

# DEEP LEARNING (DL)

1. Introduction	
<b>Fields of application</b>	<p><b>Computer Vision</b> (e.g. self-driving cars, object detection, medical diagnosis, image tagging/generation, caption generation, lip reading, lip synch from Audio)</p> <p><b>Image synthesis</b> (generation of scenes in games)</p> <p><b>Signal processing</b> (denoising signals, images)</p> <p><b>Natural Language Processing [NLP]</b> (e.g. translation, speech recognition and synthesis, Q&amp;A-bots)</p>
<b>Machine Learning (ML)</b>	Machine <b>learning</b> consists computer methods that analyse <b>observation data</b> to automatically detect patterns, and then use the uncovered patterns to perform <b>tasks</b> based on new unobserved data.
convergence	<p><b>Mathematics</b> (probability theory, statistics, regression, ...)</p> <p><b>Signal processing</b> (filtering, feature extraction, time series, fft, ...)</p> <p><b>Software Engineering</b> (very large DB, intensive CPU, distributed programming, ...)</p> <p><b>Application domain</b> (finance, medicine, energy, biometrics, ...)</p> <p><b>Algorithmic</b></p>
tasks	<p><b>Classifying:</b> Predict label of n mutually exclusive classes, e.g. classifying mails as spam or ham</p> <p><b>Predicting:</b> Predict numerical value (response), e.g. predicting house prices in given location and size</p> <p><b>Clustering:</b> Arrange set of entities into groups, e.g. group patients/clients for similar 'treatment'</p> <p><b>Robotic Tasks:</b> e.g. robot conducting certain tasks (walking, cutting the lawn, cleaning, washing dishes)</p>
process	
<b>Paradigms</b>	<p>With <b>supervised learning</b>, the goal is to extract some relevant features <b>x</b> from raw observation data <b>o</b> and to learn a <b>mapping</b> from inputs <b>x</b> to outputs <b>y</b> given a set of <b>example</b> data called the <b>training set</b>. <b>Labelled</b> data. Applications: Recognition, Planning, Diagnosis, Robot Control, Prediction</p> <p>With <b>unsupervised learning</b>, the goal is to discover interesting <b>structures</b> from inputs <b>x</b> given a set of <b>data</b> called the training set. <b>Unlabelled</b> data. Applications: Market Segmentation, Astronomical and Solar Data Analysis, Social Network Analysis.</p> <p>With <b>reinforcement learning</b>, we learn the behaviour in an environment that provides suitable <b>rewards</b>. Applications: Learn to play games, solve tasks</p>
<b>Deep Learning (DL)</b>	A sub-branch of machine learning. Neuronal architecture with many layers and neurons.
convergence	<p><b>Larger quantities of data</b> (text, audio, images, videos, ...) -&gt; scalability of learning on very large data sets</p> <p><b>New algorithms</b> (DBN, RBM, CNN, ...) -&gt; ability to learn feature extraction in unsupervised mode and classification in supervised mode</p> <p><b>Better computer performance</b> (GPU, distributed computing, ...) -&gt; train complex mapping functions</p>
<b>Machine Learning vs Deep Learning</b>	<p>Machine Learning typically <b>requires significant hand-engineering of features</b>.</p> <p>Deep Learning requires <b>less feature engineering</b> than standard machine learning.</p> <p>In DL, the machines find the features <b>automatically</b> as part of learning.</p> <p>DL has more flexible and powerful models but are more difficult and need more learn data.</p> <p>DL solves the problems in supervised learning.</p>
supervised learning problems	<ol style="list-style-type: none"> <li>1. We need <b>large quantities</b> of human validated examples! ... and this is <b>costly</b> to build.</li> <li>2. Because of the variabilities, we will need even <b>more data</b> and <b>complex mapping functions</b>.</li> <li>3. We spend a lot of time to hand-craft interesting compact features, so called <b>feature engineering</b>.</li> </ol>
Deep learning answers	<ol style="list-style-type: none"> <li>1. Let's use all the labelled data and unlabelled data</li> <li>2. Let's use deep neural networks</li> <li>3. Let's learn the feature extraction in unsupervised learning mode.</li> </ol>

# 1. Perceptron

## Biological Neural Systems

<b>Biological Neurons</b>	composed of: cell body, dendrites, axon composed in a <b>network</b> by synaptic terminals electrical impulses ' <b>signals</b> ' are sent via synapsis if a neuron receives enough signals, it fires its own signal ' <b>activation</b> '	
<b>Biological Neural Nets</b>	The connectivity in biological neural systems is huge. neurons in the brain: $\sim 10^{11}$ connections per neuron: $\sim 10^4$ $\rightarrow \sim 10^{15}$ synapses	
<b>Neuro-Science and Artificial Neural Networks</b>	Inspired by the biological brain. <b>Reverse-engineering</b> the computational principles behind the brain. Deep Learning goes beyond the neuro-scientific perspective and appeals to a more general principle of learning with multiple levels of composition. No claim to model the biological function directly. <b>Artificial Neural Nets</b> are composed of <b>Artificial Neurons</b> The network is trained to perform the task from <b>training data</b> by applying a <b>learning algorithm</b> that adjusts the networks parameters.	

## Artificial Neuron

<b>McCulloch-Pitts Neuron (1943)</b>	First artificial neuron as a model for the activation of a neuron. <table border="1" style="display: inline-table; margin-right: 20px;"> <tr> <td>Number of Neurons</td> <td><math>n</math></td> </tr> <tr> <td>Input signal</td> <td><math>x_k = 0, 1 \ 1 \leq k \leq n</math></td> </tr> <tr> <td>Weighted</td> <td><math>w_k = \begin{cases} +1 &amp; (\text{excitatory}) \\ -1 &amp; (\text{inhibitory}) \end{cases}</math></td> </tr> <tr> <td>Sum of all input signals</td> <td><math>S = \sum_{k=1}^n w_k x_k</math></td> </tr> <tr> <td>Output signal</td> <td><math>y = \begin{cases} +1 &amp; (S \geq \theta) \\ 0 &amp; (S &lt; \theta) \end{cases}</math></td> </tr> </table>		Number of Neurons	$n$	Input signal	$x_k = 0, 1 \ 1 \leq k \leq n$	Weighted	$w_k = \begin{cases} +1 & (\text{excitatory}) \\ -1 & (\text{inhibitory}) \end{cases}$	Sum of all input signals	$S = \sum_{k=1}^n w_k x_k$	Output signal	$y = \begin{cases} +1 & (S \geq \theta) \\ 0 & (S < \theta) \end{cases}$
Number of Neurons	$n$											
Input signal	$x_k = 0, 1 \ 1 \leq k \leq n$											
Weighted	$w_k = \begin{cases} +1 & (\text{excitatory}) \\ -1 & (\text{inhibitory}) \end{cases}$											
Sum of all input signals	$S = \sum_{k=1}^n w_k x_k$											
Output signal	$y = \begin{cases} +1 & (S \geq \theta) \\ 0 & (S < \theta) \end{cases}$											

Examples	AND: $y = H(x_1 + x_2 - 2)$ OR: $y = H(x_1 + x_2 - 1)$ XOR: $y = H(H(x_1 + x_2 - 0.5) + H(1.5 - x_1 - x_2) - 1.5)$	Heaviside-function $H(z) = \begin{cases} 1 & (z \geq 0) \\ 0 & (z < 0) \end{cases}$
----------	--	--

<b>Rosenblatt's Perceptron (1958)</b> = LTU (Linear Threshold Unit)	for linearly separable binary classification problems $1 = \text{yes}, 0 = \text{no}$ 0. initialize parameters (zero or random) 1. pick sample $(x^{(i)}, y^{(i)})$ 2. compute predicted value: $\hat{y}^{(i)} = H(\mathbf{w} * \mathbf{x}^{(i)} + b)$ 3. Parameter update rule: $\mathbf{w} \leftarrow \mathbf{w} - \alpha * (\hat{y}^{(i)} - y^{(i)}) * \mathbf{x}^{(i)}$ $b \leftarrow b - \alpha * (\hat{y}^{(i)} - y^{(i)})$ <p style="text-align: center; margin-left: 100px;"><small>0 or 1</small></p>	$y = H\left(\sum_{k=1}^n w_k x_k + b\right), x_k, w_k, b \in \mathbb{R}$
--	--	--

<b>Example with 2D Input Data</b>	<b>Decision Boundary</b> $H_{w,b} = w * x = 0$ $\rightarrow$ dot product is 0 when orthogonal $H_{w,b} = w_1 * x_1 + w_2 * x_2 + b = 0$ $x_2 = \frac{-b - w_1 * x_1}{w_2}$ $s = -\frac{w_1}{w_2}$	
-----------------------------------	--	--

<b>Perceptron Learning Algorithm</b>	The learning rule searches for a weights vector that defines a hyperplane that separates the points associated with the two classes. This is only possible for linearly separable input sets. <b>The solutions are not unique and not optimal.</b> The optimal solution is called <b>SVM (support vector machine)</b>	
--------------------------------------	---	--

<b>Perceptron Convergence</b>	Perceptron Learning Algorithm converges to a weights vector and bias that separates the two classes – provided that the two classes are <b>linearly separable</b> .	
-------------------------------	---	--

Artificial Neural Nets		
<b>XOR Problem</b>	Serious weakness of perceptron's: Incapable of solving some rather simple Problems These limitations can be overcome by stacking multiple perceptron's so called MLP.	
<b>Multi-Layer Perceptron (MLP)</b>	An MLP is composed of: <b>Input Layer:</b> Inputs passed through <b>Hidden Layers:</b> One or more layers of LTUs <b>Output layer:</b> Final layer of LTU  Input and hidden layers include a bias ( $x_0 = 1$ ) neuron and are <b>fully connected</b> to the next layer.	
<b>Usage</b>	typically, in areas where it is easy for humans and difficult for computers.	

## 2. Learning and Optimisation

<b>MNIST Dataset</b>	contains a lot of handwritten digits for testing and illustration -> identify which digit they are 2 versions: Original (70'000 28x28pixel images), Lightweight (1'800 8x8 pixel images, faster) Binary Classification problem (is it a e.g. 5 or not), use model-function $\hat{y} = h(x)$
----------------------	---

### Data Preparation

<b>Training and Testing</b>	Training set: Used for learning the task. Test Set: Used for testing how well the learned model performs. Split in test and train data -> randomly shuffle dataset before splitting split ratio depends on available data (large set: 99% to train, small set: 70% to train)	
<b>Data Normalisation (Feature Scaling)</b>	Bring your values to similar scales (range and importance). Apply to input and output data. 1. Scaling: improves convergence speed and accuracy of the learning algorithm 2. Centring: improves the robustness of the learning algorithm	
2 schemas	<b>Z-Normalisation</b> Shifting and rescaling the data so that a zero mean and a unit-variance is obtained $x_k^{(i)} = \frac{x_k^{(i)} - \mu_k}{\sigma_k}$ compute on training set: $\mu_k = \frac{i}{N} \sum_{k=1}^N x_k^{(i)} \rightarrow \text{mean deviation}$ $\sigma_k^2 = \frac{i}{N} \sum_{k=1}^N (x_k^{(i)} - \mu_k)^2 \rightarrow \text{standard deviation}$	<b>Min-Max Rescaling</b> $x_k^{(i)} = \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})} \rightarrow [0,1]$ <b>Min-Max Normalisation</b> $x_k^{(i)} = 2 * \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})} - 1 \rightarrow [-1,1]$ calculate min/max on training set

<b>Notations</b>	$m$	number of samples in the input dataset
	$n_x$	number of input features, dimension of the input feature vector
	$\mathbf{x}$	input feature vector of dimension $n_x$
	$x_k$	k-th component of the input feature vector ( $1 \leq k \leq n_x$ )
	$\mathbf{x}^{(i)}$	input feature vector of the i-th training sample ( $1 \leq i \leq m$ )
	$x_k^{(i)}$	k-th component of the input feature vector of the i-th training sample
	$y$	scalar output variable, also called target output or label
	$\mathbf{y}$	output vector (or target output vector) of dimension $n_y$ vector
	$y_k$	k-th component of the output vector (or target output vector) ( $1 \leq k \leq n_y$ )
	$\hat{y}$	predicted output, as computed by the mapping function
	$\hat{\mathbf{y}}$	predicted output vector, as computed by the mapping function
	$(\mathbf{x}, \mathbf{y})$	input sample (pair of input feature vector and corresponding label vector)
$(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$	i-th input sample of the input dataset ( $1 \leq i \leq m$ )	

### Model and Cost

<b>Generalised from Rosenblatt's Perceptron</b>		1) Smooth activation function: Sigmoid $\sigma(z) = \frac{1}{1 + e^{-z}}$ 2) Gradient Descent Optimization Algorithm - widely used in practice - gives the direction of the steepest ascent - works 'locally' and finds local wells (gradient=0) -> not designed to find global minima
<b>Learning Rate <math>\alpha</math></b> $\alpha > 0$	Determines the learning speed. Needs to be tuned to a given problem. If too large, it is not guaranteed to converge. If too small -> slow convergence of cost and error rate	

<b>Epochs</b>	Iterations run until the bottom of the valley reached. If too small -> not optimal values. If too big -> overfitting.	
<b>Probabilistic interpretation</b>	$p(y = 1 x, \theta) = h_{\theta}(x)$ $p(y = 0 x, \theta) = 1 - h_{\theta}(x)$	
<b>Cost functions</b>	<b>Mean Square Error Cost Function</b> $J_{MSE}(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$ Square to get positiv result. Training can get stuck.	<b>Cross-Entropy Cost Function</b> $J_{CE}(\theta) = -\frac{1}{m} \sum_{i=1}^m \log(p(y^{(i)} x^{(i)}, \theta))$ <b>Cross Entropy in generalised perceptron</b> $\nabla J_{CE}(\theta) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \begin{pmatrix} x^{(i)} \\ 1 \end{pmatrix}$ For classification tasks. Probabilistic consideration.
	<b>Cross-Entropy Loss Function for binary classification problem</b> $\mathcal{L}(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$ <b>Cross-Entropy Cost Function for binary classification problem</b> $J_{CE}(\theta) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$	
<b>Mathematical Formulation</b>	<b>Partial Derivative of Function</b> $J(\theta) = J(\theta_1, \dots, \theta_n)$ $\frac{\delta J}{\delta \theta_k} = \lim_{\epsilon \rightarrow 0} \frac{1}{\Delta \theta_k} (J(\theta_1, \dots, \theta_k + \Delta \theta_k, \dots, \theta_n) - J(\theta_1, \dots, \theta_k, \dots, \theta_n))$	
	<b>Gradient</b> $\nabla_{\theta} J = \frac{\delta J}{\delta \theta} = \begin{pmatrix} \frac{\delta J}{\delta \theta_1} \\ \dots \\ \frac{\delta J}{\delta \theta_n} \end{pmatrix}$	<b>General Gradient Descent Update-Rules</b> vector notation: $\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$ coordinates: $\theta_k \leftarrow \theta_k - \alpha \frac{\delta J(\theta)}{\delta \theta_k}$ <b>Update rules for generalised perceptron</b> $w \leftarrow w - \frac{\alpha}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x^{(i)}$ $b \leftarrow b - \frac{\alpha}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})$

### 3. Shallow Networks (MLP with a single hidden layer)

<b>SoftMax for Multi-Class</b>	$h_{\theta,l}(x) = \frac{\exp(z_l)}{\sum_{j=1} \exp(z_j)}$ , where $z_j = w_j * x + b$ It peaks at the largest $z_l$ and smoothly approximates $\max\{z_1, \dots, z_{k-1}\}$ if one element is much larger than all the others. Typically, as final layer. m inputs and m outputs.	
Training	Gradient $\frac{\delta}{\delta w_j} J_{CE}(\theta) = -\frac{1}{m} \sum_{i=1}^m (\delta_{j,y^{(i)}} - h_{\theta,j}(x^{(i)})) x^{(i)}$ $\frac{\delta}{\delta b_j} J_{CE}(\theta) = -\frac{1}{m} \sum_{i=1}^m (\delta_{j,y^{(i)}} - h_{\theta,j}(x^{(i)}))$	
	Update rules $w_j \leftarrow w_j - \alpha \frac{\delta}{\delta w_j} J_{CE}(\theta)$ $b_j \leftarrow b_j - \alpha \frac{\delta}{\delta b_j} J_{CE}(\theta)$	
<b>Adding Hidden Layers</b>	One additional layer with n neurons, sigmoid activation function. Can improve the performance. But in MNIST one layer already capture the correlation between pixels. <b>How do the formulas look like when applying gradient descent networks with hidden layers?</b>	
<b>Role of activation function</b>	Non linearities in the mapping between input and output of a neural network are crucial for gaining enough power for learning a task with enough accuracy. The choice of activation functions has an impact of robustness and performance.	
<b>Universal Approximation Theorem</b>	A feedforward network with a linear output layer and at least one hidden layer with a non-linear ("squashing") activation function (e.g. sigmoid) can approximate a large class of functions with arbitrary accuracy - provided that the network is given a sufficient number of hidden units and the parameters are suitably chosen.	
e.g. problem	combine 2 sigmoid to generate a step function -> with this we can approximate a large class of functions with a shallow network, any function can be represented, but with problems: - for improving accuracy, more and more neurons are needed (exponentially growing number) - with more dimensions (e.g. image or audio) more sampled data is needed -> <b>curse of dimensionality</b> - if we just use the available data and interpolate with step-functions between data points -> we overfit	

Overfitting			
<b>Overfitting</b>	<b>Underfitting</b> - High Bias: Strong bias in the way the data deviate from the linear model.	<b>Good Fit</b> - "Just Right": Model seems to capture just right the underlying structure in the data.	<b>Overfitting</b> - "High Variance": Model matches samples perfectly: too well given the number of samples. seems to capture also statistical fluctuations
example			
$h_{\theta}(x) =$	$g(\theta_0 + \theta_1x_1 + \theta_2x_2)$	$g\left(\begin{matrix} \theta_0 + \theta_1x_1 + \theta_2x_2 + \\ \theta_3x_1^2 + \theta_4x_2^2 + \theta_5x_1x_2 \end{matrix}\right)$	$g\left(\begin{matrix} \theta_0 + \theta_1x_1 + \theta_2x_1^2 + \theta_3x_1^2x_2 \\ \theta_4x_1^2x_2^2 + \theta_5x_1^2x_2^3 + \theta_6x_1^3x_2 \\ + \dots \end{matrix}\right)$
Occam's Razor	"Among competing hypothesis, the simplest is the best." William of Ockham (1285-1347)		
Definition	Overfitting occurs when the learned hypothesis (trained model) fits the training data set very well - but fails to generalise to new examples.		
can occur when	<ul style="list-style-type: none"> <li>- the training set is too noisy</li> <li>- its size is too small in comparison with the dimensionality of the input data</li> <li>- the number of parameters of the model is too large, i.e. the model is too "flexible"</li> </ul>		
examine overfitting			<p>Overfitting occurs when the learned hypothesis (trained model) fits the training data set very well - but fails to generalise to new examples.</p>

### 4. Model Selection Process

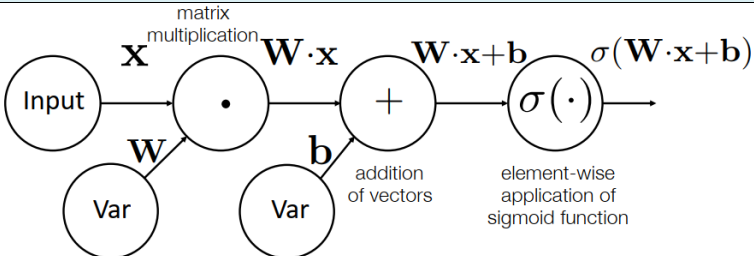
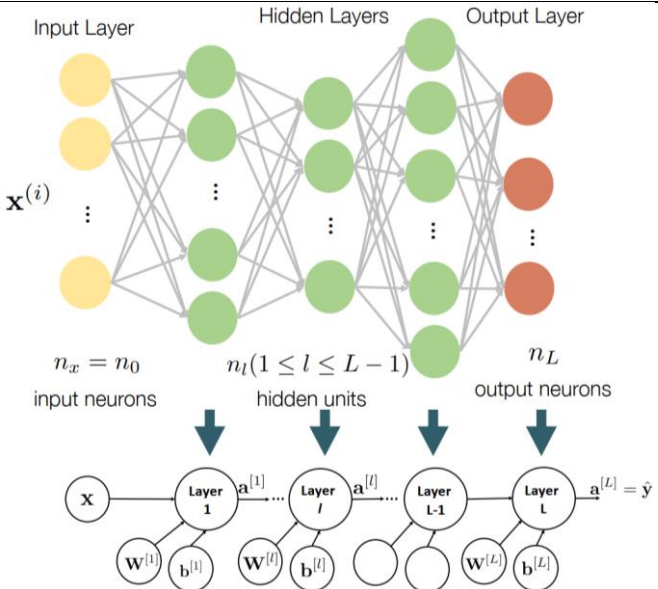
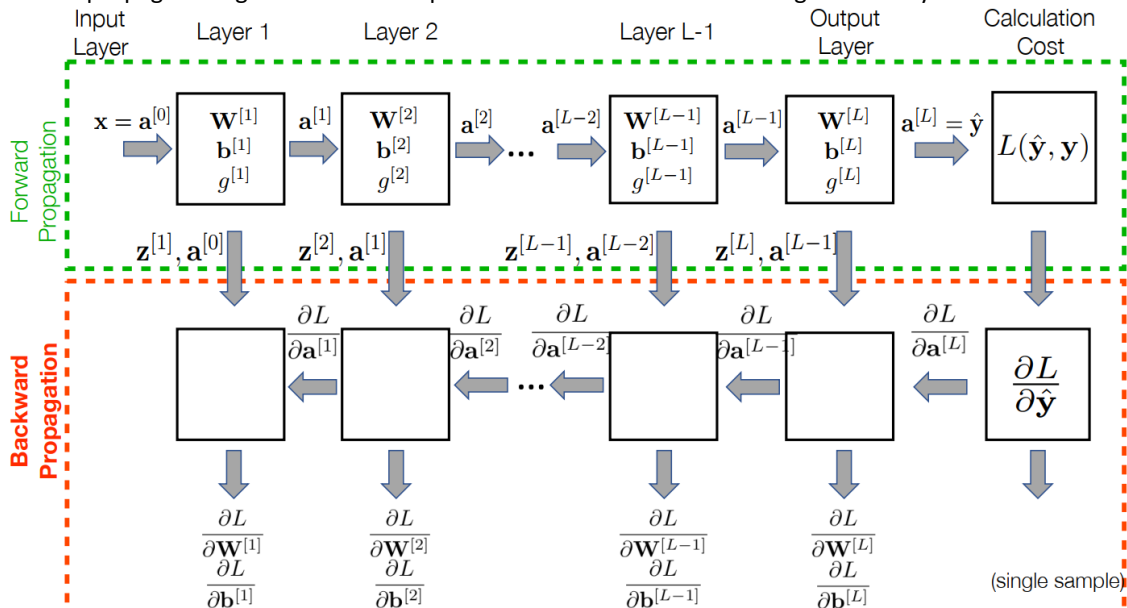
Goal	Select from a family of models the model with the best performance.
Problem	To achieve this goal, we need to evaluate the performance for different models. But to avoid overfitting on the test set we must not use information from the test set to tune params.
Solution	Split the original training set into <b>training set</b> and <b>validation set</b> .
Hyper-Parameter Tuning	Hyper-Parameters specify higher level properties of the mapping function (model) and/or the learning process. These parameters are optimised on the validation set. e.g. learning rate, batch size, number of hidden layers, number of neurons in a layer
learning curves	With more training data it gets more difficult to perfectly fit a model the model trained with more data captures more details about the underlying problem.
Cross-Validation	<p>Split data set in different junks (folds, 5-10) with equal size. -&gt; not that important, see ML courses</p>



**Performance Measures**

<p><b>Confusion Matrix</b></p>	<p>Given a classification system, a confusion matrix evaluates the performance of such system through an <math>m \times m</math> matrix, where <math>m</math> is the number of classes. It shows how many of class "a" were confused as class "b", hence its name "confusion matrix". Very useful to understand the type of errors the system is doing. from sklearn.metrics import confusion_matrix confmat = confusion_matrix(y_actual, y_pred)</p> $\text{Overall Accuracy} = \frac{\sum \text{diagonal elements}}{\# \text{samples}}$ $\text{Error rate} = 1 - \text{Accuracy}$	<p><b>Confusion Matrix Example</b></p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td></td><td></td><th colspan="4">predicted class</th></tr> <tr><td></td><td></td><th>a</th><th>b</th><th>c</th><th>d</th></tr> <tr><th rowspan="4">actual class</th><th>a</th><td>120</td><td>21</td><td>7</td><td>8</td></tr> <tr><th>b</th><td>8</td><td>131</td><td>63</td><td>19</td></tr> <tr><th>c</th><td>12</td><td>30</td><td>80</td><td>11</td></tr> <tr><th>d</th><td>1</td><td>11</td><td>8</td><td>40</td></tr> </table> $\frac{120 + 131 + 80 + 40}{570} = \frac{371}{570} = 65\%$ $1 - 0.65 = 35\%$			predicted class						a	b	c	d	actual class	a	120	21	7	8	b	8	131	63	19	c	12	30	80	11	d	1	11	8	40															
		predicted class																																																
		a	b	c	d																																													
actual class	a	120	21	7	8																																													
	b	8	131	63	19																																													
	c	12	30	80	11																																													
	d	1	11	8	40																																													
<p><b>Confusion Table</b></p>	<p>A <b>Confusion Table</b> is used to measure the classification performance of a <b>two-class</b> system.</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td></td><td></td><th colspan="3">Predicted</th></tr> <tr><td></td><td></td><th>Positive</th><th>Negative</th><th>Total</th></tr> <tr><th rowspan="2">Actual</th><th>Positive</th><td>True Positive</td><td>False Negative</td><td>(TP+FN)</td></tr> <tr><th>Negative</th><td>False Positive</td><td>True Negative</td><td>(FP+TN)</td></tr> <tr><td></td><th>Total</th><td>(TP+FP)</td><td>(FN+TN)</td><td>N</td></tr> </table> $\text{Class accuracy} = \frac{TP + TN}{\# \text{sample}}$ <p>correct classification considering a given class against the others.</p> $\text{Recall} = \text{class sensitivity} = \frac{TP}{TP + FN}$ <p>correct classification for a given class</p> $\text{class precision} = \frac{TP}{TP + FP}$ <p>correct classification in the predicted outputs for a given class</p> $F - \text{Score} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$ $\text{specificity} = \frac{TN}{TN + FP}$			Predicted					Positive	Negative	Total	Actual	Positive	True Positive	False Negative	(TP+FN)	Negative	False Positive	True Negative	(FP+TN)		Total	(TP+FP)	(FN+TN)	N	<p><b>Confusion Table Example (digit 5)</b></p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td></td><td></td><th colspan="3">Predicted</th></tr> <tr><td></td><td></td><th>P</th><th>N</th><th>Total</th></tr> <tr><th rowspan="2">Actual</th><th>P</th><td>809</td><td>112</td><td>921</td></tr> <tr><th>N</th><td>128</td><td>8951</td><td>9079</td></tr> <tr><td></td><th>Total</th><td>937</td><td>9063</td><td>10'000</td></tr> </table> $= \frac{809 + 8951}{10'000} = 97.6\%$ $= \frac{809}{809 + 112} = 87.84\%$ $= \frac{809}{809 + 128} = 86.34\%$ $= \frac{2 * 0.8634 * 0.8784}{0.8634 + 0.8784} = 87.08\%$ $\frac{8951}{8951 + 128} = 98.59\%$			Predicted					P	N	Total	Actual	P	809	112	921	N	128	8951	9079		Total	937	9063	10'000
		Predicted																																																
		Positive	Negative	Total																																														
Actual	Positive	True Positive	False Negative	(TP+FN)																																														
	Negative	False Positive	True Negative	(FP+TN)																																														
	Total	(TP+FP)	(FN+TN)	N																																														
		Predicted																																																
		P	N	Total																																														
Actual	P	809	112	921																																														
	N	128	8951	9079																																														
	Total	937	9063	10'000																																														
<p><b>General Rule for Model Selection</b></p>	<p>Plot curves with these performance measures computed on the training and the validation set, to avoid overfitting, look at the performance on the validation set.</p>																																																	
<p><b>Dimensionality, Curse of Dimensionality</b></p>	<p>This is best understood by how many parameters we need for each point in the Figure to the right.</p> <ul style="list-style-type: none"> <li>In 1D we only have, say 100 points.</li> <li>In 2D we have <math>100^2 = 10'000</math> points</li> <li>In 3D we have <math>100^3 = 1'000'000</math> points</li> <li>Real world application: an image with 300 dpi of 2400x2400 in RGB (3 channels) = 17'280'000 params. And that's just ONE image from your training data.</li> </ul> <p>In machine learning, the more features our data has, the higher dimensionality problem we'll have. A solution for this is called "dimensionality reduction", and different methods like Principal Component Analysis are used.</p>	<p style="text-align: center;"><math>m \propto N^d</math></p> <p><math>m</math>: Number of samples needed  <math>N</math>: Number of points along each dimension  <math>d</math>: Number of parameters of the model</p>																																																
<p><b>Feature Variation, Levels of Hierarchy</b></p>	<p>This brings us to the concept of feature variation and levels of hierarchy. In Feature variation, the idea and solution to the Curse of Dimensionality, is to use a machine learning method that does not really "learn" all the parameters, but instead is robust to variation and finds similar features across the data. images, as the best example, have similar features, and variations from rotations, different pose, moustaches and beards, etc. So, we would build a hierarchy of features as: edges, contours, objects.</p>																																																	

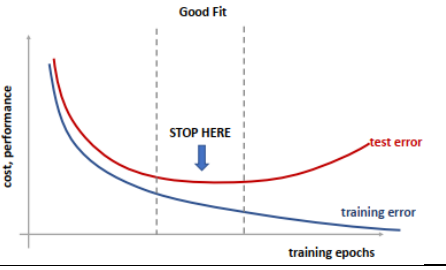
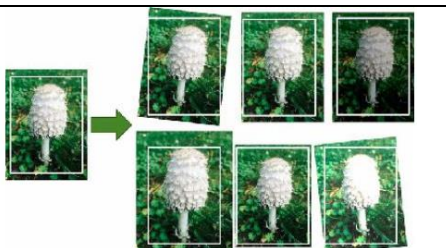
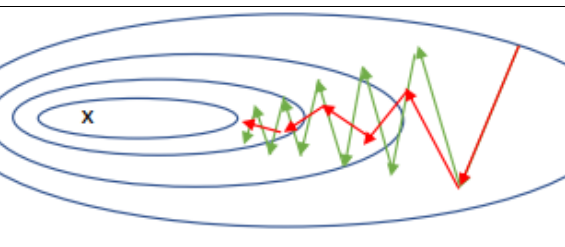
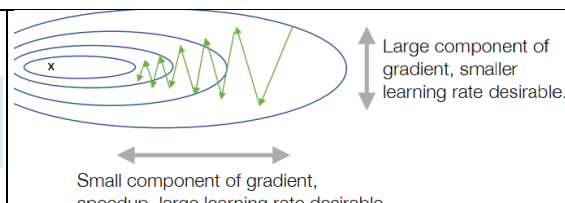
5. Back Propagation

<p><b>Computational Graphs</b></p>	<p>A computational graph is a directed graph where the nodes correspond to operations or variables. Variables can feed their value into operations, and operations can feed their output into other operations. This way every node in the graph defines a function of the variables.</p>	
<p>e.g.</p>		<p>Input: e.g. MNIST pictures                  Var w: weights                  Var b: biases                  Function: e.g. sigmoid</p>
<p>multi-layer one sample                   calculus per layer  <math>x = a^{[0]}; n_x \times 1</math></p>		<p>weights: <math>n_l \times n_{l-1}</math> (1 row for 1 neuron)  <math display="block">W^{[l]} = \begin{pmatrix} w_{11} &amp; \dots &amp; w_{1n_{l-1}} \\ \dots &amp; &amp; \dots \\ w_{n_l1} &amp; \dots &amp; w_{n_l n_{l-1}} \end{pmatrix}</math>                  bias: <math>n_l \times 1</math>  <math display="block">b^{[l]} = \begin{pmatrix} b_1^{[l]} \\ \dots \\ b_{n_l}^{[l]} \end{pmatrix}</math>                  activation: <math>n_{l-1} \times 1</math>  <math display="block">a^{[l-1]} = \begin{pmatrix} a_1^{[l-1]} \\ \dots \\ a_{n_{l-1}}^{[l-1]} \end{pmatrix}</math>                  activation: <math>n_l \times 1</math>  <math display="block">a^{[l]} = \begin{pmatrix} a_1^{[l]} \\ \dots \\ a_{n_l}^{[l]} \end{pmatrix}</math>                  temporary variable: <math>n_l \times 1</math>  <math display="block">z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}</math>                  element wise activation function: <math>n_l \times 1</math>  <math display="block">a^{[l]} = g^{[l]}(z^{[l]})</math></p>
<p>multi-layer multiple sample: m                   calculus per layer  <math>X = A^{[0]}; n_x \times m</math></p>	<p>handle multiple samples with one calculation (add 3rd dimension)</p> $\overset{W}{n_l \times n_x} * \overset{A}{n_x \times m} + \overset{b}{n_l \times m} = \overset{z}{n_l \times m}$	$A^{[l-1]}: n_{l-1} \times m$ $W^{[l]}: n_l \times n_{l-1}$ $z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}: n_l \times m$ $A^{[l]} = g^{[l]}(z^{[l]}): n_l \times m$ $b^{[l]}: n_l \times m \text{ (broadcasting)}$
<p><b>Chain rule</b></p>	<p>Compute the amount of change of a function by changing one of its variables.</p>	$L = L(g(h(x))) \rightarrow dL = \frac{\delta L}{\delta g} \frac{\delta g}{\delta h} \frac{\delta h}{\delta x} dx$
<p><b>Backward Propagation</b>                   based on chain rule</p>	<p>Learning consists in <b>minimising the cost</b> with respect to the model parameters. Gradient descent is an iterative scheme to accomplish that. Update rule with learning rate <math>\alpha</math>. We can propagate all gradients with respect to the activations back through all the layers.</p> 	

## 6. Regularisation and Optimisation

<b>Problem</b>	When learning with backprop, we observe that: - learning is slower in earlier layers (~length of gradients) - learning becomes even slower when having more layers But small gradients here don't imply that we are close to the minimum. Generally, the gradients in deep neural networks are unstable, tending to either vanish (prevalent = dominant) or explode in earlier layers.													
<b>Reason</b>	- <b>Multiplicative Structure</b> -> consequence of the chain rule - <b>Product Term</b> - <b>Coordinating Updates</b> -> updates of the weights in different layers are highly coupled.													
<b>Solution</b>	The problem cannot be completely solved -> but alleviated - <b>Parameter initialisation</b> -> proper initialization of the weights so that the logit z does not grow too large in magnitude. <ul style="list-style-type: none"> <li>randomly initialise weights -&gt; to break symmetries at start of learning</li> <li>put initial weights at proper weights -&gt; normalize z-values in different layer (Xavier)</li> </ul> - <b>Batch Normalisation, Gradient Clipping</b> : Make sure that the weights do not grow too large in magnitude also during training <ul style="list-style-type: none"> <li>Normalise the input to each layer per mini-batch by adding an operation in the model just before the activation function of each layer: zero-centre and scale the inputs by estimating mean and stdev from the current mini-batch.</li> </ul> <table border="1" data-bbox="448 770 1311 902"> <tr> <td style="padding: 5px;">                     Normalisation  <math display="block">\left(Z_{norm}^{\{r\},[l]}\right)_{k,i} = \frac{Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}}{\sigma_k^{\{r\},[l]} + \epsilon}</math> </td> <td style="padding: 5px;"> <math display="block">\mu_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} Z_{k,i}^{\{r\},[l]}</math> <math display="block">\sigma_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} \left(Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}\right)^2</math> </td> </tr> </table> <ul style="list-style-type: none"> <li>Then let the model learn the optimal scale and mean of the inputs for each layer.</li> </ul> <table border="1" data-bbox="448 936 1311 1025"> <tr> <td style="padding: 5px;">                     Scaling and Shifting  <math display="block">\hat{Z}_{k,i}^{\{r\},[l]} = \gamma_k^{[l]} \left(Z_{norm}^{\{r\},[l]}\right)_{k,i} + \beta_k^{[l]}</math> </td> <td style="padding: 5px;"> <math>\gamma_k^{[l]}</math> and <math>\beta_k^{[l]}</math> have to be learned                 </td> </tr> </table> <p>-&gt; can be applied to any input or hidden layer in a network.</p> <ul style="list-style-type: none"> <li>Or clip the gradient in length during backprop so that they never exceed some threshold. Because neural networks or recurrent networks often have an extremely steep cliff structure.</li> </ul> - <b>Suitable Activation Functions</b> : Use activation functions that cannot saturate. <ul style="list-style-type: none"> <li>use ReLU, LeakyReLU or ELU</li> <li><b>ReLU</b> suffers the dying unit's problem: During training, if a neuron's weight gets updated such that the weighted sum of the neuron's inputs is negative, it is outputting 0. Since the gradient at <math>z &lt; 0</math> is 0, there is no weight update for this neuron and the neuron is likely to stay dead.</li> <li>With the <b>leaky ReLU</b> (or its smooth version, the <b>ELU</b>) this problem cannot occur.</li> </ul>	Normalisation $\left(Z_{norm}^{\{r\},[l]}\right)_{k,i} = \frac{Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}}{\sigma_k^{\{r\},[l]} + \epsilon}$	$\mu_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} Z_{k,i}^{\{r\},[l]}$ $\sigma_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} \left(Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}\right)^2$	Scaling and Shifting $\hat{Z}_{k,i}^{\{r\},[l]} = \gamma_k^{[l]} \left(Z_{norm}^{\{r\},[l]}\right)_{k,i} + \beta_k^{[l]}$	$\gamma_k^{[l]}$ and $\beta_k^{[l]}$ have to be learned									
Normalisation $\left(Z_{norm}^{\{r\},[l]}\right)_{k,i} = \frac{Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}}{\sigma_k^{\{r\},[l]} + \epsilon}$	$\mu_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} Z_{k,i}^{\{r\},[l]}$ $\sigma_k^{\{r\},[l]} = \frac{1}{N_B} \sum_{i=1}^{N_B} \left(Z_{k,i}^{\{r\},[l]} - \mu_k^{\{r\},[l]}\right)^2$													
Scaling and Shifting $\hat{Z}_{k,i}^{\{r\},[l]} = \gamma_k^{[l]} \left(Z_{norm}^{\{r\},[l]}\right)_{k,i} + \beta_k^{[l]}$	$\gamma_k^{[l]}$ and $\beta_k^{[l]}$ have to be learned													
<b>Regularisation</b>	<b>Problem</b> : Deep neural networks typically have many parameters to fit a huge variety of complex datasets. But bear the risk to <b>overfitting</b> the training set. "Regularisation is any modification to a learning algorithm that is intended to reduce its generalisation error but not its training error."													
Weight Decay (Weights Penalties)	Add constraints to the parameters to give preference to simple models - restriction in the length of the parameter vector or in number of parameters (sparsity). <table border="1" data-bbox="352 1525 1481 1688"> <tr> <td style="padding: 5px;"><math>J = J_{data} + J_{reg}</math></td> <td style="padding: 5px;"> <math>J_{data}</math>: Performance term: on the data, how far are we from the ground truth  <math>J_{reg}</math>: Regularisation term: we should not let our model become too complex                 </td> </tr> <tr> <td style="padding: 5px;"><math>J = J_0 + \lambda \Omega(\mathbf{W})</math></td> <td style="padding: 5px;"> <math>J_0</math>: original loss function (e.g. RMS loss)  <math>\Omega(\mathbf{W})</math>: penalty term that favours models with smaller weights  <math>\lambda</math>: regularisation parameter                 </td> </tr> </table> <p>Two forms of penalties are common</p> <table border="1" data-bbox="352 1722 1273 1955"> <thead> <tr> <th></th> <th>Penalty term</th> <th>Gradient for the regularised loss function:</th> </tr> </thead> <tbody> <tr> <td style="text-align: center;"><math>L_1</math></td> <td style="padding: 5px;"><math>\Omega(\mathbf{W}) = \ \mathbf{W}\ _1 = \sum_{l,k,j}  W_{kj}^{[l]} </math></td> <td style="padding: 5px;"><math>\nabla(J_0(\mathbf{W}) + \lambda \ \mathbf{W}\ _1) = \nabla J_0(\mathbf{W}) + \lambda \text{sign}(\mathbf{W})</math></td> </tr> <tr> <td style="text-align: center;"><math>L_2</math></td> <td style="padding: 5px;"><math>\Omega(\mathbf{W}) = \ \mathbf{W}\ _2^2 = \sum_{l,k,j}  W_{kj}^{[l]} ^2</math></td> <td style="padding: 5px;"> <math>\nabla \left( J_0(\mathbf{W}) + \frac{\lambda}{2} \ \mathbf{W}\ _2^2 \right) = \nabla J_0(\mathbf{W}) + \lambda \mathbf{W}</math>                      update rule for gradient descent  <math>\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla J(\mathbf{W})</math> </td> </tr> </tbody> </table>	$J = J_{data} + J_{reg}$	$J_{data}$ : Performance term: on the data, how far are we from the ground truth $J_{reg}$ : Regularisation term: we should not let our model become too complex	$J = J_0 + \lambda \Omega(\mathbf{W})$	$J_0$ : original loss function (e.g. RMS loss) $\Omega(\mathbf{W})$ : penalty term that favours models with smaller weights $\lambda$ : regularisation parameter		Penalty term	Gradient for the regularised loss function:	$L_1$	$\Omega(\mathbf{W}) = \ \mathbf{W}\ _1 = \sum_{l,k,j}  W_{kj}^{[l]} $	$\nabla(J_0(\mathbf{W}) + \lambda \ \mathbf{W}\ _1) = \nabla J_0(\mathbf{W}) + \lambda \text{sign}(\mathbf{W})$	$L_2$	$\Omega(\mathbf{W}) = \ \mathbf{W}\ _2^2 = \sum_{l,k,j}  W_{kj}^{[l]} ^2$	$\nabla \left( J_0(\mathbf{W}) + \frac{\lambda}{2} \ \mathbf{W}\ _2^2 \right) = \nabla J_0(\mathbf{W}) + \lambda \mathbf{W}$ update rule for gradient descent $\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla J(\mathbf{W})$
$J = J_{data} + J_{reg}$	$J_{data}$ : Performance term: on the data, how far are we from the ground truth $J_{reg}$ : Regularisation term: we should not let our model become too complex													
$J = J_0 + \lambda \Omega(\mathbf{W})$	$J_0$ : original loss function (e.g. RMS loss) $\Omega(\mathbf{W})$ : penalty term that favours models with smaller weights $\lambda$ : regularisation parameter													
	Penalty term	Gradient for the regularised loss function:												
$L_1$	$\Omega(\mathbf{W}) = \ \mathbf{W}\ _1 = \sum_{l,k,j}  W_{kj}^{[l]} $	$\nabla(J_0(\mathbf{W}) + \lambda \ \mathbf{W}\ _1) = \nabla J_0(\mathbf{W}) + \lambda \text{sign}(\mathbf{W})$												
$L_2$	$\Omega(\mathbf{W}) = \ \mathbf{W}\ _2^2 = \sum_{l,k,j}  W_{kj}^{[l]} ^2$	$\nabla \left( J_0(\mathbf{W}) + \frac{\lambda}{2} \ \mathbf{W}\ _2^2 \right) = \nabla J_0(\mathbf{W}) + \lambda \mathbf{W}$ update rule for gradient descent $\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla J(\mathbf{W})$												



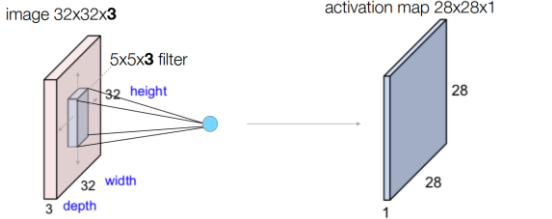
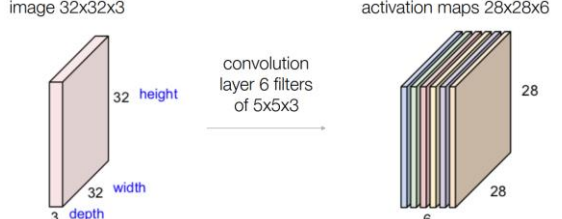
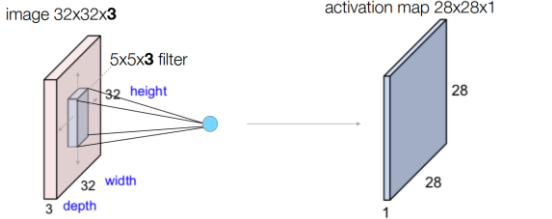
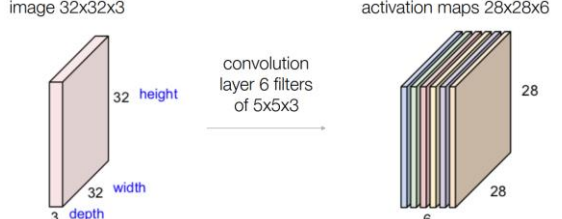
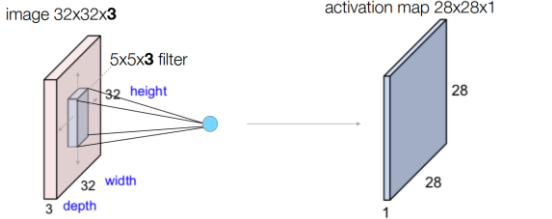
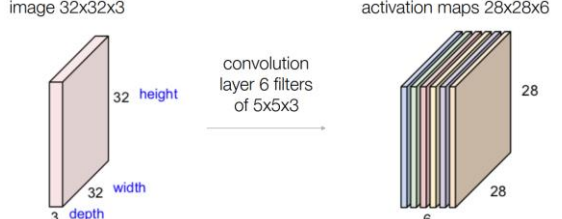
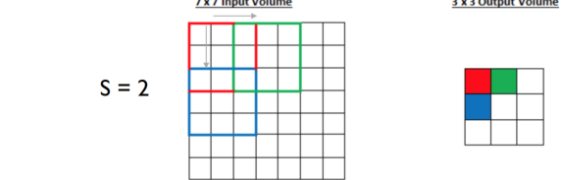
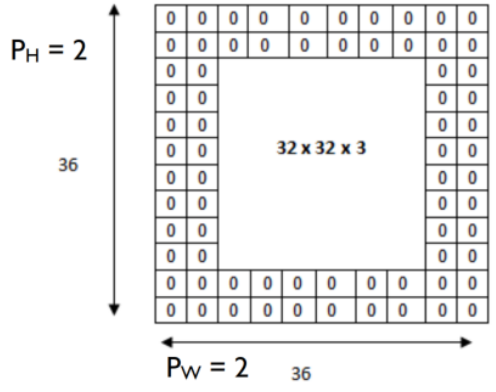
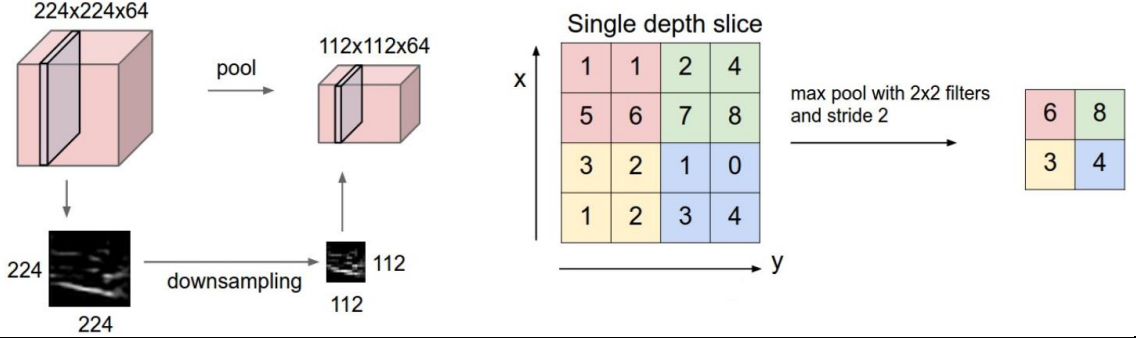


Dropout	<p>Randomly drop neurons during training steps to make the solution less dependent on individual neurons</p> <p>Popular, successful, imply simpler models, more robust, cheap computation at training, very versatile (=vielfältig), can be combined, need to be bit larger</p> <p>At every training step, every neuron (input by 20%, hidden by 50%) has a probability <math>p</math> of being temporarily -&gt; masked activations. Hyper-parameter <math>p</math> is called dropout rate.</p> <p>Correct weights after training, because they have been tuned with dropout rate.</p>																
Early Stopping	<p>Stop training at the minimum of the cost function on the validation set to avoid overfitting.</p> <p>Often, the training error monotonically decreases while the validation error begins to increase after a certain number of epochs. A behaviour that can only be observed when training large models with sufficient representational capacity so that overfitting is possible.</p> <ul style="list-style-type: none"> <li>- Run optimisation algorithm to train the model - simultaneously compute the validation set error.</li> <li>- Store a copy of the model parameters as long as the validation set error improves.</li> <li>- Iterate until validation set error stops improving (e.g. for <math>k</math> steps)</li> <li>- Return the parameters where the smallest validation set error is observed.</li> </ul> <p>efficient, non-intrusive (aufdringlich), parallelizable, combinable</p>																
Data Augmentation	<p>Generate more training data to introduce additional characteristic features (e.g. symmetries) the solution should have.</p> <p>Reduce model complexity.</p> <p>Increase the amount of training data (realistic).</p> <p>Model is forced to be more tolerant to position, orientation, size, contrast, etc.</p> <p>Typically used for classification tasks.</p> <p>Be careful with flips and 180° rotation (d-b, 9-6)</p> <p>May be considered rather as pre-processing step.</p> <p>But can be applied on the fly. Reduce network bandwidth</p>																
Advanced Optimisation Methods	<p>Gradient Descent</p> <p>+ works for convex function -&gt; but cost function in deep neural networks are typically non-convex.</p> <p>- can get stuck in a local minimum or saddle point -&gt; hard to find in high dimensional spaces.</p> <p>- can get very slow -&gt; faster optimisers desirable</p>																
Momentum, Nesterov	<p>Allows to surpass flat regions or saddle points - like a ball that keeps on rolling down if it has an initial speed when entering flat regions.</p> <p>Compute an exponentially decay moving average of past gradients and move in the direction of the moving average. 'Momentum' hyper parameter which controls the decay and the friction.</p> <p>Momentum <math>\mathbf{m} \leftarrow \beta_1 \mathbf{m} + \alpha \nabla J(\theta)</math></p> <p>Nesterov <math>\mathbf{m} \leftarrow \beta_1 \mathbf{m} + \alpha \nabla J(\theta + \beta_1 \mathbf{m})</math></p> <p><math>\theta \leftarrow \theta - \mathbf{m}</math></p> <p>good beta: <math>\beta_1 = 0.9</math></p>																
RMS Prop	<p>Adaptively adjusts the learning rate to incorporate differences in the steepness along different directions in parameter space.</p> <p>Increase learning rate in direction of slow progress and decrease in direction of fast progress.</p> <p><math>s \leftarrow \beta_2 s + (1 - \beta_2) \nabla J(\theta) \odot \nabla J(\theta)</math></p> <p><math>\theta \leftarrow \theta - \frac{\alpha}{\sqrt{s + \epsilon}} \odot \nabla J(\theta)</math></p> <p><math>\odot</math> elementwise operation</p>	 <p>Large component of gradient, smaller learning rate desirable.</p> <p>Small component of gradient, speedup, large learning rate desirable.</p>															
Adam (state of the art)	<p>Combination of Momentum / Nesterov and RMSprop</p> <div style="display: flex; align-items: flex-start;"> <div style="flex: 1;"> <p><math>\mathbf{m} \leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \nabla J(\theta)</math></p> <p><math>s \leftarrow \beta_2 s + (1 - \beta_2) \nabla J(\theta) \odot \nabla J(\theta)</math></p> <p><math>\hat{\mathbf{m}} = \frac{\mathbf{m}}{1 - \beta_1}</math> (init with 0)</p> <p><math>\hat{s} = \frac{s}{1 - \beta_2}</math> (init with 0)</p> <p><math>\theta \leftarrow \theta - \alpha \frac{\hat{\mathbf{m}}}{\sqrt{\hat{s} + \epsilon}} \odot \nabla J(\theta)</math></p> </div> <table border="1" style="margin-left: 10px;"> <tr> <td>Learning Rate</td> <td><math>\alpha</math></td> <td>0.001</td> </tr> <tr> <td>Momentum</td> <td><math>\beta_1</math></td> <td>0.9</td> </tr> <tr> <td>RMS Decay</td> <td><math>\beta_2</math></td> <td>0.999</td> </tr> <tr> <td>Numerical Stabilisation</td> <td><math>\epsilon</math></td> <td>1.0E-08</td> </tr> <tr> <td>Learning rate</td> <td><math>\alpha</math></td> <td>(reduce after epochs)</td> </tr> </table> </div>		Learning Rate	$\alpha$	0.001	Momentum	$\beta_1$	0.9	RMS Decay	$\beta_2$	0.999	Numerical Stabilisation	$\epsilon$	1.0E-08	Learning rate	$\alpha$	(reduce after epochs)
Learning Rate	$\alpha$	0.001															
Momentum	$\beta_1$	0.9															
RMS Decay	$\beta_2$	0.999															
Numerical Stabilisation	$\epsilon$	1.0E-08															
Learning rate	$\alpha$	(reduce after epochs)															

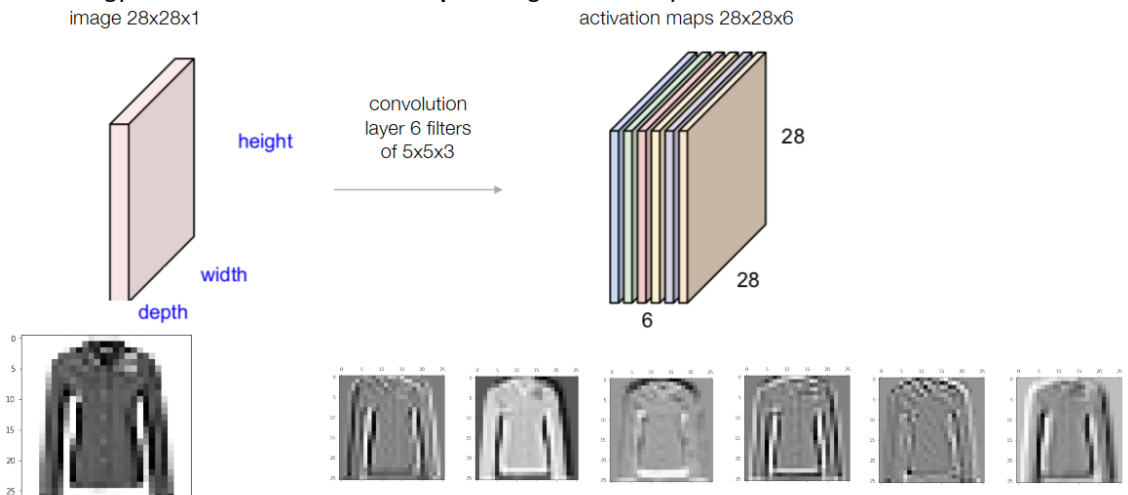
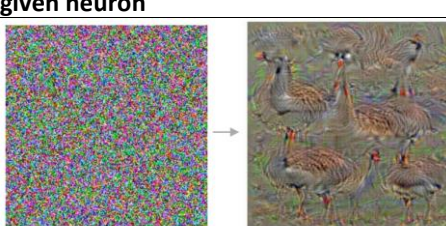





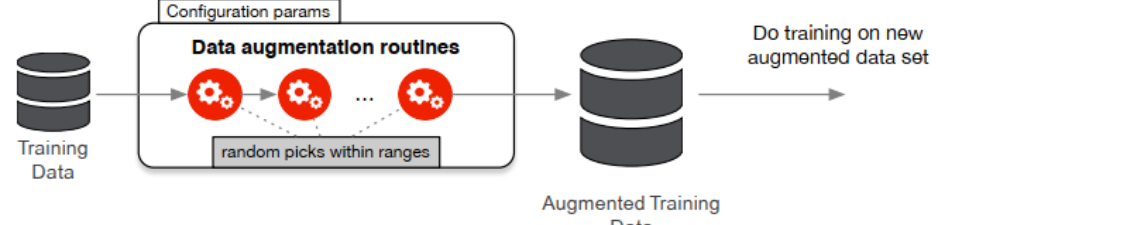

**7. DL-frameworks**

<p><b>CPU vs GPU</b></p> <p><b>TPU</b></p> <p><b>GPU</b></p>	<p>Central Processing Unit: Few cores (~10), fast (~4GHz), lots of cache, few parallel processes</p> <p>Graphical Processing Unit: Many cores (~1'000), slow (~1.5GHz), few caches, many parallel processes</p> <p>GPU is suitable for matrix multiplication -&gt; because it's highly parallelizable.</p> <p>Tensor Processing Unit: even more specialised hardware to perform matrix and convolution operations</p>											
<p><b>GPU</b></p>	<p>CUDA - for NVIDIA only - low level API for programming GPU</p> <p>OpenCL - like CUDA but runs on anything - usually slower on NVIDIA hardware</p> <p>HIP - Write code once in HIP C++ and port on NVIDIA or AMD</p>											
<p><b>Forward</b></p> <p><b>Backward</b></p> <p><b>Node principles</b></p> <p>add = distributor max = router mul = switcher</p>		<table border="1" data-bbox="346 750 1497 981"> <tr> <td data-bbox="346 750 422 846">→</td> <td data-bbox="422 750 667 846"> <math>x = 2</math>  <math>w = 3</math>  <math>b = -4</math> </td> <td data-bbox="667 750 970 846"> <math>q = x * w = 2 * 3</math>  <math>q = 6</math> </td> <td data-bbox="970 750 1273 846"> <math>z = q + b = 6 + (-4)</math>  <math>z = 2</math> </td> <td data-bbox="1273 750 1497 846"> <math>f = z^2 = 2^2</math>  <math>f = 4</math> </td> </tr> <tr> <td data-bbox="346 846 422 981">←</td> <td data-bbox="422 846 667 981"> <math>\frac{\delta q}{\delta x} \frac{\delta z}{\delta q} = 3 * 4 = 12</math>  <math>\frac{\delta q}{\delta w} \frac{\delta z}{\delta q} = 2 * 4 = 8</math> </td> <td data-bbox="667 846 970 981"> <math>\frac{\delta z}{\delta q} \frac{\delta f}{\delta z} = 1 * 4 = 4</math>  <math>\frac{\delta z}{\delta b} \frac{\delta f}{\delta z} = 1 * 4 = 4</math> </td> <td data-bbox="970 846 1273 981"> <math>\frac{\delta f}{\delta z} \frac{\delta f}{\delta z} = 2z * 1 = 4</math> </td> <td data-bbox="1273 846 1497 981"> <math>\frac{\delta f}{\delta f} = 1</math> </td> </tr> </table>	→	$x = 2$ $w = 3$ $b = -4$	$q = x * w = 2 * 3$ $q = 6$	$z = q + b = 6 + (-4)$ $z = 2$	$f = z^2 = 2^2$ $f = 4$	←	$\frac{\delta q}{\delta x} \frac{\delta z}{\delta q} = 3 * 4 = 12$ $\frac{\delta q}{\delta w} \frac{\delta z}{\delta q} = 2 * 4 = 8$	$\frac{\delta z}{\delta q} \frac{\delta f}{\delta z} = 1 * 4 = 4$ $\frac{\delta z}{\delta b} \frac{\delta f}{\delta z} = 1 * 4 = 4$	$\frac{\delta f}{\delta z} \frac{\delta f}{\delta z} = 2z * 1 = 4$	$\frac{\delta f}{\delta f} = 1$
→	$x = 2$ $w = 3$ $b = -4$	$q = x * w = 2 * 3$ $q = 6$	$z = q + b = 6 + (-4)$ $z = 2$	$f = z^2 = 2^2$ $f = 4$								
←	$\frac{\delta q}{\delta x} \frac{\delta z}{\delta q} = 3 * 4 = 12$ $\frac{\delta q}{\delta w} \frac{\delta z}{\delta q} = 2 * 4 = 8$	$\frac{\delta z}{\delta q} \frac{\delta f}{\delta z} = 1 * 4 = 4$ $\frac{\delta z}{\delta b} \frac{\delta f}{\delta z} = 1 * 4 = 4$	$\frac{\delta f}{\delta z} \frac{\delta f}{\delta z} = 2z * 1 = 4$	$\frac{\delta f}{\delta f} = 1$								
<p><b>Node composition or factorisation</b></p>	<p>A sub-graph composed of nodes can be reimplemented in a single node if we can compute an analytic form of the gradients. Any complex learning architecture can be composed from atomic nodes. No need to compute complex global gradients. The loss function can be seen as extra nodes.</p>	<p><math>\frac{\delta \sigma}{\delta z}</math></p>										
<p><b>A zoo of frameworks</b></p>	<p><b>Academia</b></p> <p><b>Theano</b> - U Montreal</p> <p><b>Torch</b> - IDIAP/NYU/...</p> <p><b>Caffe</b> - UC Berkeley</p> <p><b>MXNet</b> - CMU, MIT, U Wash, HK UST</p>	<p><b>Companies</b></p> <p><b>TensorFlow</b> (based on Theano, static)- Google</p> <p><b>Pytorch</b> (based on Torch, dynamic)- Facebook</p> <p><b>Caffe2</b> - Facebook</p> <p><b>MXNet</b> - Amazon, Intel</p> <p><b>DeepLearning4J</b> - Skymind</p> <p><b>CNTK</b> - Microsoft</p> <p><b>PaddlePaddle</b> - Baidu</p>										
<p><b>Advantages</b></p>	<ul style="list-style-type: none"> <li>- Easily build big computational graphs</li> <li>- Easily compute losses and gradients in computation graphs for update rules</li> <li>- Have at hand all the state-of-art strategies for regularisations and optimisations</li> <li>- Switch easily from cpu to gpu when needed.</li> </ul>											
<p><b>Problem</b></p>	<p>The loss is not going down -&gt; assign calls are not executed</p> <p>-&gt; add dummy graph node that depends on the update</p>											

**8. Keras + CNN**

<p><b>Keras</b></p> 	<p>is a high-level open-source neural networks API written in Python and capable of running on top of TensorFlow, CNTK, Theano. It was developed with a focus on enabling fast experimentation.</p> <ul style="list-style-type: none"> <li>- minimalistic, straight forward, extensive, python, simple, good documentation, large community</li> </ul> <p><a href="https://keras.io">https://keras.io</a></p>							
<p>pipeline</p>								
<p>Models</p>	<p>A model in keras is the way to organise the layers of neurons: sequential or graph          The <b>sequential</b> model corresponds to a regular stack of layers (1 layer = 1 object that feeds to the next)          The <b>graph</b> model is used for non sequential architectures (diverge or merge networks)</p>							
<p><b>CNN - Convolutional Neural Networks</b></p>	<p>General idea: let's define new type of layers and connections that will bring</p> <ul style="list-style-type: none"> <li>- preservation of spatial (=räumlich) structure</li> <li>- hierarchical feature detection - objects are composed of features that are themselves composed of other features</li> <li>- robustness to object variabilities such as viewpoint, occlusion (=Verdeckung), etc</li> </ul>							
<p><b>Convolution layers</b> feature extracting!  convolve (=falten)</p>	<p>have the property to preserve the spatial structure and discover the local connectivity.          A given filter is translated on <b>receptive fields</b> of the input and produces as output an <b>activation map</b>.</p> <table border="1" data-bbox="352 674 1493 936"> <tr> <td data-bbox="352 674 916 936"> <p>convolution filter results in an activation map.</p>  </td> <td data-bbox="920 674 1493 936"> <p>several filters producing several activation maps</p>  </td> </tr> </table>		<p>convolution filter results in an activation map.</p> 	<p>several filters producing several activation maps</p> 				
<p>convolution filter results in an activation map.</p> 	<p>several filters producing several activation maps</p> 							
<p>stride s</p>	<p>The stride specifies a step size when moving the filter across the signal. Larger stride means less overlap and reduction of the output volume. some s might be incompatible.</p>							
<p>padding p</p>	<p>is the size of a zeroed frame added around the input.</p> <table border="1" data-bbox="352 1173 981 1518"> <tr> <td><math>P_W = \frac{w - 1}{2}</math></td> <td rowspan="5"> <math>P_W</math>: padding in width  <math>P_H</math>: padding in height  <math>h</math>: height of filter  <math>w</math>: width of filter  <math>H</math>: height of image  <math>W</math>: width of image  <math>o_W</math>: output width  <math>o_H</math>: output height  <math>C</math>: # of channels  <math>D</math>: # of filters  <math>N_{param}</math>: # of params                 </td> </tr> <tr> <td><math>P_H = \frac{h - 1}{2}</math></td> </tr> <tr> <td><math>o_W = \frac{(W - w + 2P_W)}{s_w} + 1</math></td> </tr> <tr> <td><math>o_H = \frac{(H - h + 2P_H)}{s_h} + 1</math></td> </tr> <tr> <td><math>N_{param} = D \left( \underbrace{whC}_{weight} + \underbrace{1}_{bias} \right)</math></td> </tr> </table>	$P_W = \frac{w - 1}{2}$	$P_W$ : padding in width $P_H$ : padding in height $h$ : height of filter $w$ : width of filter $H$ : height of image $W$ : width of image $o_W$ : output width $o_H$ : output height $C$ : # of channels $D$ : # of filters $N_{param}$ : # of params	$P_H = \frac{h - 1}{2}$	$o_W = \frac{(W - w + 2P_W)}{s_w} + 1$	$o_H = \frac{(H - h + 2P_H)}{s_h} + 1$	$N_{param} = D \left( \underbrace{whC}_{weight} + \underbrace{1}_{bias} \right)$	
$P_W = \frac{w - 1}{2}$	$P_W$ : padding in width $P_H$ : padding in height $h$ : height of filter $w$ : width of filter $H$ : height of image $W$ : width of image $o_W$ : output width $o_H$ : output height $C$ : # of channels $D$ : # of filters $N_{param}$ : # of params							
$P_H = \frac{h - 1}{2}$								
$o_W = \frac{(W - w + 2P_W)}{s_w} + 1$								
$o_H = \frac{(H - h + 2P_H)}{s_h} + 1$								
$N_{param} = D \left( \underbrace{whC}_{weight} + \underbrace{1}_{bias} \right)$								
<p><b>Max pooling layer</b></p>	<p>Reduce the spatial size of the representation. Applies independently to every depth.          Defined by a stride S and a padding P.          Most significant activations are kept, reduce the amount of computation and control overfitting.</p> 							
<p><b>Dense layers</b></p>	<p>regular fully connected layers usually used as the last layers in the architecture to take classification decisions</p>							

## 9. Deep Architectures

<p><b>CNN for FashionMNIST</b></p>	<p>60'000 train images, 10'000 test images, 28x28x1 grayscale, 10 classes                  Alternative to MNIST which is a bit more challenging, <a href="https://github.com/zalandoresearch/fashion-mnist">https://github.com/zalandoresearch/fashion-mnist</a>                  Conv2D+Relu, MaxPooling2D, Dropout, Flatten, FC + ReLU, FC + SoftMax -&gt; 91.04%</p>				
<p><b>Going deeper</b></p>	<p>Means stacking CONV-RELU-POOL(+DROPOUT) layers and playing with the different hyper-parameters.                  The deeper the network                  - the bigger the number of filters in the conv layers                  - the smaller are the activation maps due to max pooling                  Early layers will extract lower-level features -&gt; late layers will extract higher-level features.</p>				
<p><b>Deep Learning</b></p>	<p>The whole process (pre-processing + feature extraction) is learned not only to classify/train</p>				
<p><b>Visualisation</b></p>	<p>1st Strategy: <b>Visualize the activation maps</b> for a given conv layer to see feature extraction.</p>  <p>2nd Strategy: Find <b>input images that maximise the action of a given neuron</b></p> <ol style="list-style-type: none"> <li>1. Start with a random image <math>x</math></li> <li>2. Forward: compute the activation of the neuron <math>a_{i,j}(x)</math></li> <li>3. Backward: perform the backprop of the gradient of <math>a_{i,j}(x)</math> up to the pixel values: <math>\frac{\delta a_{i,j}(x)}{\delta x}</math></li> <li>4. This gradient tells us how to change the pixel values to increase the activation of the neuron. We can then apply an update rule in the form of a gradient ascent:</li> </ol> $x \leftarrow x + a \frac{\delta a_{i,j}(x)}{\delta x} + Regularisation Term$ 				
<p><b>Data Augmentation</b></p>	<p>takes the approach of generating artificially more training data from existing training samples.</p>				
<p>Rotation (angle)</p>	<p>Translation (shift)</p>	<p>Flip (not to all applicable)</p>	<p>Shear</p>	<p>Zoom</p>	<p>Colour</p>
					<p>???</p>
<p><b>Pre-Augmentation (offline):</b> Do training on new augmented data set.</p>					
					
<p><b>Online Augmentation:</b> Do training using images going out of the augmentation pipeline.</p>					
					
<p>Overall improvement is not so high (because stability of MNIST data), but we avoid overfitting.</p>					

<b>Convolutional layer patterns</b>	Common form: INPUT -> [[CONV -> RELU]*N -> POOL?]*M -> [FC -> RELU]*K -> FC $0 \leq N \leq 3, \quad 0 \leq M, \quad 0 \leq K \leq 3$ ILSVRC: ImageNet Large Scale Visual Recognition -> Challenge				
famous				<b>accuracy</b>	<b>layers</b>
	LeNet-5	1998	CONV-POOL-CONV-POOL-FC-FC		shallow
	AlexNet	2012	the "boot" of deep architectures	16.4%	8
	ZFNet	2013	AlexNet with Hyper parameter tuning	11.7%	8
	VGGNet	2014	going deeper with simpler filter	7.3%	19
	GoogLeNet	2014	Network in the Network (Inception), no FC, only 5M param	6.7%	22
	ResNet	2015	way deeper. more difficult to optimize -> "fall back" layers to go deeper than 50 -> use bottleneck layers	3.57%	152
	Shao et al	2016		3%	152
	SENet	2017		2.3%	152
<b>fall back layers</b>	<b>Hypothesis:</b> Deeper models are more difficult to optimize <b>Ideas:</b> Deeper layers should be able to "fall back" to identity mapping if difficulties to converge. Easier to model a "delta" from one layer to the other than a full feature. <b>Solution:</b> Use network layers to fit a <b>residual</b> (=Rest) mapping instead of directly trying to fit a desired underlying mapping.				
<b>bottleneck layers</b>	Use bottleneck layers 1x1x64 to reduce the number of operations and parameters				






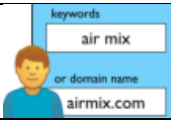

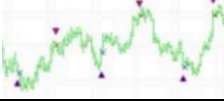
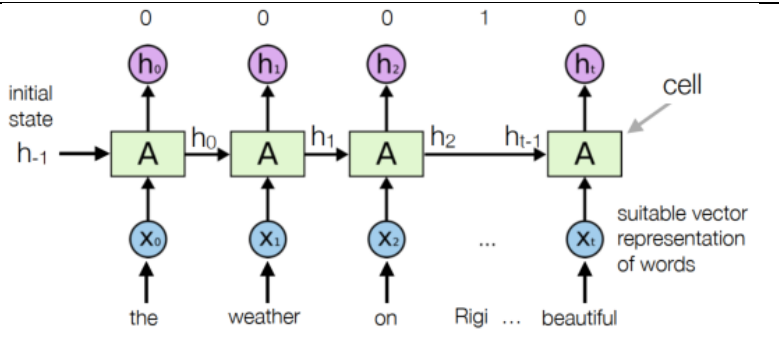
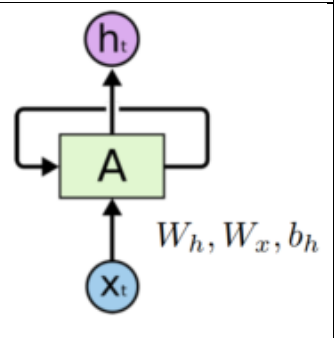


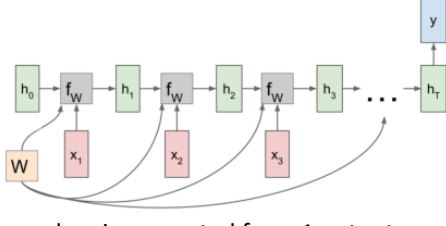
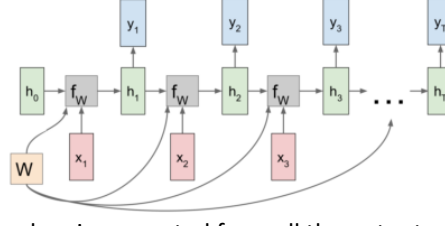
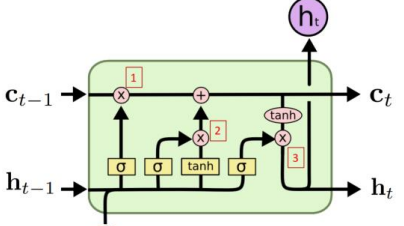
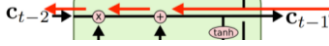
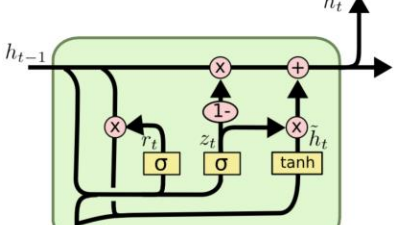
10. CNN3 Keras Functional API, Transfer Learning, Autoencoders

<p><b>Keras Model types</b></p>	<p>The <b>sequential</b> model corresponds to a regular stack of layers.</p> <div data-bbox="319 212 957 638"> <p><b>Keras Sequential Model - a sequence of layers</b></p> </div>	<p>The <b>functional</b> API is used for non-sequential architectures (multiple paths/inputs/outputs).</p> <div data-bbox="1021 212 1436 728"> </div>
<p>python callable syntax</p>	<pre>add1 = Adder(1) // calls __init__ and __call__ add1(2) // calls __call__</pre>	<pre>def __init__(self,x=0):     self.__memory = x def __call__(self,x):     self.__memory += x     return self</pre>
<p><b>Functional API</b></p>	<pre>visible = Input(shape=(28,28,1)) conv1 = Conv2D(32, kernel_size=3, activation='relu')(visible) pool1 = MaxPooling2D(pool_size=(2,2))(conv1)</pre>	
<p>split (multiple path)</p>	<pre>conv1 = Conv2D(32, kernel_size=3, activation='relu')(visible) conv2 = Conv2D(32, kernel_size=5, activation='relu')(visible)</pre>	
<p>merge</p>	<pre>merge = concatenate([conv1, conv2])</pre> <p>e.g. multiple images of one object, stereo speech recording, different parameters of acquisition device</p>	
<p>callbacks</p>	<p>Callback are functions to be applied at given stages of the training procedure. To get a view on internal states and statistics. Declare them in fit() function.</p> <pre>checkpoint = ModelCheckpoint('model-(epoch:03d).h5', verbose=1, monitor='val_acc', save_best_only=True, mode='auto') log = model5.fit(..., calbacks=[checkpoint])</pre>	
<p>Best practices</p>	<p>Use <b>Consistent Variable Names</b>. Use same variable name for input (visible) and output layers (output), hidden layers (hidden1, hidden2). It will help to connect things together correctly.</p> <p><b>Review Layer Summary</b>. Always print the model summary and review the layer outputs to ensure that the model was connected as you expected.</p> <p><b>Review Graph Plots</b>. Create a plot of the model graph and review it to ensure that everything was put together as you intended.</p> <p><b>Name the layers</b>. You can assign names to layers that are used when reviewing summaries and plots of the model graph. For example: Dense(1, name='hidden1').</p> <p><b>Use Callbacks</b>. You should use callbacks to save the (best) models from epochs to epochs as a safety against interrupt and overfitting</p>	

<p><b>Transfer learning</b></p>	<p>is about using knowledge learned from tasks for which a lot of labelled data is available in settings where only little labelled data is available. -&gt; Re-use the feature extraction part. Not too different.</p>	
<p>workflow</p>	<p>The diagram illustrates the transfer learning workflow in three stages:</p> <ol style="list-style-type: none"> <li><b>1. Train on Imagenet:</b> A full neural network is trained on the Imagenet dataset. The architecture includes an input image, followed by two stages of convolutional layers (conv-64, conv-128, conv-256) and maxpooling, and finally three fully connected layers (FC-4096, FC-4096, FC-1000) with a softmax output.</li> <li><b>2. Small dataset: feature extractor:</b> The feature extractor part (conv-64 to conv-512) is frozen. Only the final fully connected layers (FC-4096, FC-4096, FC-1000) and the softmax are trained on a small dataset.</li> <li><b>3. Medium dataset: finetuning:</b> The feature extractor part is frozen. The final fully connected layers and softmax are trained on a medium dataset. A note indicates that more data allows for retraining more of the network (or all of it).</li> </ol>	
<p>Best practice</p>	<p>Freeze reused components and train new layers (classification part). Many pre-trained models for image recognition available in Keras.</p>	
<p>Code</p>	<pre>from keras.applications.mobilenet_v2 import MobileNetV2 model = MobilNetV2(weights='imagenet', include_top=False, input_shape=(...))</pre>	
<p>bottleneck layers</p>	<p>To reduce number of operations, 1x1 conv-layer are used.</p>	
<p>Bayes Law</p>	$P(C_k x) = \frac{p(x C_k)P(C_k)}{p(x)}$	<p>Select as winning category the one having the largest a posteriori probability. Doing so we guaranteed to maximise the accuracy.</p> <p><math>P(C_k x)</math>: a posteriori probability, probability of class j given observation x</p> <p><math>p(x C_k)</math>: likelihood, probability of observing x given class j</p> <p><math>P(C_k)</math>: a priori probability, probability of class j</p> <p><math>p(x)</math>: evidence = probability of x, ...unconditional to any class...</p>
<p><b>Auto-encoders</b></p>	<p>An autoencoder is a neural network trained to reproduce its input and able to discover "structures" and "efficient coding's" of the input space. The simplest form is a feedforward, non-recurrent neural network.</p>	
<p>"diabolo" network</p>	<p>The diagram shows a feedforward neural network with an encoding phase and a decoding phase. The input <math>x</math> is processed by an encoder <math>E_{\theta_E}</math> to produce a latent representation <math>z</math>. This <math>z</math> is then processed by a decoder <math>D_{\theta_D}</math> to produce the reconstructed output <math>\hat{x}</math>. The overall mapping function is <math>\hat{x} = h_{\theta}(x)</math>.</p>	<p>The mapping function <math>h_{\theta}(x)</math> is trained to reconstruct its own inputs, instead of predicting a target value.</p> <p><b>Usage:</b></p> <ul style="list-style-type: none"> <li>Data compression / dimensionality reduction</li> <li>Use encoder to obtain features</li> <li>De-noising images </li> <li>Image in-painting </li> <li>Initialise deep network for supervised learning</li> </ul> <p><b>Loss function</b></p> $J(\theta) = J(\theta_E, \theta_D) = \frac{1}{2N} \sum_{n=1}^N (h_{\theta}(x_n) - x_n)^2 = \frac{1}{2N} \sum_{n=1}^N (\hat{x}_n - x_n)^2$

**11.+12. Recurrent Neural Networks (RNN)**

<b>RNN</b>	RNNs learn to memorise context information that may be important to draw conclusions on observations made later on. RNN helps wherever we need context from the previous input.			
<b>Applications</b>		Input sequence $x = (x_1 \dots x_{T_x})$	Output sequence $y = (y_1 \dots y_{T_y})$	
	<b>Speech recognition</b> Assistants (Siri, Alexa, ...) Video captioning Transcription to text		→	Today, the weather on Rigi above the fog is beautiful.
	<b>Sentiment classification</b> Classification of text Emojification	Today, the weather on Rigi above the fog is beautiful.	→	
	<b>Machine translation</b> DeepL, Google Translate	Today, the weather on Rigi above the fog is beautiful.	→	Heute ist das Wetter auf der Rigi oberhalb des Nebels schön.
	<b>Captioning, Subtitling</b> images, YouTube videos		→	Sun shining above the clouds.
	<b>Chatbots, Q/A</b> Siri, Alexa		→	task, answer
	<b>Named entity recognition (NER)</b> flag names in sentences	Today, the weather above the fog on Rigi is beautiful. 0 0 0 0 0 0 0 0 1 0 0		
	<b>Music generation</b> Magenta DeepJazz		→	
	<b>Word generation</b>		→	
<b>time sequence modelling, prediction</b>				
<b>Approaches</b>	<b>One-to-many</b>	<b>Many-to-one</b>	<b>many-to-many</b>	<b>many-to-many</b>
	image tagging, music generation	sentiment analysis	named entity recognition	language translation, speech recognition, chatbots
	$X \rightarrow YYY$	$XXX \rightarrow Y$	$XXX \rightarrow YYY$	$XX \rightarrow YYYYY$
	$T_x = 1$ $T_y \geq 1$	$T_x \geq 1$ $T_y = 1$	$T_x \geq 1, T_y \geq 1$ $T_x = T_y$	$T_x \geq 1, T_y \geq 1$ $T_x \neq T_y$
<b>Simple RNNs</b>  e.g. $T_x = T_y$ like for NER	un-rolled, un-folded		rolled, folded	
				
	<p><b>State</b> h is update through a succession of steps. The state <b>memorises</b> (part of) the "history". Init with 0.  <math>h_t = g(W_x * x_t + W_h * h_{t-1} + b_h) \rightarrow W</math> and <math>b</math> are the same for all time step</p>			
code	<pre>from keras.layers import SimpleRNN model = Sequential() model.add(SimpleRNN(units=..., input_shape=..., ...))</pre>			
Example	Predict gender / home country from first name			

<p><b>Generative RNNs</b></p>	<p>A generative system is a system able to generate new consistent data from a seed. By consistent, we mean respecting temporal or spatial "structure" that have been learned from the input space.</p>															
<p>example</p>	<p>SEED "Once upon a time" ↓ Text generation system ↓ OUTPUT "Once upon a time a princess kissed a frog." ...</p>	<p>domain name suggestion Shakespeare sonnets (poetic verse form) C code generation</p>														
<p>Approaches</p>	<p><b>Many to one</b></p>  <p>loss is computed from 1 output</p>	<p><b>Many to many</b></p>  <p>loss is computed from all the outputs</p>														
<p>The difference is mainly at training time where the RNN cells are emitting an output at each step.</p>																
<p>back propagation</p>	<p>run forward and backward through <b>chunks</b> of the sequence instead of whole sequence.</p>															
<p><b>Long Term Dependencies Problem</b></p>	<p>Simple RNN's fail if a wider context is needed and the gap between the words grows too large. Long-range dependencies are hard to learn due to vanishing and exploding gradients problem.</p>															
<p>Suitable Initialisation of Weights</p>	<p><b>MLPs/ CNNs:</b></p> <ul style="list-style-type: none"> <li>- Properly initialise weights (Xavier/He) by using random numbers (uniform or normal) with mean=0 and suitably scaled stdev.</li> <li>- Batch normalisation</li> </ul> <p><b>RNNs:</b></p> <ul style="list-style-type: none"> <li>with 'vanilla' recurrent units</li> <li>- IRNN: identity matrix - 67.0 % accuracy - highly sensitive to parameters</li> <li>- np-RNN: normalized-positive definite matrix - 75.2% accuracy - low sensitive to parameters</li> </ul> <p>both</p> <ul style="list-style-type: none"> <li>- Use non-saturating activation functions (ReLU or LeakyReLU) to alleviate (=mindern) the vanishing gradients problem</li> <li>- Clipping gradients</li> </ul>															
<p>Long-Short-Term Memory (LSTM) Cells or Gated Recurrent Units (GRU)</p>	<p>- LSTM 78.5% accuracy - low sensitive to parameters</p> <p>LSTM network have the same general structure of cyclically updated cells. It is designed to keep a <b>long-term memory</b> that is kept as additional state variable <math>c_t</math> to be updated.</p> <p>LSTM Cell</p>  <p>long-term: <math>c_t = f_t * c_{t-1} + I - T * \tilde{c}_t</math> short-term: <math>h_t = o_t * \tanh(c_t)</math></p>	<p>Long-term memory is updated through <b>Gates</b> (marked with <math>\otimes</math>).</p> <p>Gates control how the information flows between short-term state, input and long-term state</p> <ol style="list-style-type: none"> <li>1: Forget Gate: <math>f_t = \sigma(W_{f,x} * x_t + W_{f,h} * h_{t-1} + b_f)</math></li> <li>2: Input Gate: <math>i_t = \sigma(W_{i,x} * x_t + W_{i,h} * h_{t-1} + b_i)</math> <math>\tilde{c}_t = \tanh(W_{c,x} * x_t + W_{c,h} * h_{t-1} + b_c)</math></li> <li>3: Output Gate: <math>o_t = \sigma(W_{o,x} * x_t + W_{o,h} * h_{t-1} + b_o)</math></li> </ol> <p>The gates are implemented by</p> <ul style="list-style-type: none"> <li>- affine transformation of inputs</li> <li>- activation function</li> <li>- element-wise multiplication</li> </ul>														
<p>backprop: "Super-Highway for Backprop": </p>																
<p>Word Embedding</p>	<p>Gated Recurrent Unit (GRU)</p>  <p>most popular variation of LSTM. No separate long-term memory forget gate and input gate merged, new relevance gate</p> <p>Relevance Gate: <math>r_t = \sigma(W_{rh} * h_{t-1} + W_{rx} * x_t + b_r)</math> Candidate State: <math>\tilde{h}_t = \tanh(W_{ch} * r_t * h_{t-1} + W_{cx} * x_t + b_c)</math> Update Gate: <math>u_t = \sigma(W_{uh} * h_{t-1} + W_{ux} * x_t + b_u)</math> Update: <math>h_t = (1 - u_t) * h_{t-1} + u_t * \tilde{h}_t</math> less parameters, since only 3 weight matrices (<math>n_h \times (n_h + n_x)</math>)</p>	<table border="1"> <tr> <td>apple</td> <td>juice</td> </tr> <tr> <td>orange</td> <td></td> </tr> <tr> <td>banana</td> <td>rice</td> </tr> <tr> <td></td> <td>milk</td> </tr> <tr> <td></td> <td>bus</td> </tr> <tr> <td></td> <td>train</td> </tr> <tr> <td></td> <td>car</td> </tr> </table>	apple	juice	orange		banana	rice		milk		bus		train		car
apple	juice															
orange																
banana	rice															
	milk															
	bus															
	train															
	car															
<p>Word embedding is a projection of a word into vectors of real numbers, i.e. a mapping of a space where each word has its own dimension to a space that is of lower dimensionality.</p> <p><b>word2vec</b> relates to models producing word embedding into space with contextual similarity or words, i.e. words that share a common context are near. Bag of Words eat-apple, driving-car</p> <p>Continuous Bag of Word: Eat an ... every day (apple) -&gt; see google search</p>																