

Neural network architectures (part one: deep NNs & CNNs)

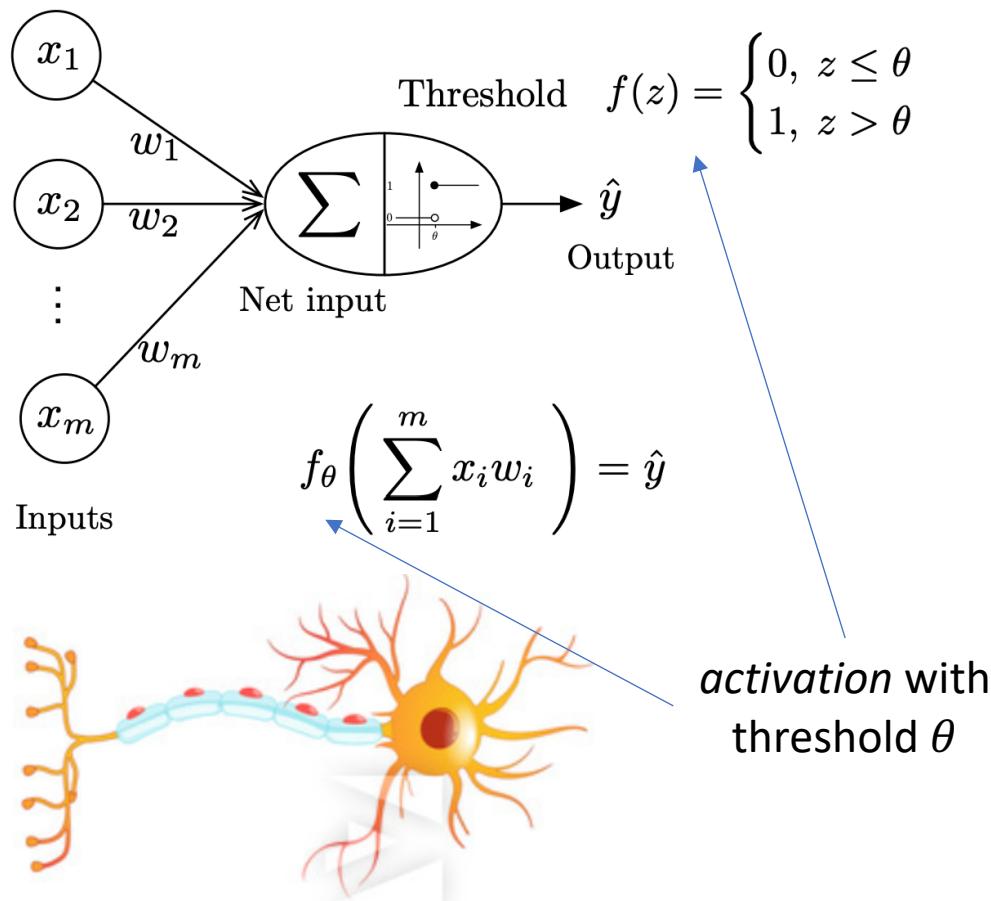
Materials for this Section:

- I. Goodfellow , Y. Bengio, and A. Courville,
Deep Learning, Chapters 6-10:
<https://www.deeplearningbook.org>
- <https://playground.tensorflow.org>

Deep neural networks

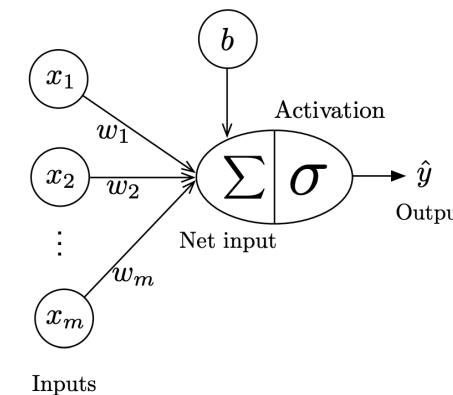
simplest neural network, *perceptron*

- a computational model of a biological neuron
- single input *layer* with multiple input *units* and single output



activation with threshold θ

- It is convenient to include threshold as additional *bias unit*



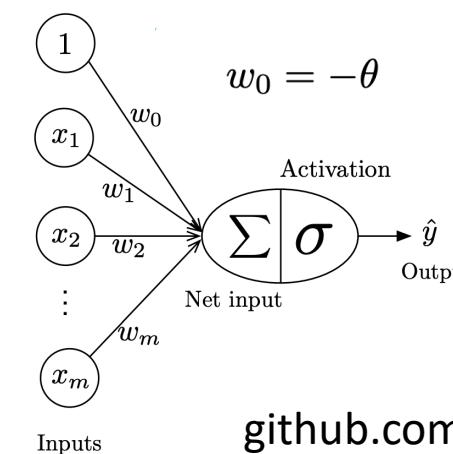
$$\hat{y} = \begin{cases} 0, & z - \theta \leq 0 \\ 1, & z - \theta > 0 \end{cases}$$

"separate" bias unit

$$\sigma\left(\sum_{i=1}^m x_i w_i + b\right) = \sigma(\mathbf{x}^T \mathbf{w} + b) = \hat{y}$$

$$\sigma(z) = \begin{cases} 0, & z \leq 0 \\ 1, & z > 0 \end{cases}$$

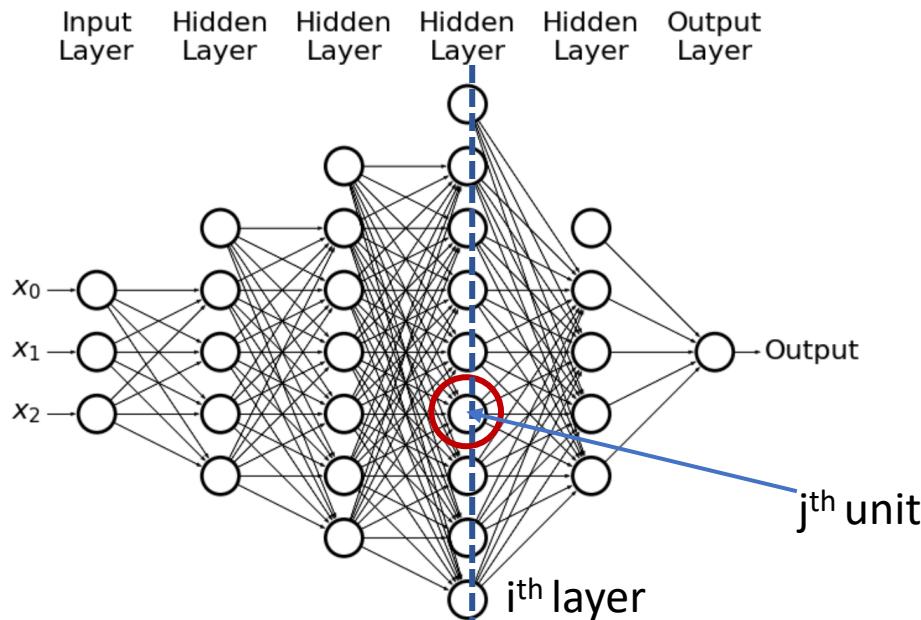
$$b = -\theta$$



- or additional weight

Deep neural networks

- perceptron can not represent XOR function
-- to do so we need additional *hidden* layers

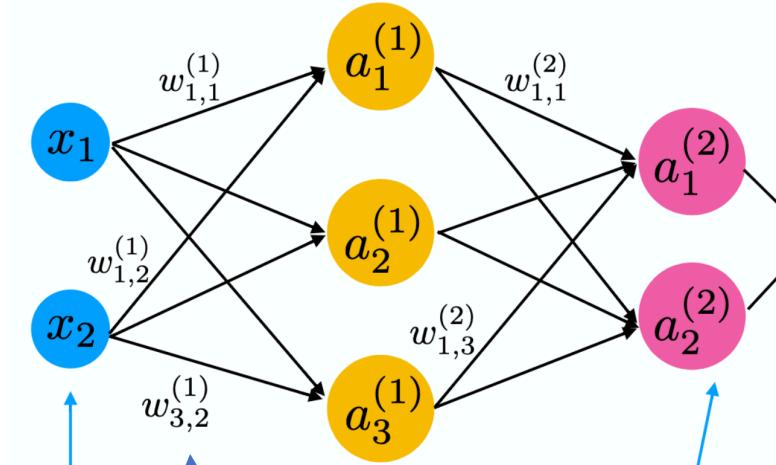


- fully-connected multi-layer perceptron

By noting i the i^{th} layer of the network and j the j^{th} hidden unit of the layer, we have:

$$x_j^{(i)} = f\left(\sum_k w_k^{(i)} x_k^{(i-1)}\right)$$

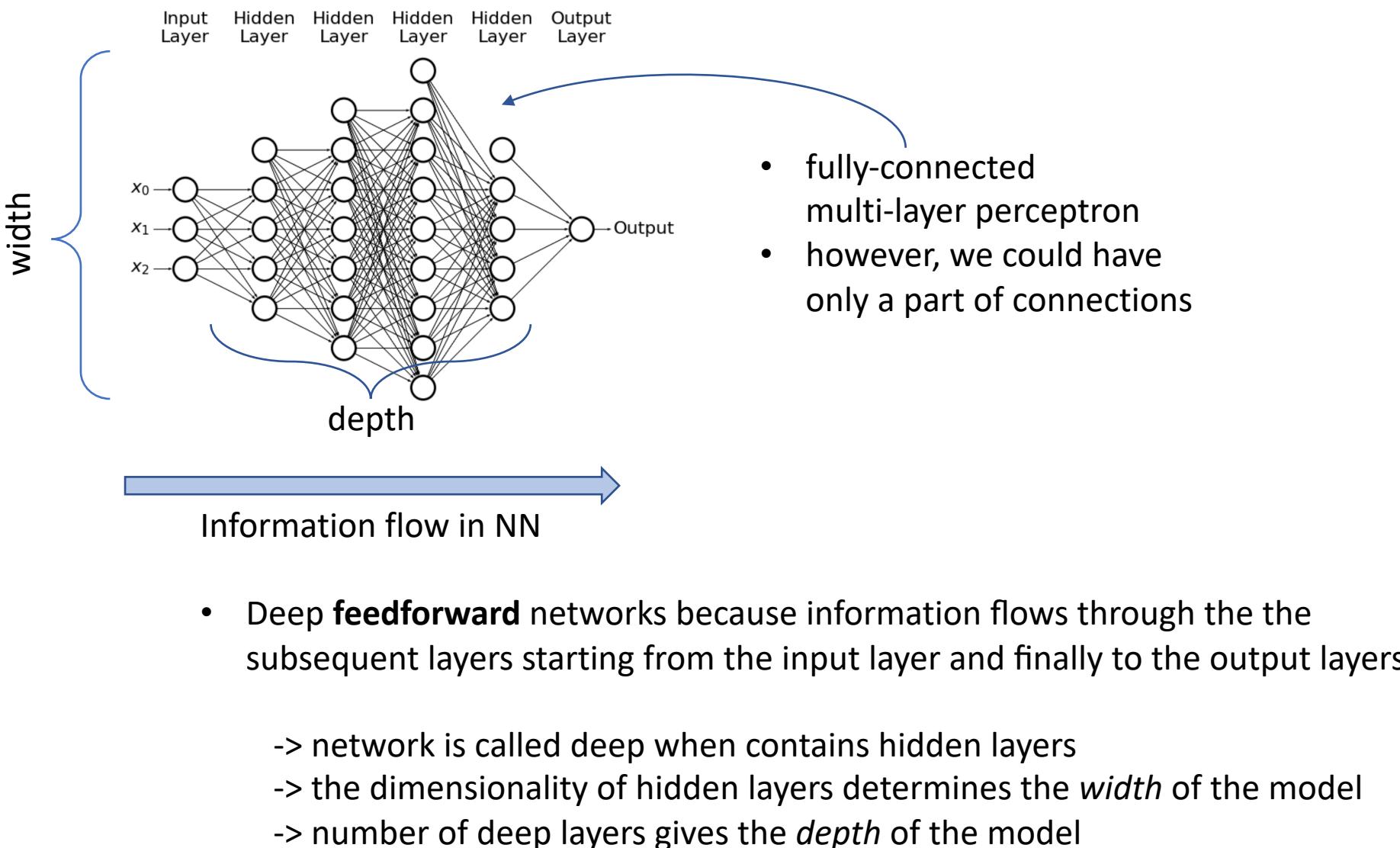
where we note w , b , z the weight, bias and output respectively.



- each connection is parametrized by its weight
- each neuron has its own activation function

- network is called **deep** when it contains hidden layers
- all weights are fitted during *training*

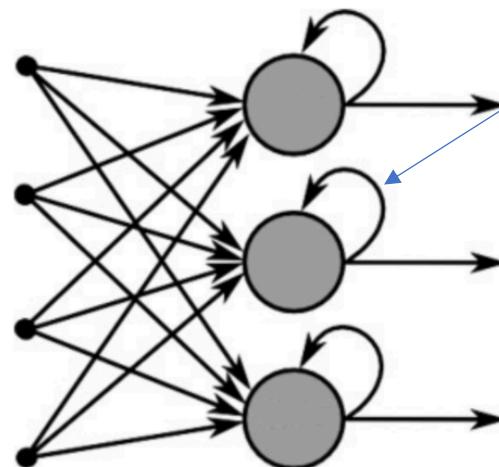
Deep neural networks



Deep neural networks

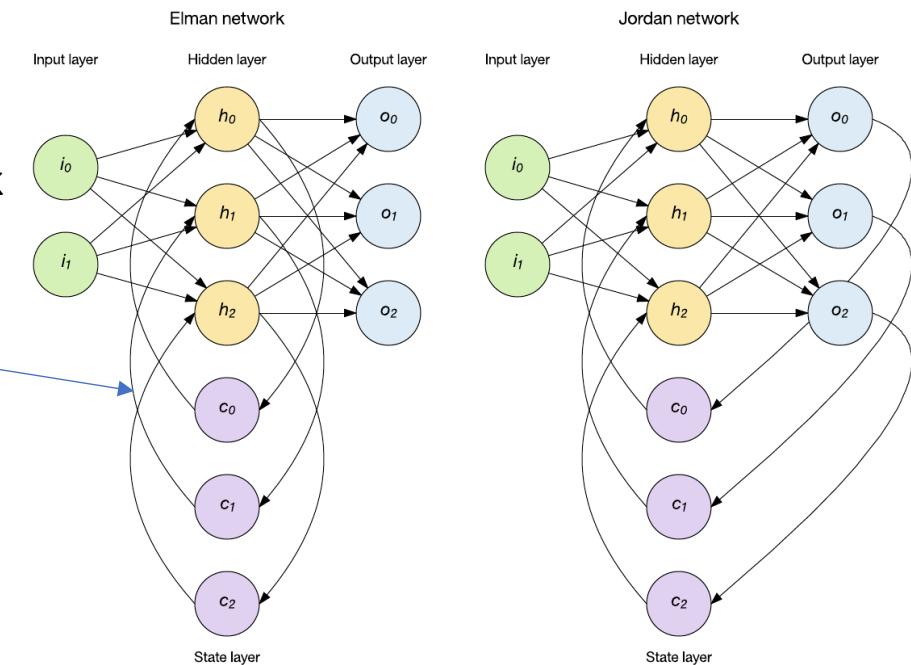
Information flow in NN

- When feedforward neural networks are extended to include feedback connections, they are called **recurrent** neural networks

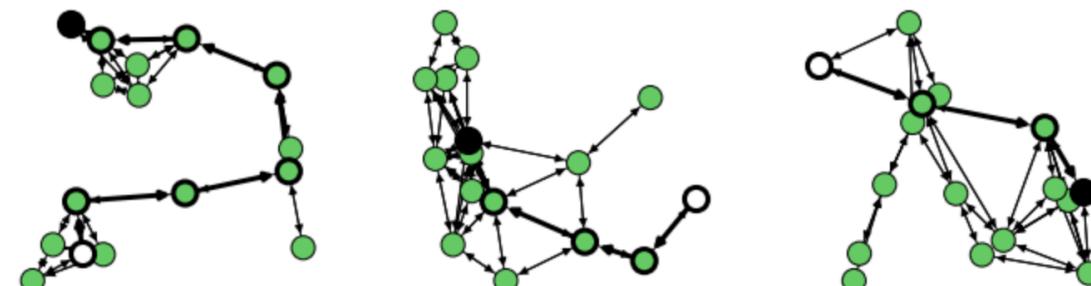


Recurrent Neural Network

feedback connections enables
information flow in reverse
direction



- graph** neural networks
-> no distinguished direction of information flow through the network

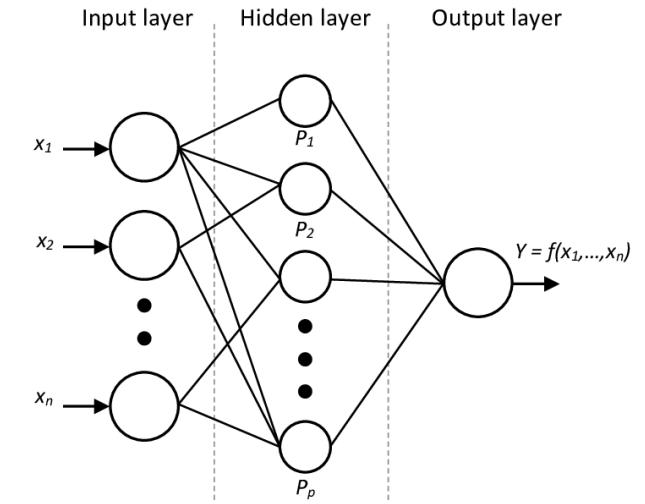


Deep neural networks: universality

Universal Approximation Theorem

- A feedforward network with a linear output layer and at least **one hidden layer** with any nonlinear bounded activation function (such as logistic activation) can **approximate any** continuous **function** (on a closed and bounded subset of \mathbb{R}^n), from one finite-dimensional space to another, with **any** desired **precision**
- such network should have enough hidden units (enough width)
- however the layer may be infeasibly large and may fail to generalize correctly
- using **deeper** models can reduce the number of units required to represent the desired, function and reduce the amount of generalization error

any $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$



In a nutshell:

- a feedforward network with a single hidden layer is sufficient to represent any function, but may overfit,
- using deeper models can reduce the required number of hidden units and reduce the generalization error.

Deep neural networks: training

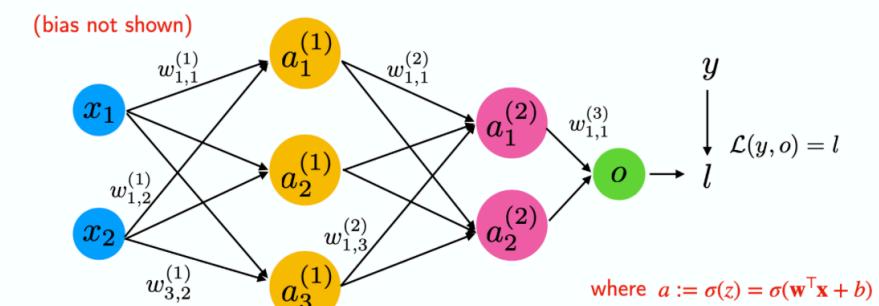
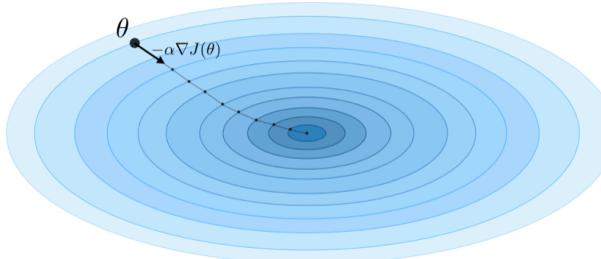
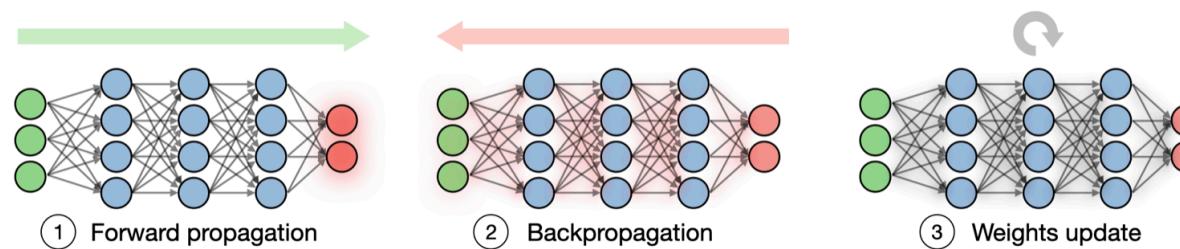
Training NN is just optimization problem:

- The task is to minimize loss function in parameters (weights) space

$$w \leftarrow w - \alpha \frac{\partial L(z, y)}{\partial w}$$

□ **Updating weights** — In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data and perform forward propagation to compute the loss.
- Step 2: Backpropagate the loss to get the gradient of the loss with respect to each weight.
- Step 3: Use the gradients to update the weights of the network.



$$\begin{aligned}\frac{\partial l}{\partial w_{1,1}^{(1)}} &= \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}} \\ &\quad + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}^{(1)}}\end{aligned}$$

(Assume network for binary classification)

backpropagation utilizes *chain rule*, and reuses partial derivatives calculated before.

good review: brilliant.org/wiki/backpropagation/

Deep neural networks: training

Tutorials > Deep Learning with PyTorch: A 60 Minute Blitz

DEEP LEARNING WITH PYTORCH: A 60 MINUTE BLITZ ↗

Author: Soumith Chintala

Goal of this tutorial:

- Understand PyTorch's Tensor library and neural networks at a high level.
- Train a small neural network to classify images

This tutorial assumes that you have a basic familiarity of numpy

* NOTE
Make sure you have the `torch` and `torchvision` packages installed.

What is PyTorch? Autograd: Automatic Differentiation Neural Networks

DEEP LEARNING WITH PYTORCH: A 60 MINUTE BLITZ ↗

https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html

`torch.autograd` is PyTorch's **automatic differentiation** engine that powers neural network training

Generally speaking, `torch.autograd` is an engine for computing vector-Jacobian product. That is, given any vector $v = (v_1 \ v_2 \ \dots \ v_m)^T$, compute the product $v^T \cdot J$. If v happens to be the gradient of a scalar function $l = g(y)$, that is, $v = (\frac{\partial l}{\partial y_1} \ \dots \ \frac{\partial l}{\partial y_m})^T$, then by the chain rule, the vector-Jacobian product would be the gradient of l with respect to x :

$$J^T \cdot v = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \dots & \frac{\partial y_m}{\partial x_n} \end{pmatrix} \begin{pmatrix} \frac{\partial l}{\partial y_1} \\ \vdots \\ \frac{\partial l}{\partial y_m} \end{pmatrix} = \begin{pmatrix} \frac{\partial l}{\partial x_1} \\ \vdots \\ \frac{\partial l}{\partial x_n} \end{pmatrix}$$

Text source: https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html#sphx-glr-beginner-blitz-autograd-tutorial-py

Deep neural networks: training

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$

Gradient Descent optimization

- In the **gradient descent** (GD) method weights are updated after a full loop over training data:

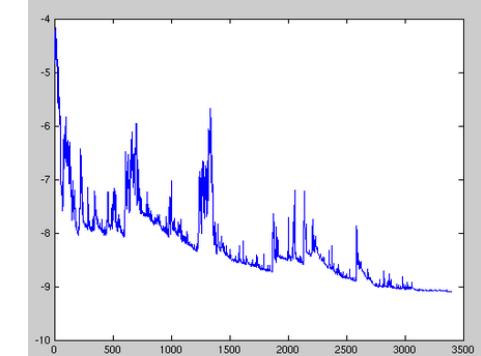
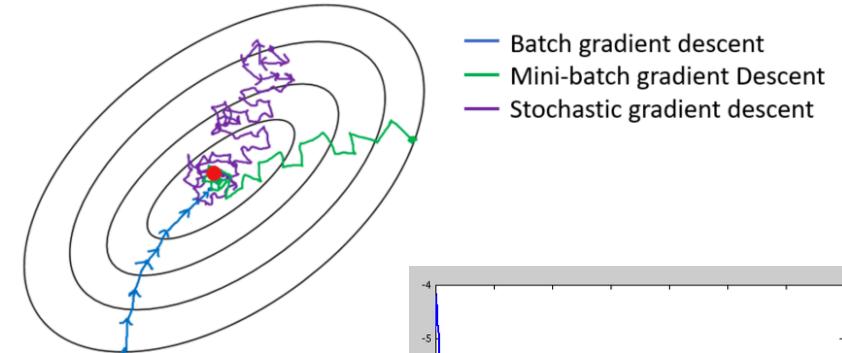
$$W \rightarrow W - \alpha \cdot \nabla_W L(W)$$

- In the **stochastic gradient descent** (SGD) method weights are updated for each training sample:

$$W \rightarrow W - \alpha \cdot \nabla_W L(W; x^{(i)}; y^{(i)})$$

- The common method, which is somehow between GD and SGD, is **mini-batch gradient descent** (MBGD) - the one we used for MNIST example:

$$W \rightarrow W - \alpha \cdot \nabla_W L(W; x^{(i:i+n)}; y^{(i:i+n)})$$



SGD:

- fluctuations may help to escape local minima
- problem with overshooting exact minimum
- decreasing** learning rate during optimization enables to converge

Deep neural networks: training

CLASS `torch.optim.SGD(params, lr=<required parameter>, momentum=0, dampening=0, weight_decay=0, nesterov=False)`

[SOURCE]

Implements stochastic gradient descent (optionally with momentum).

Adaptive optimization

SGD: constant learning rate, which needs to be tuned by hand

Adaptive models modifies the learning rate α using information about gradients:

- Adagrad:
- Adadelta
- Adam

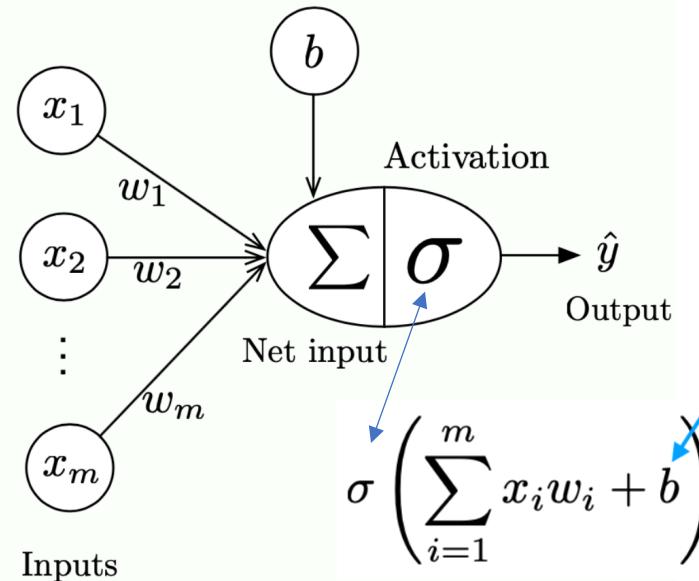
$$W_{t+1,i} = W_{t,i} - \frac{\alpha}{\sqrt{G_{t,ii} + \epsilon}} \nabla_W L(W_{t,i})$$

Another technique designed to accelerate learning: *momentum*

$$\begin{aligned} v_t &= \gamma \cdot v_{t-1} + \alpha \cdot \nabla_W L(W; x^{(i)}; y^{(i)}) \\ W &\rightarrow W - v_t \end{aligned}$$

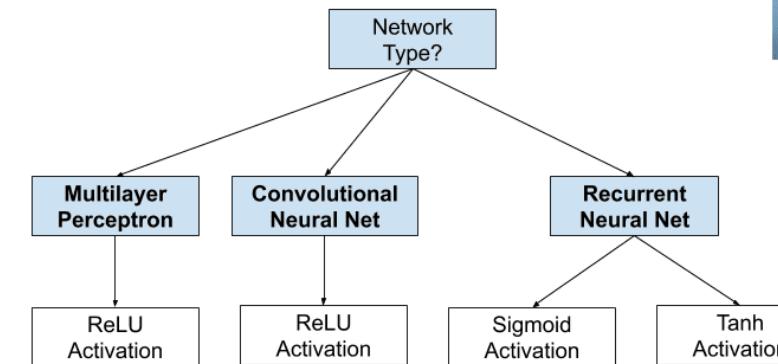
The SGD optimizer in PyTorch is just gradient descent. User decide how to pass data: single sample mini-batch or batch.

Deep neural networks: activation

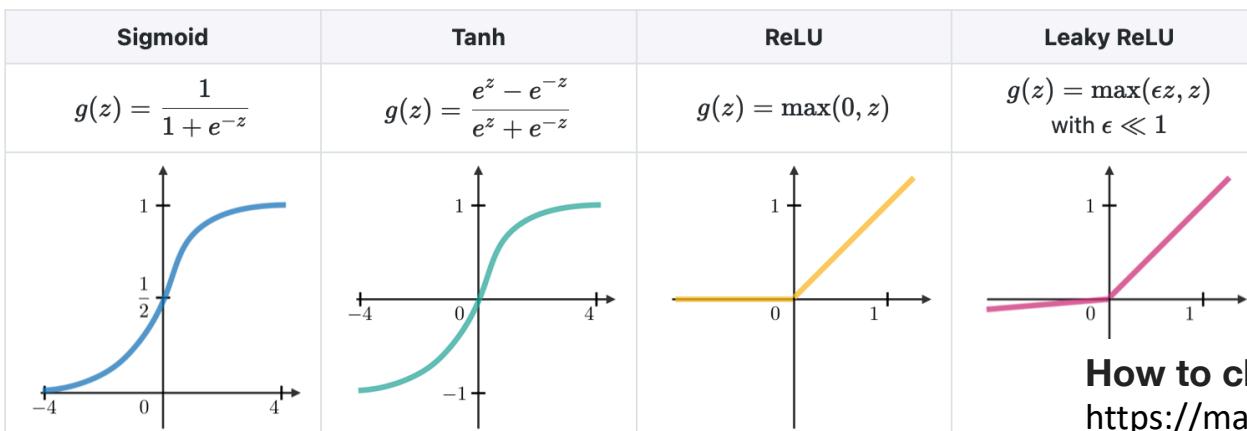


We usually chose activation for hidden units and output layer separately:

How to Choose an Hidden Layer Activation Function



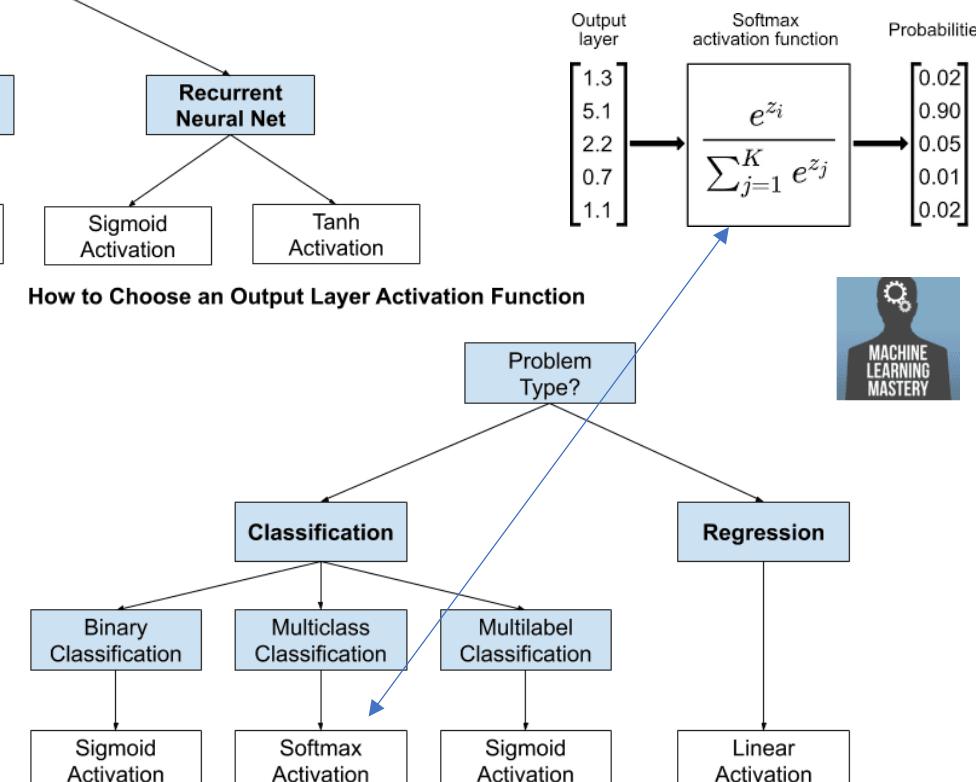
□ **Activation function** — Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:



How to chose Activation for Hidden and output Layers:

<https://machinelearningmastery.com/choose-an-activation-function-for-deep-learning/>

How to Choose an Output Layer Activation Function



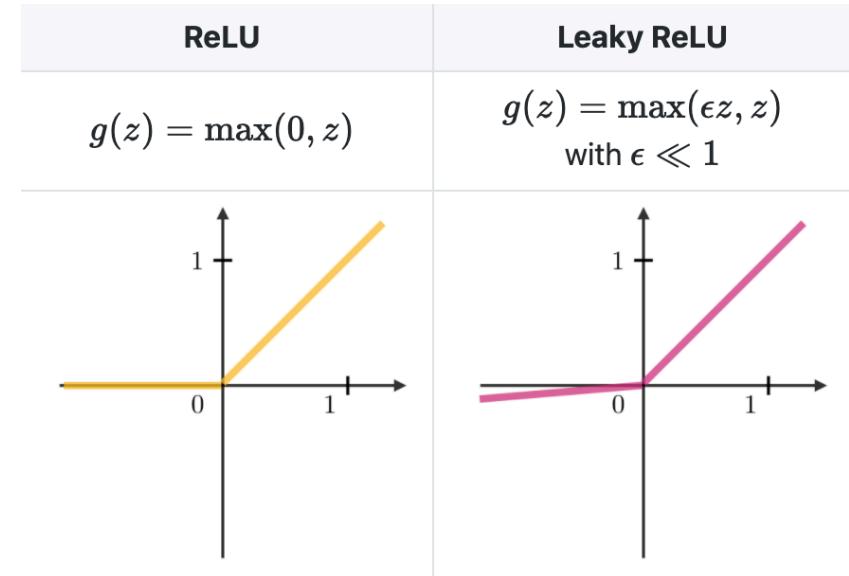
Deep neural networks: activation

ReLU is a default choice for activation of hidden units:

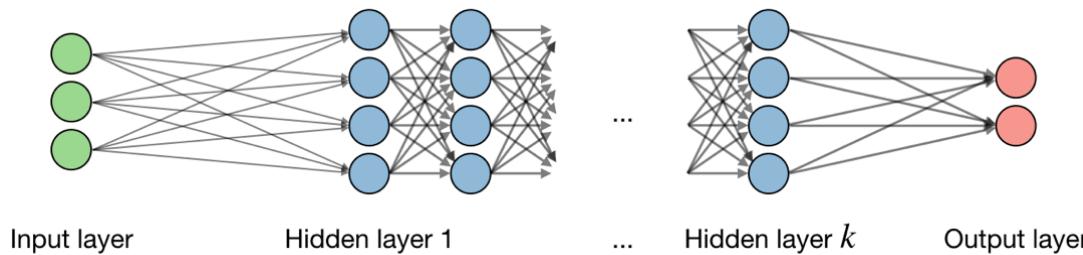
- easy to optimize because it is similar to linear unit
- reduces *vanishing gradient* problem
 - backpropagation computes gradients by the chain rule
 - to compute gradients of the early layers we multiply small numbers many *times*
 - loss signal decreases exponentially when propagating backwards therefore early layers train very slowly

some drawbacks:

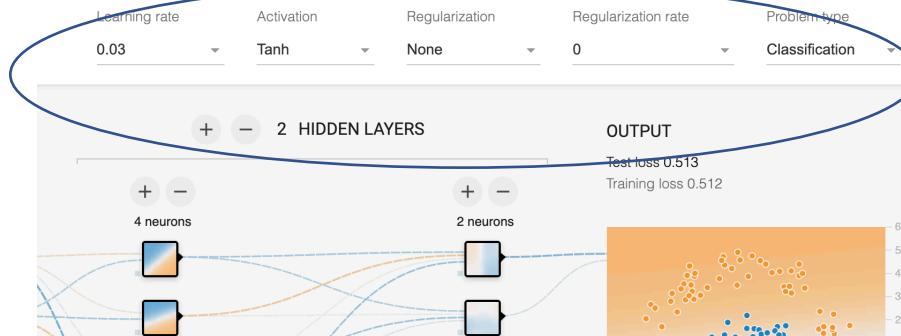
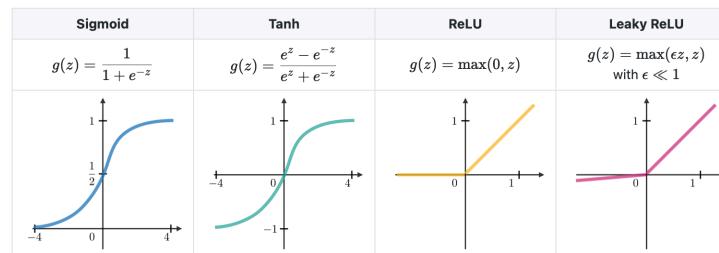
- non-differentiable at zero
- *dying ReLU* problem
 - some neurons essentially *die* for all inputs and remain *inactive* no matter what input is supplied
 - corrected in *leaky ReLU*



Deep neural networks: hyperparameters



Activation function — Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:



Note, that NN has much more hyperparameters than other ML models:

- number of hidden layers
- number of neurons in layer
- activation functions
- learning algorithm settings: SGD variation, learning rate, number of epochs, momentum, etc.
- batch size in mini-batch gradient descent
- regularization method and its parameters

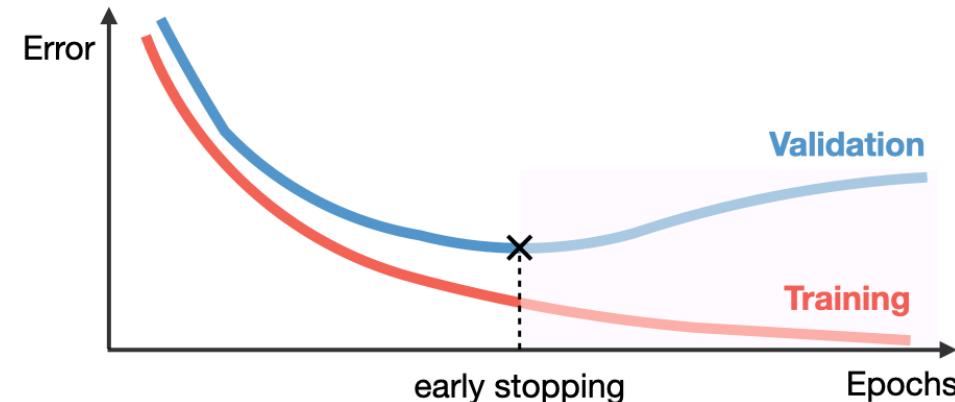
It is crucial to use cross-validation (at least single validation subset) for hyperparameter tuning.

- Let's play with hyperparameters: <https://playground.tensorflow.org>

Deep neural networks: regularization

A central problem in machine learning is how to make an algorithm that generalize well

- regularization is any strategy designed to prevent overfitting
- deep networks uses three main types of regularizations:
 - loss function modification
 - dropout technique
 - early stopping:



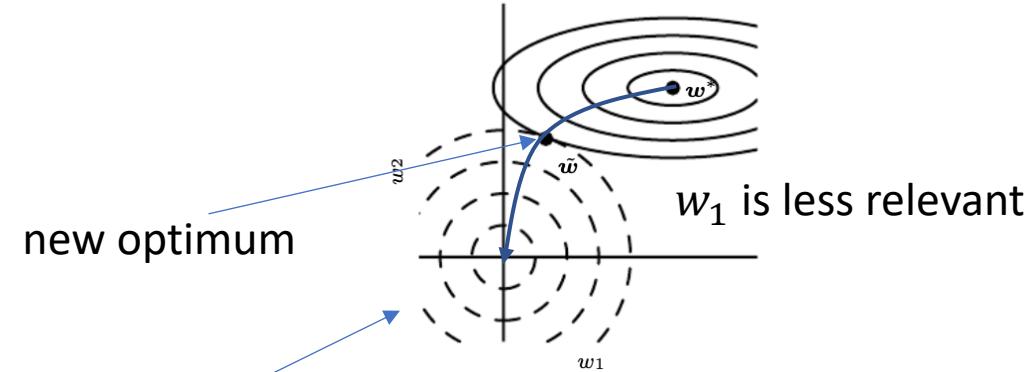
Deep neural networks: regularization

Loss function modification:

we add **regularization term (regularizer)** to the loss function $L(\mathbf{w})$:

$$L(\mathbf{w}, \mathbf{X}, \mathbf{y}) \rightarrow L(\mathbf{w}, \mathbf{X}, \mathbf{y}) + \lambda R(\mathbf{w})$$

- **ridge regression**, uses L_2 norm: $R(\mathbf{w}) = \sum_i w_i^2$
 - we put some constraints on the magnitude of weights
 - regularizer suppress stronger less relevant weights
-> modify loss function but not mess it up.
- **LASSO** uses L_1 : $R(\mathbf{w}) = \sum_i |w_i|$
 - The L_1 penalty causes a subset of the weights to become zero (suppressing is stronger than L_2),
 - can be used to features selection



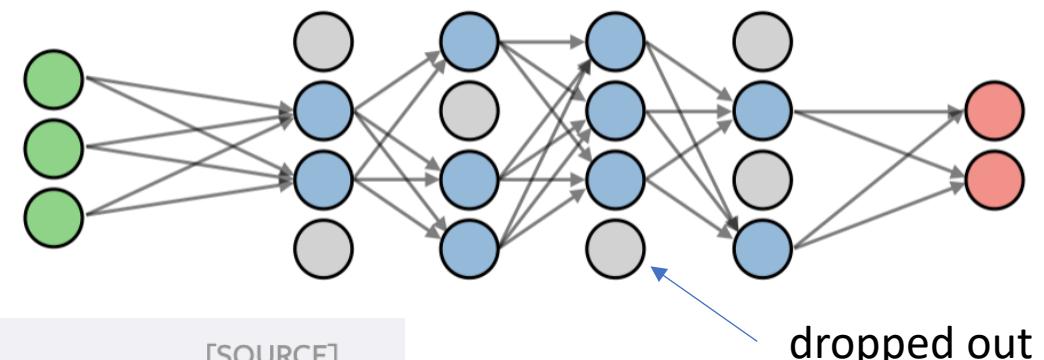
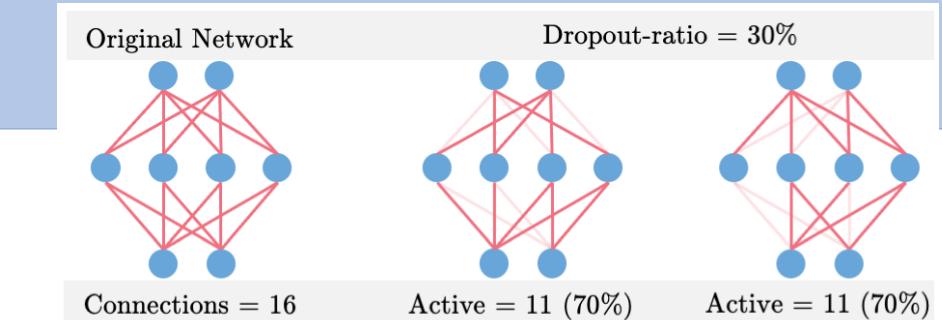
LASSO	Ridge
<ul style="list-style-type: none">• Shrinks coefficients to 0• Good for variable selection	Makes coefficients smaller
$\dots + \lambda \ \theta\ _1$ $\lambda \in \mathbb{R}$	$\dots + \lambda \ \theta\ _2^2$ $\lambda \in \mathbb{R}$

Deep neural networks: regularization

Dropout

- dropout is similar to ensemble learning:
-> involves training multiple network models on each test example
- This improves generalization -- force different neurons to learn the same *concept*

CLASS `torch.nn.Dropout(p=0.5, inplace=False)`



[SOURCE]

During training, randomly zeroes some of the elements of the input tensor with probability `p` using samples from a Bernoulli distribution. Each channel will be zeroed out independently on every forward call.

This has proven to be an effective technique for regularization and preventing the co-adaptation of neurons as described in the paper [Improving neural networks by preventing co-adaptation of feature detectors](#).

Furthermore, the outputs are scaled by a factor of $\frac{1}{1-p}$ during training. This means that during evaluation the module simply computes an identity function.

- dropout deactivates random set of neurons during training phase (usually about 50%)
- usually, dropout is applied to fully connected layers
- during inference stage dropout is turned off

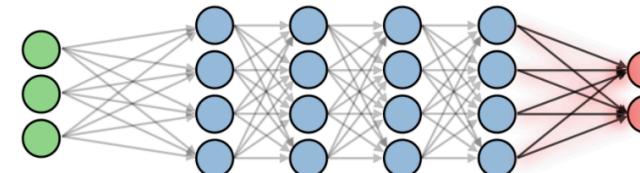
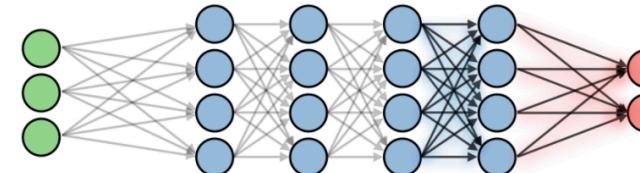
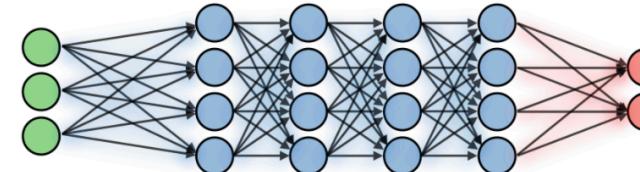
Deep neural networks: transfer learning

Transfer Learning

- It is useful to use initial weights pretrained on similar but much larger dataset
 - > it can dramatically shorten training
- it may be also useful to train only (last) part of a network
 - > especially when we have less data

Idea:
to train VGG (very deep convolutional NN)
from scratch we need weeks

-> to re-train pre-trained VGG to fit *our*
data we may need only hours or just days

Training size	Illustration	Explanation
Small		Freezes all layers, trains weights on softmax
Medium		Freezes most layers, trains weights on last layers and softmax
Large		Trains weights on layers and softmax by initializing weights on pre-trained ones

<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-deep-learning-tips-and-tricks>