architecture [10].

### 1.3 Transformer

The sequential nature of RNNs precludes parallelization with training examples, heavily slowing down the training of these models. Another type of model for sequential data is the **Transformer**, which eschews recurrence and instead relies on an **attention** mechanism to draw dependencies between input and output. The Transformer allows for more parallelization and can reach state of the art results [9].

#### 1.3.1 Attention

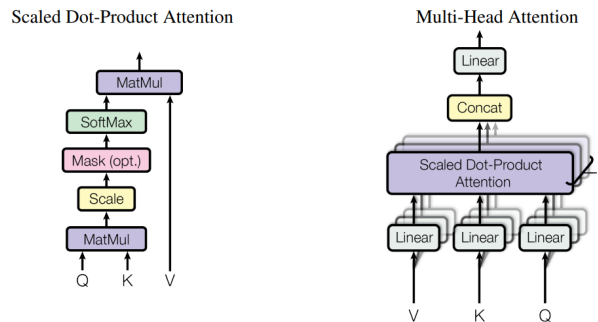


Figure 4: Illustration of (left) the Scaled Dot-Product Attention, and (right) the Multi-Head Attention, which consists of several attention layers running in parallel [9].

An attention function is a function that maps a given query and a given set of key-value pairs to an output. In this function the query, keys, values, and output are all vectors, and the output is a weighted sum of the values, where the weight of each value is computed by a compatibility function of the query with the corresponding key [9]. Figure 4 illustrates the two types of

attention functions used in the Transformer.

**Scaled Dot-Product Attention:** The scaled dot-product attention is the core attention function of the transform. The function takes queries and keys of dimensions  $d_k$  as input, and values of dimension  $d_v$ . In practice, the attention function is computed on a set of queries, keys and values simultaneously, by packing them into matrices  $Q$ ,  $K$ , and  $V$ , respectively. Thus, the scaled dot-product attention is computed as

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (14)$$

where the scaling of  $QK^T$  is done to counteract the cases with a large  $d_k$ , which would result in the softmax-function having an extremely small gradient [9].

**Multi-Head Attention:** The **multi-head attention** is the used attention function by the Transformer and is an extension to the scaled dot-product attention. This attention function linearly projects the queries, keys and values  $h$  times with different, learned linear projections to  $d_k$ ,  $d_k$ , and  $d_v$  dimensions, respectively. For each of the projected queries, keys and values, the scaled dot-product attention is applied, yielding  $d_v$ -dimensional output values. These output values are finally concatenated and further projected, resulting in the final values. By using multi-head attention, the model is allowed to jointly attend to information from different representation subspaces at different positions [9].

### 1.3.2 The Architecture

The Transformer follow an encoder-decoder structure, where the encoder maps an input sequence  $x$  to a sequence of continuous representations  $z$ . The decoder then uses  $z$  to generate an output sequence  $y$ , one element at a time. At each step the model consumes the previously generated output element as additional input when generating the next output [9]. Figure 3 illustrates the overall architecture of the Transformer.

**Encoder:** The encoder consists of  $N$  identical layers, where each layer consists of two sub-layers. The first sub-layer is a multi-head **self-attention** layer, and the second sub-layer is a position-wise fully-connected feedforward network. Around each sub-layer is a **residual connection**, where the the input of the sub-layer is added to the output of the sub-layer. This residual connection is then followed by a round of **layer normalization**. In this self-attention layer the keys, values and queries come from the output of the previous layer in the encoder. Each position in the encoder can attend to all positions in the previous layer of the encoder [9].

**Decoder:** The decoder also consists of  $N$  identical layers. In addition to the two sub-layers in each encoder layer, the decoder also consists of a third sub-layer, which performs multi-head attention over the output of the encoder stack (called **cross-attention** [6]). Also here, is residual connections used around each of the sub-layers, followed by layer normalization. To ensure, that the predictions for position  $i$  only depends on the known outputs at positions less than  $i$ , the self-attention sub-layer in the decoder stack is modified. This self-attention sub-layer, allows each position in the decoder to attend to all positions in the decoder up to and including that position. The self-attention sub-layer is modified by masking out all values in the input of the softmax which correspond to illegal connections [9].

**Feedforward Networks:** Each of the layers in the encoder and decoder contains a fully-connected feedforward layer, consisting of two linear transformations with the ReLU activation-function

applied in between, as described below [9]

$$FFN(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2. \quad (15)$$

**Positional Encoding:** As the model does not contain any recurrences nor any convolutions, it has to carry some other type of information to know about the relative or absolute position of the elements in the input sequence. This is done by adding **positional encodings** to the input of the encoder and decoder stacks. The positional encoding is a vector of size  $d_{model}$  and is defined as the following

$$PE_{pos,2i} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad (16)$$

$$PE_{pos,2i+1} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad (17)$$

$$(18)$$

where  $pos$  is the position and  $i$  is the dimension [9].

## 1.4 General Machine Learning Terminology

The following section covers the general terminology and some algorithms, that will be used through this paper.

**Loss Function:** Training a machine learning model is generally done by minimizing a pre-defined **loss function**, which is used for measuring the error of the model. One common loss function for regression is the **Mean Squared error (MSE)**. Let  $f^*(\mathbf{x})$  be a vector of ground truth observations and  $f(\mathbf{x}; \boldsymbol{\theta})$  be an estimation of  $f^*(\mathbf{x})$ . Then, MSE is defined as

$$MSE = \frac{1}{n} \sum_{i=1}^n \left( f^* \left( \mathbf{x}^{(i)} \right) - f \left( \mathbf{x}^{(i)}; \boldsymbol{\theta} \right) \right)^2. \quad (19)$$

Thus, MSE measures the average squared difference between the ground truth observations and the estimated observations [5].

**Optimizer:** The minimization of the loss function is usually done by an **optimizer**, which is a type of algorithm, that is typically some type of variation of **gradient descent**. Gradient descent is a process that iteratively updates the parameters of a function by

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L \left( f \left( \mathbf{x}^{(i)}; \boldsymbol{\theta} \right), f^* \left( \mathbf{x}^{(i)} \right) \right) \quad (20)$$

where  $\eta$  is called the **learning-rate**, which controls the stepsize of each update, and  $L$  is a loss function such as MSE.

One of the most used optimizers is **ADAM**, which is an optimizer that uses **momentum** and **adaptive learning-rate** to accelerate its learning. Momentum works by accumulating an exponentially decaying moving average of past gradients and continuing to move in their direction, which especially accelerates learning in the face of high curvature, small but consistent gradients, or noisy gradients. Adaptive learning rate is a mechanism for adapting individual learning rates for the model parameters. The idea is as follows; if the partial derivative of the

loss, with respect to a given model parameter, remains the same sign, then the learning rate should increase. If that partial derivate changes sign, then the learning rate should decrease. The algorithm uses a hyperparameter  $\rho$  for controlling the length scale of the moving average. The pseudocode of the algorithm has been illustrated in Algorithm 1. [4].

**Online, mini-batch and batch gradient methods:** Unlike gradient descent, ADAM bases its computations on a **mini-batch** of samples instead of basing it on a single sample. Thus, there are three ways of sampling data for the optimizer. The first way is done by sampling a single sample at each iteration. This class of methods is called **online gradient methods**. The advantage of these methods is, that each iteration is very quick, as we only have to compute the gradient of a single sample, however, the main disadvantage is, that the algorithm uses a lot of iterations. On the other hand we have **batch gradient methods** which uses the average gradient of all of the  $n$  samples of the dataset to compute  $\nabla L(f(\mathbf{x}, \mathbf{y}))$ . By doing so the algorithm only requires one iteration, however, this iteration is much slower than each iteration of the online gradient methods. The last class of methods is **mini-batch gradient methods** which uses the average gradient of  $|B|$  samples of the dataset for computing  $\nabla L(f(\mathbf{x}, \mathbf{y}))$ , where  $B$  is a predefined hyperparameter called the **mini-batch size**. This method lies in between the online gradient methods and batch gradient methods, as it uses fewer iterations than the online gradient methods, however, it uses more iterations than the batch gradient methods. Likewise, each iteration is faster to compute than it is for the batch gradient methods, but slower than it is for the online gradient methods. Further, online gradient methods can lead to a very noisy learning process due to the high variance of the computed gradients of the samples, which can make it more difficult for the optimizer to reach the minimum of the loss function. mini-batch gradient methods and batch gradient methods removes this noise by averaging the computed gradient at each iteration, essentially leading to a quicker learning process [4].

---

**Algorithm 1** Adam [4]

---

**Require:** Learning rate  $\eta$

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$

**Require:** Small constant  $\delta$  used for numerical stabilization

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  random observations from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

    Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

    Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

    Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

    Compute update:  $\Delta \theta = -\eta \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$

    Apply update:  $\theta = \theta + \Delta \theta$

---

**Layer Normalization:** One problem when optimizing a neural network is, that the gradients with respect to the weights of one layer depends on the outputs of the neurons of the previous layer. If these outputs change in a highly correlated way, the model could experience some undesired "covariate shift", which would make the training take longer. One way



to reduce these "covariate shifts" is to make use of **layer normalization**, where the layer normalization statistics over all the hidden units in the same layer  $l$  is computed as the following

$$\mu^{(l)} = \frac{1}{H} \sum_{i=1}^H \mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)}, \quad \sigma^{(l)} = \sqrt{\frac{1}{H} \sum_{i=1}^H \left( \mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)} - \mu^{(l)} \right)^2} \quad (21)$$

where  $H$  is the amount of hidden units of the  $l$ th layer,  $\mathbf{w}_i^{(l)}$  is the incoming learnable weight vector to the  $i$ th hidden unit of the  $l$ th layer, and  $\mathbf{h}^{(l-1)}$  is the output of the preceding layer. Thus, the summed inputs  $\mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)}$  is normalized and rescaled by

$$\bar{a}_i = \frac{\mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)} - \mu^{(l)}}{\sigma^{(l)}}. \quad (22)$$

Further, the input to the current layer  $l$  also learns an adaptive bias  $b$  and gain  $g$  for each neuron after the normalization. Thus, the  $i$ th input to layer  $l$  is defined as the following

$$h_i^{(l)} = f \left( \frac{g_i(a_i - \mu^{(l)}) + b_i}{\sigma^{(l)}} \right) \quad (23)$$

where  $f$  is a non-linear activation function [3].

**Evaluation:** We often want to know how well a machine learning model performs on an unseen sample. This is commonly done by splitting the dataset into three non-overlapping subsets to avoid a biased evaluation. The first dataset is the **training** dataset which is used for training the machine learning model. The second dataset is the **validation** dataset which is used during training for evaluating various settings of the the used (hyper)parameters. The last dataset is the **testing** dataset, which is then used for evaluating the final machine learning model.

The evaluation of a machine learning model is done by making use of an **evaluation metric**, which depends on the task of the machine learning model. For pose estimation, the **Percentage of Correct Keypoints (PCK)** is commonly used. This metric measures the ratio of predicted keypoints that are closer to their corresponding ground-truth keypoints than some threshold. Often **PCK@0.2**, **PCK@0.1**, **PCK@0.05** are used, which uses 20%, 10% and 5% of the torso diameter as its threshold, respectively [2].

**Dropout:** A machine learning model is said to be **overfitting** if the loss of the training dataset is decreasing, while the loss of the validation dataset is increasing. One can make use of a group of techniques called **regularization** to avoid the overfitting of the network. One common regularization technique is **dropout**, where the each parameter of the machine learning model is randomly disabled with probability  $p$  during training. By doing so the parameters are being trained on the dataset fewer times, leading to less of a probability of the model overfitting.

**Epoch:** Often, the optimizer needs multiple "rounds" of using the whole training data for reaching its minimum. Each of these rounds is called an **epoch**. One can pick the number of epochs to iterate through prior to training a model and just stop the training once the optimizer has used all of its epochs. Another method is to use **early-stopping**, where one keeps track of the loss of the model on the validation-set for each epoch. Once the validation-loss has not decreased for  $n$  consecutive epochs, the training is terminated.