


“A beginning is the time for taking the most delicate care that the balances are correct.”

Frank Herbert, *Dune*



CMPUT 628
Deep RL

Marlos C. Machado

Class 3/ 25

Plan

Overview / Refresher of (everything?) Deep Learning

Warning! *This will be quick. It is meant to start establishing a common language between us, but it is too fast for you if you are seeing this for the first time.*

Reminder: You can still leave

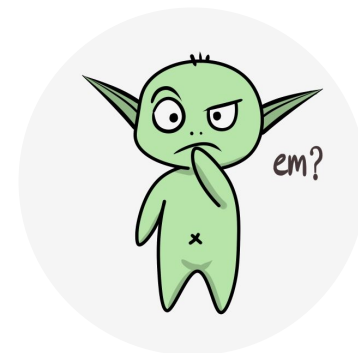
- I know, I know, *Deep Reinforcement Learning* sounds fun, modern, and hyp-ey

But...

- But this course won't be so well-structured as you (or I) would hope
- I won't teach you how to code fancy deep RL algorithms
- I'm not as much fun as you might think
- I don't care about grades – I might have a reputation :-)
 - *There won't be a practice midterm*
- I don't care if this course ends up being difficult



Please, interrupt me at any time!



MNIST and LeNet [Lecun et al.,1998]

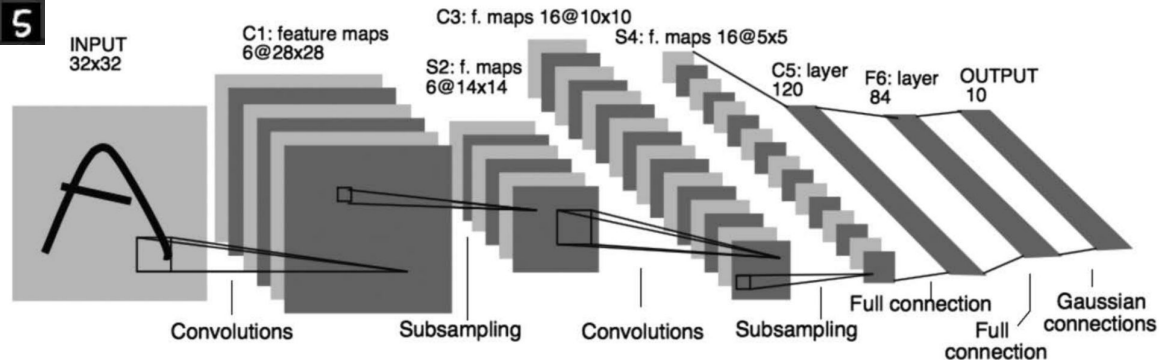
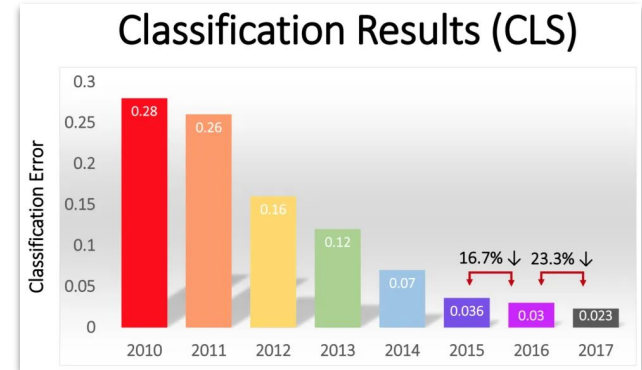
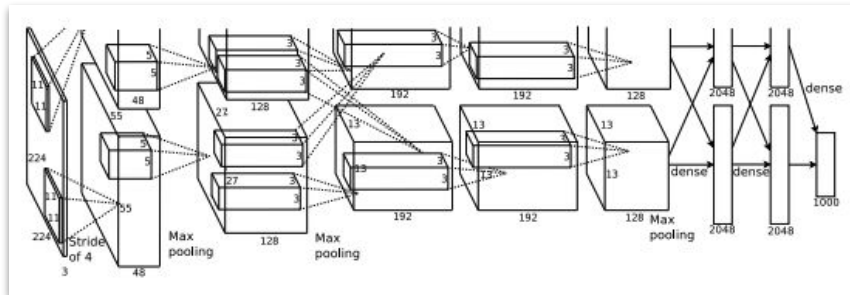


Figure 9.15: Deep Convolutional Network. Republished with permission of Proceedings of the IEEE, from Gradient-based learning applied to document recognition, LeCun, Bottou, Bengio, and Haffner, volume 86, 1998; permission conveyed through Copyright Clearance Center, Inc.

ImageNet



<https://medium.com/@prudhvi.gvr/imagenet-challenge-advancement-in-deep-learning-and-computer-vision-124fd33cb948>

Krizhevsky et al. (2012)

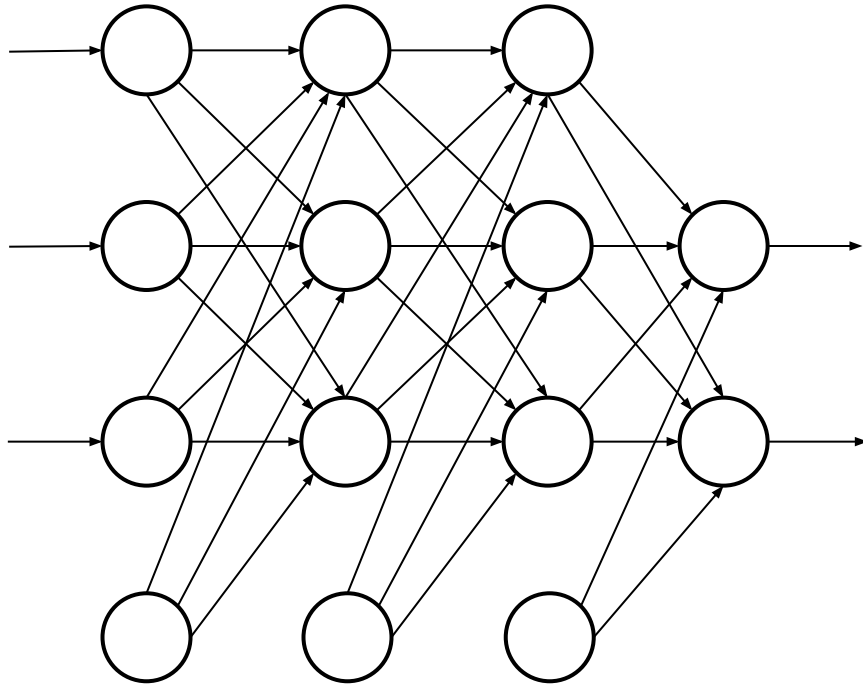
Nonlinear Function Approximation: Artificial Neural Networks

- The basics of deep reinforcement learning.
- Idea: Instead of using linear features, we feed the “raw” input to a neural network and ask it to predict the state (or state-action) value function.

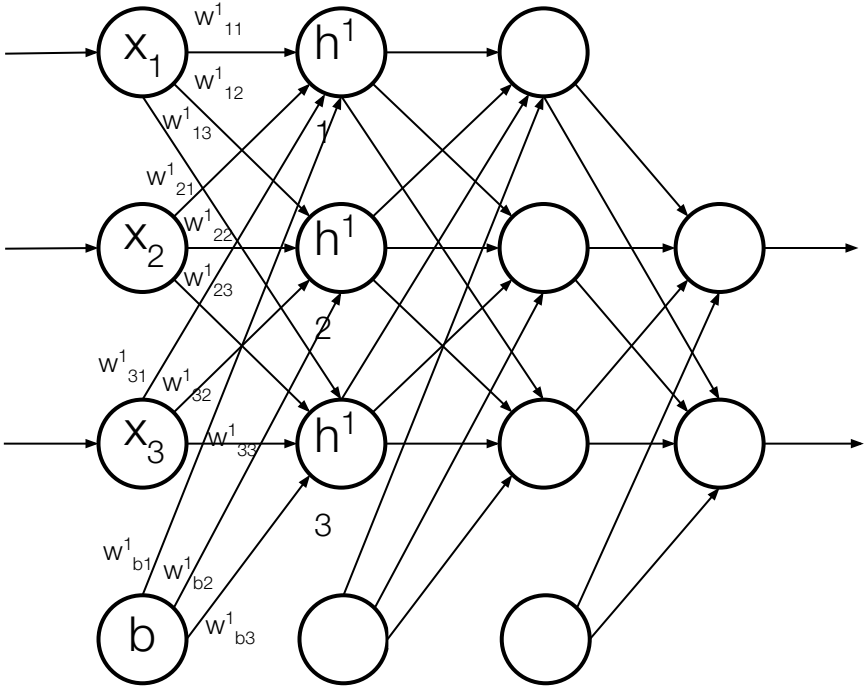


Neural Networks

Neural Networks



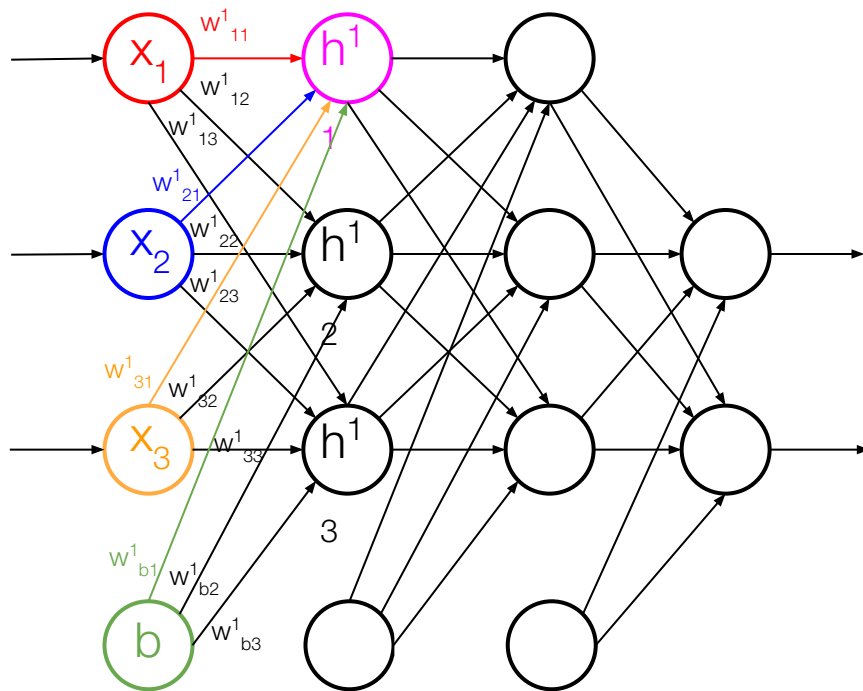
Neural Networks



$$\mathbf{h}^1 = \text{act}(\mathbf{x}\mathbf{W}^1)$$

$$\text{s.t. } h^1_1 = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

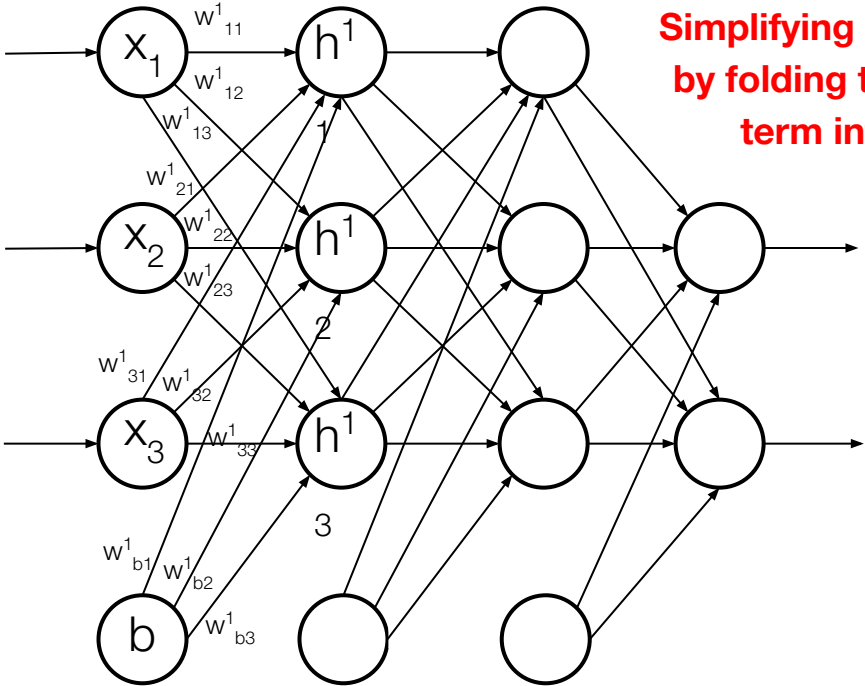
Neural Networks



$$\mathbf{h}^1 = \text{act}(\mathbf{x}\mathbf{W}^1)$$

$$\text{s.t. } h^1_1 = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

Neural Networks



**Simplifying notation
by folding the bias
term into x**

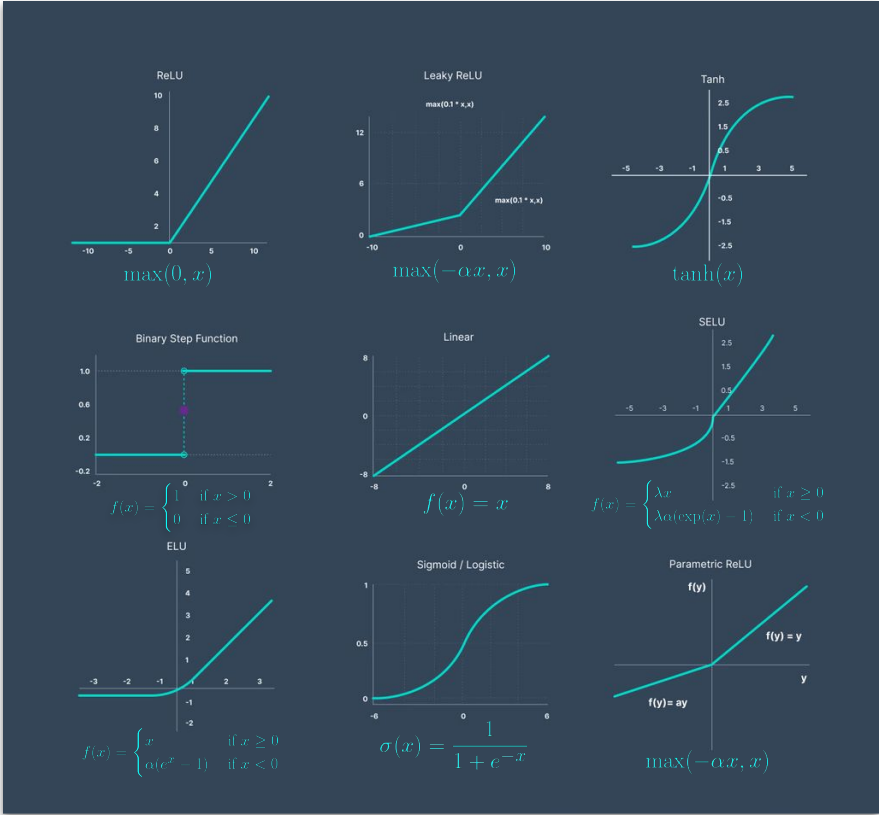
$$h^1 = \text{act}(xW^1)$$

$$\text{s.t. } h^1_{} = x_1w^1_{11} + x_2w^1_{21} + x_3w^1_{31} + bw^1_{b1}$$

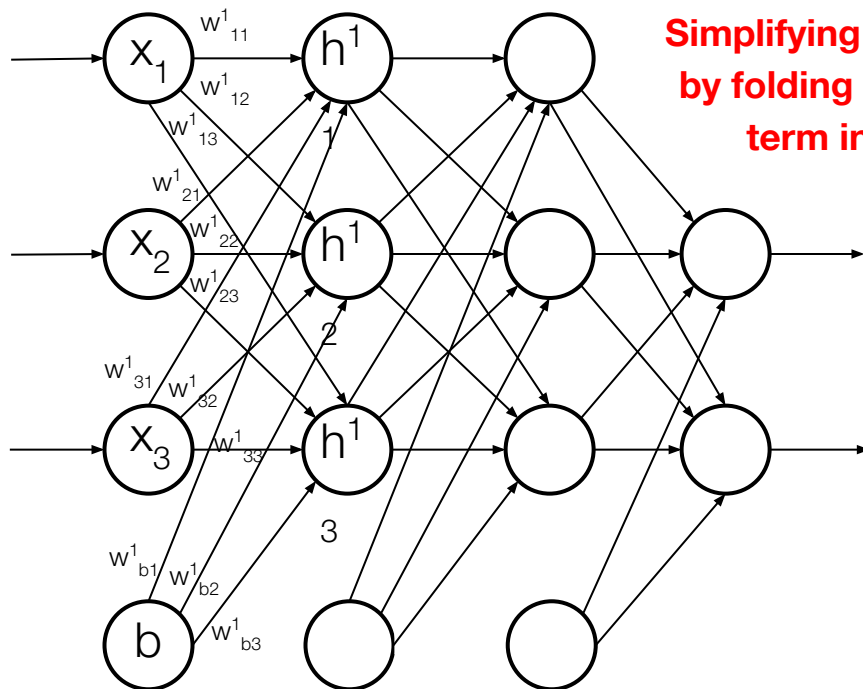
**The activation function
introduces non-linearity**

E.g.: $f(x) = \max(0, x)$

Main Types of Activation Functions



Neural Networks



**Simplifying notation
by folding the bias
term into x**

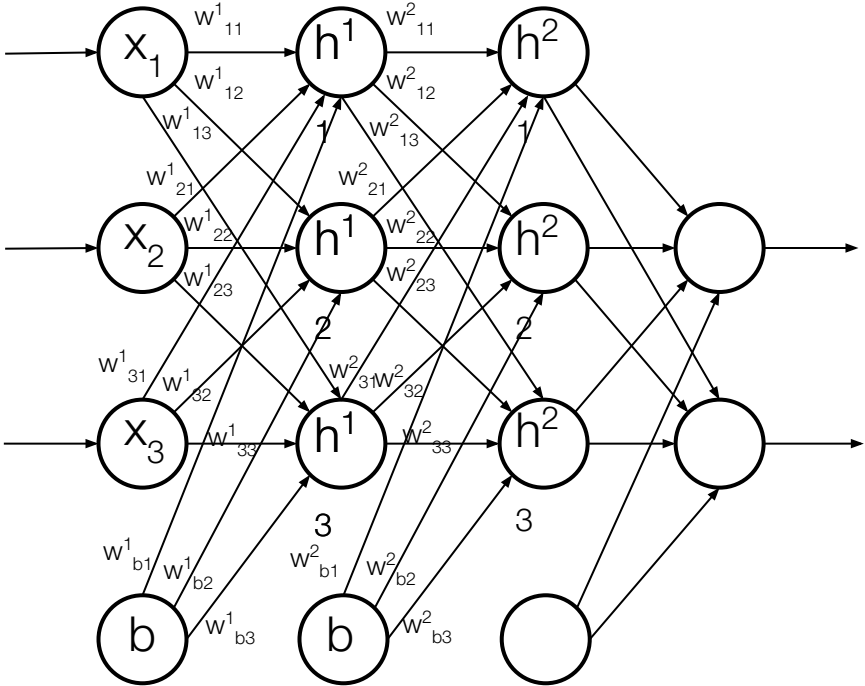
$$h^1 = \text{act}(xW^1)$$

$$\text{s.t. } h^1_{\cdot} = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

**The activation function
introduces non-linearity**

E.g.: $f(x) = \max(0, x)$

Neural Networks



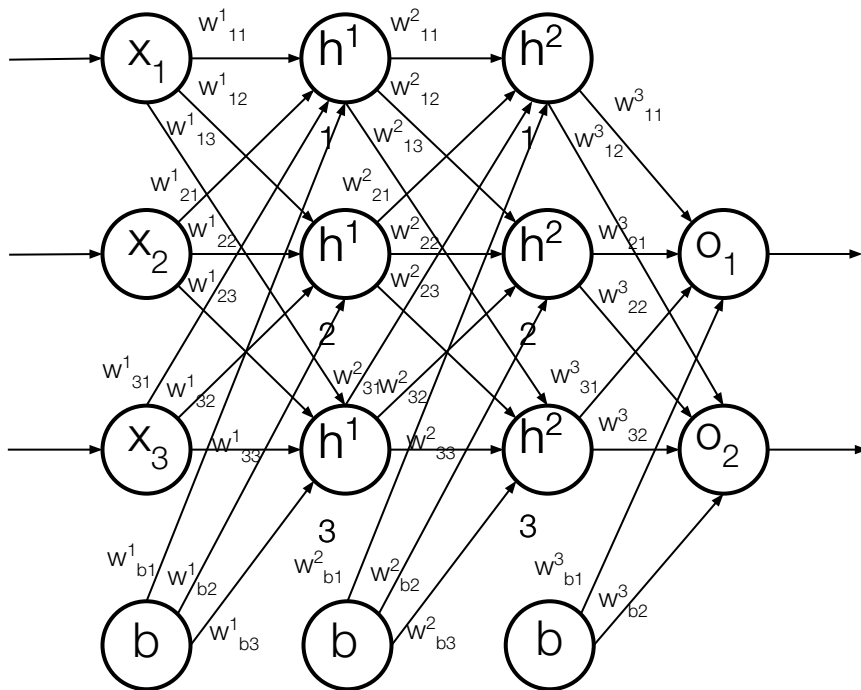
$$\mathbf{h}^1 = \text{act}(\mathbf{x}\mathbf{W}^1)$$

$$\text{s.t. } h^1_1 = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

$$\mathbf{h}^2 = \text{act}(\mathbf{h}^1\mathbf{W}^2)$$

$$\text{s.t. } h^2_1 = h^1_1 w^2_{11} + h^1_2 w^2_{21} + h^1_3 w^2_{31} + b w^2_{b1}$$

Neural Networks



$$\mathbf{h}^1 = \text{act}(\mathbf{x}\mathbf{W}^1)$$

$$\text{s.t. } h^1_1 = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

$$\mathbf{h}^2 = \text{act}(\mathbf{h}^1\mathbf{W}^2)$$

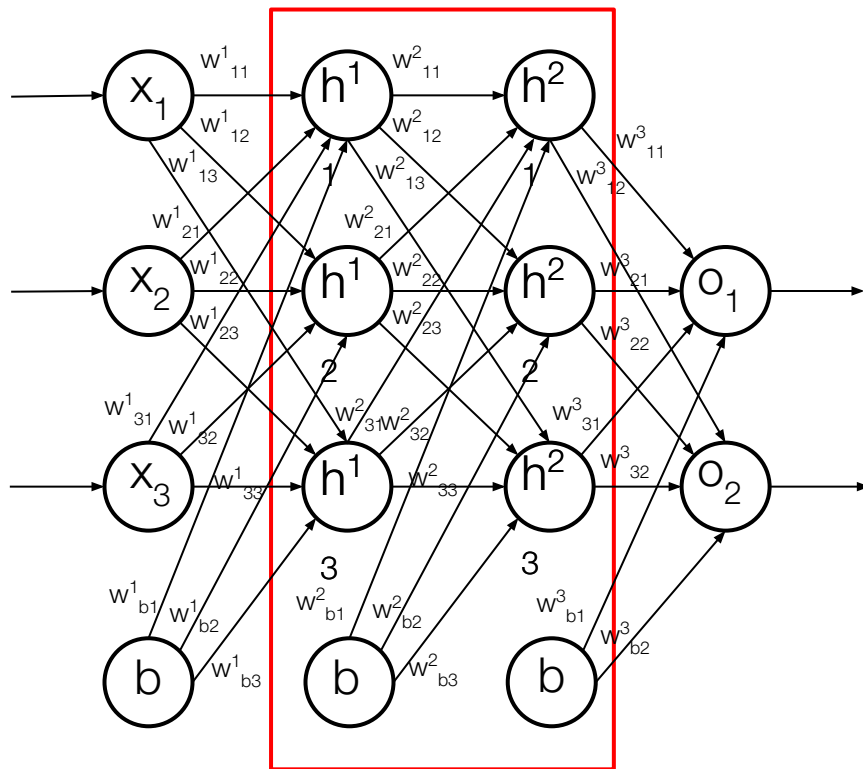
$$\text{s.t. } h^2_1 = h^1_1 w^2_{11} + h^1_2 w^2_{21} + h^1_3 w^2_{31} + b w^2_{b1}$$

$$\mathbf{o} = \text{act}(\mathbf{h}^2\mathbf{W}^3)$$

$$\text{s.t. } o_1 = h^2_1 w^3_{11} + h^2_2 w^3_{21} + h^2_3 w^3_{31} + b w^3_{b1}$$

$$\mathbf{o} = \text{act}(\text{act}(\text{act}(\mathbf{x}\mathbf{W}^1)\mathbf{W}^2)\mathbf{W}^3)$$

Neural Networks



**Representation
(Learned features)**

$$\mathbf{h}^1 = \text{act}(\mathbf{x}\mathbf{W}^1)$$

$$\text{s.t. } h^1_1 = x_1 w^1_{11} + x_2 w^1_{21} + x_3 w^1_{31} + b w^1_{b1}$$

$$\mathbf{h}^2 = \text{act}(\mathbf{h}^1\mathbf{W}^2)$$

$$\text{s.t. } h^2_1 = h^1_1 w^2_{11} + h^1_2 w^2_{21} + h^1_3 w^2_{31} + b w^2_{b1}$$

$$\mathbf{o} = \text{act}(\mathbf{h}^2\mathbf{W}^3)$$

$$\text{s.t. } o_1 = h^2_1 w^3_{11} + h^2_2 w^3_{21} + h^2_3 w^3_{31} + b w^3_{b1}$$

$$\mathbf{o} = \text{act}(\text{act}(\text{act}(\mathbf{x}\mathbf{W}^1)\mathbf{W}^2)\mathbf{W}^3)$$



How do we adjust the weights?

Stochastic Gradient Descent

- The nonlinearities make the problem non-convex (there's no algorithm with global convergence guarantees).
- The stochastic part of it allows us to consider one (or more) samples at a time, but it does not require us to process *all* samples before an update.
- We consider the impact that particular weight had in the final error we observed.
 - If we had no error, there's nothing to change :-)
 - If there was an error, we adjust the weights proportional to how big of a role they played.

But what error should we consider?

How do we figure out how big of a role a specific weight had?

But what error should we consider?

Examples of Loss Functions

Regression

Mini-batch

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N \left(y_i - \text{act}(\text{act}(\text{act}(\mathbf{x}\mathbf{W}^1)\mathbf{W}^2)\mathbf{W}^3) \right)^2$$

$\langle 5, 4 \rangle, \langle 1, 0.5 \rangle, \langle 3, 2.7 \rangle$

$$\text{MSE} = (1^2 + 0.5^2 + 0.3^2)/3 = 1.34/3 = \sim 0.45$$

Classification

$$\text{H}(x) = -\frac{1}{N} \sum_{i=1}^N \sum_x p(x_i) \log q(x_i)$$

True label

Estimate

Example, 3 classes):

True: $\langle 1, 0, 0 \rangle$ Estimate $\langle 0.7, 0.2, 0.1 \rangle$

$$\text{H}(x) = \sim 1 \times 0.36 + 0 \times 1.6 + 0 \times 2.3 = \sim 0.36$$

How do we figure out how big of a role a specific weight had?

Chain rule / Backpropagation [Rumelhart et al., 1986]

- To go from \mathbf{x} to \hat{y} we do a forward pass, or forward propagation.
- Backpropagation (or backprop) allows the information to flow back from \hat{y} in order to compute the gradient.
- Backpropagation is not the whole learning algorithm, but just the method to compute the gradients.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N \left(y_i - \text{act}(\text{act}(\text{act}(\mathbf{x}\mathbf{W}^1)\mathbf{W}^2)\mathbf{W}^3) \right)^2$$

$$\nabla_{\mathbf{x}} J(\mathbf{W}) ?$$

Chain rule / Backpropagation [Rumelhart et al., 1986]

Explanation from Goodfellow, Bengio, and Courville (2016)

- Backpropagation is an algorithm that computes the chain rule, with a specific order of operations that is highly efficient.
- Suppose that $y = g(x)$ and $z = f(g(x)) = f(y)$. The chain rule states that $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$.

- Beyond the scalar case, let $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$, with g mapping from \mathbb{R}^m to \mathbb{R}^n , and f mapping from \mathbb{R}^n to \mathbb{R} . If $\mathbf{y} = g(\mathbf{x})$ and $z = f(\mathbf{y})$, then $\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$.

In vector notation: $\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z$

where $\partial \mathbf{y} / \partial \mathbf{x}$ is the $n \times m$ Jacobian matrix of g .

- The gradient of variable \mathbf{x} can be obtained by multiplying a Jacobian matrix by a gradient. Backpropagation does so for each operation in the computation graph.

Chain rule / Backpropagation [Rumelhart et al., 1986]

Explanation from Goodfellow, Bengio, and Courville (2016)

- Suppose that $x = f(w)$, $y = f(x)$ and $z = f(y)$, that is, $z = f(f(f(w)))$.

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = f'(y) f'(x) f'(w) = f'(f(f(w))) f'(f(w)) f'(w).$$

Which one is better?

The first computes $f(w)$ only once and store it in x . This is back-propagation.

When we don't have that much memory, we can use the second.

- “The back-propagation algorithm is designed to reduce the number of common subexpressions without regard to memory (...). Back-propagation thus avoids the exponential explosion in repeated subexpressions.” [Goodfellow et al., 2016].



Optimization approaches

Stochastic gradient descent (SGD): $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha_k \hat{\mathbf{g}}$

- In practice, we need to decrease the step size* over time (thus the k) to reduce the noise introduced by the random sampling of minibatches
- Sometimes people use specific decay schemes: $\alpha_k \leftarrow (1 - \beta)\alpha_0 + \beta\alpha_{\tau}$
- There's no way around it, you need to tune the step size (and the other auxiliary variables) $\setminus_{(\text{ツ})}_{/}$
- Ideally, we would like a step size for each weight, and to adapt them automatically

*

I suggest avoiding the expression *learning rate*

RMSProp [Hinton, 2012]

$$\nu \leftarrow \rho \nu + (1 - \rho) \hat{g} \odot \hat{g}$$

Exponentially weighted moving average accumulating the gradient

$$\Delta \theta \leftarrow - \frac{\eta}{\sqrt{\nu + \epsilon}} \odot \hat{g}$$

We individually scale the step-size by the size of the accumulated gradient. Weights with large partial derivatives have a rapid decrease in the step-size

$$\theta \leftarrow \theta + \Delta \theta$$

This is what was first used by DQN. More on this later.

Adaptive Moments – Adam [Kingma and Ba, 2014]

$$\mathbf{m} \leftarrow \beta_1 \mathbf{m} + (1 - \beta_1) \hat{\mathbf{g}} \quad \leftarrow \text{Biased first moment estimate}$$

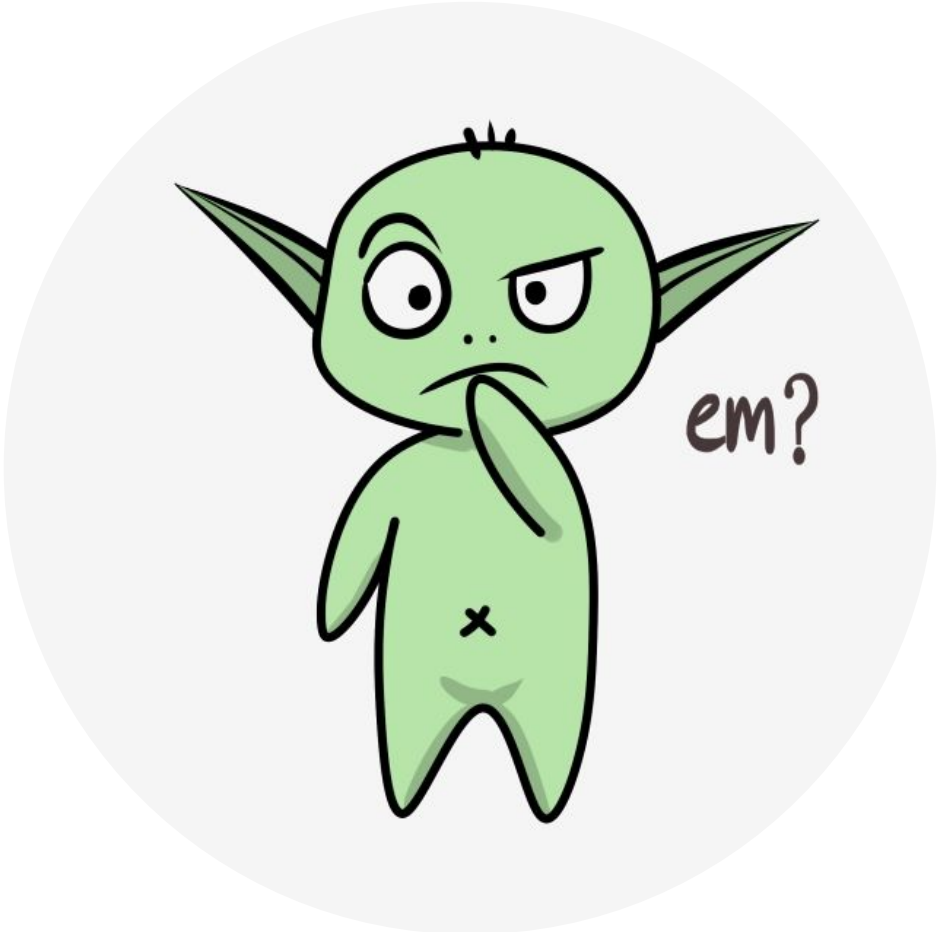
$$\boldsymbol{\nu} \leftarrow \beta_2 \boldsymbol{\nu} + (1 - \beta_2) \hat{\mathbf{g}} \odot \hat{\mathbf{g}} \quad \leftarrow \text{Biased second moment estimate}$$

$$\hat{\mathbf{m}} \leftarrow \frac{\mathbf{m}}{1 - \beta_1^t} \quad \leftarrow \text{Correct bias in first moment}$$

$$\hat{\boldsymbol{\nu}} \leftarrow \frac{\boldsymbol{\nu}}{1 - \beta_2^t} \quad \leftarrow \text{Correct bias in second moment}$$

$$\Delta \boldsymbol{\theta} \leftarrow -\alpha \frac{\hat{\mathbf{m}}}{\sqrt{\hat{\boldsymbol{\nu}} + \epsilon}} \quad \leftarrow \text{RMSProp-like update}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$$



More Optimization for Deep Learning

Initialization *Really* Matters

- Many approaches try to ensure the weight matrix has nice properties before optimization begins

Uniform initialization

$$W_{i,j} \sim U\left(-\frac{1}{\sqrt{\text{fan}_{in}}}, \frac{1}{\sqrt{\text{fan}_{in}}}\right)$$

Xavier (or Glorot) initialization [Glorot and Bengio, 2010]

$$W_{i,j} \sim U\left(-\sqrt{\frac{6}{\text{fan}_{in} + \text{fan}_{out}}}, \sqrt{\frac{6}{\text{fan}_{in} + \text{fan}_{out}}}\right)$$

He (or Kaiming) initialization [He et al., 2015]

$$W_{i,j} \sim \mathcal{N}\left(0, 2/\text{fan}_{in}\right)$$

Regularization: Parameter Norm Penalties

$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \alpha\Omega(\boldsymbol{\theta})$$

Generally we don't regularize the bias term

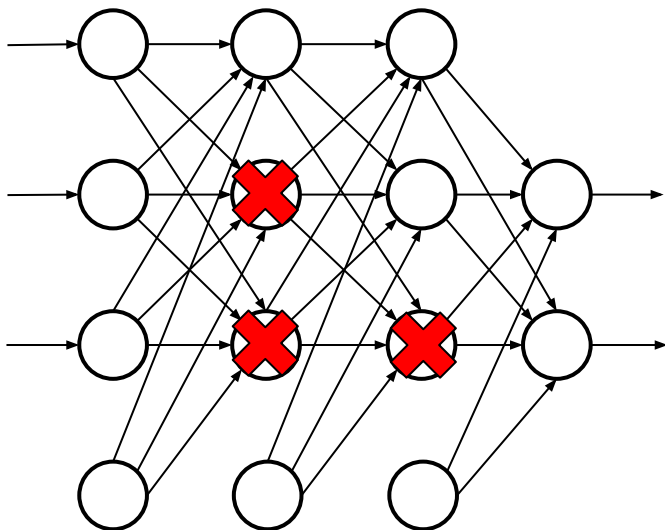
$$\Omega(\boldsymbol{\theta}) = \frac{1}{2}\|\mathbf{w}\|_2^2 \quad \Omega(\boldsymbol{\theta}) = \|\mathbf{w}\|_1 = \sum_i |w_i|$$

ℓ_2 , weight decay

ℓ_1 , sparsity inducing

Regularization: Dropout [Srivastava et al., 2014]

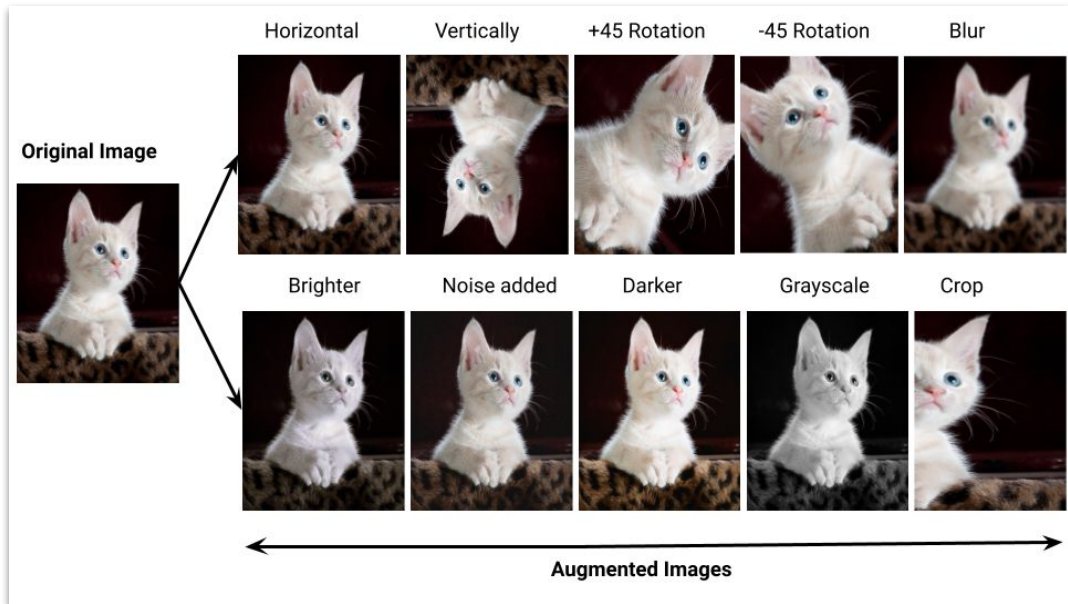
- Dropout is a cheap approximation for training multiple models at the same time.
- Dropout consists in training subnetworks obtained by removing nonoutput units from the original network. A unit needs to perform well regardless of the others.



- For each new minibatch, randomly sample a different binary mask.
E.g., Hidden units are generally dropped w.p. 0.5.
- For inference, we multiply each weight by the probability of them not being dropped (weight scaling inference rule).
E.g., Divide weights by 2 at the end of training.

Regularization: Data Augmentation

- How do we generate more data without actually having more data? Fake data!
 - Depending on the problem, it is relatively easy to generate fake data.



Batch Normalization Layer [Ioffe and Szegedy, 2015]

- Instead of normalizing only the inputs, we normalize the input to every layer.
- “To shift and rescale each activation so that its mean and variance across the batch become values that are learned during training” [Prince, 2023].
 - Ensures variance is stable during forward pass at initialization.
 - Loss surface and its gradient change more smoothly, thus we can use larger step sizes.
 - Regularization. Batch normalization adds noise because of the differences across batch statistics.

$$\textcircled{1} \quad \hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mu_B^{(k)}}{\sqrt{(\sigma_B^{(k)})^2 + \epsilon}}$$

$$\textcircled{2} \quad y_i^{(k)} = \gamma^{(k)} \hat{x}_i^{(k)} + \beta^{(k)}$$

**learnable
parameters**

$\textcircled{3}$ We keep track of the exponential moving average so we can use that in inference (remember, we have one sample now)

Layer Normalization [Ba et al., 2016]

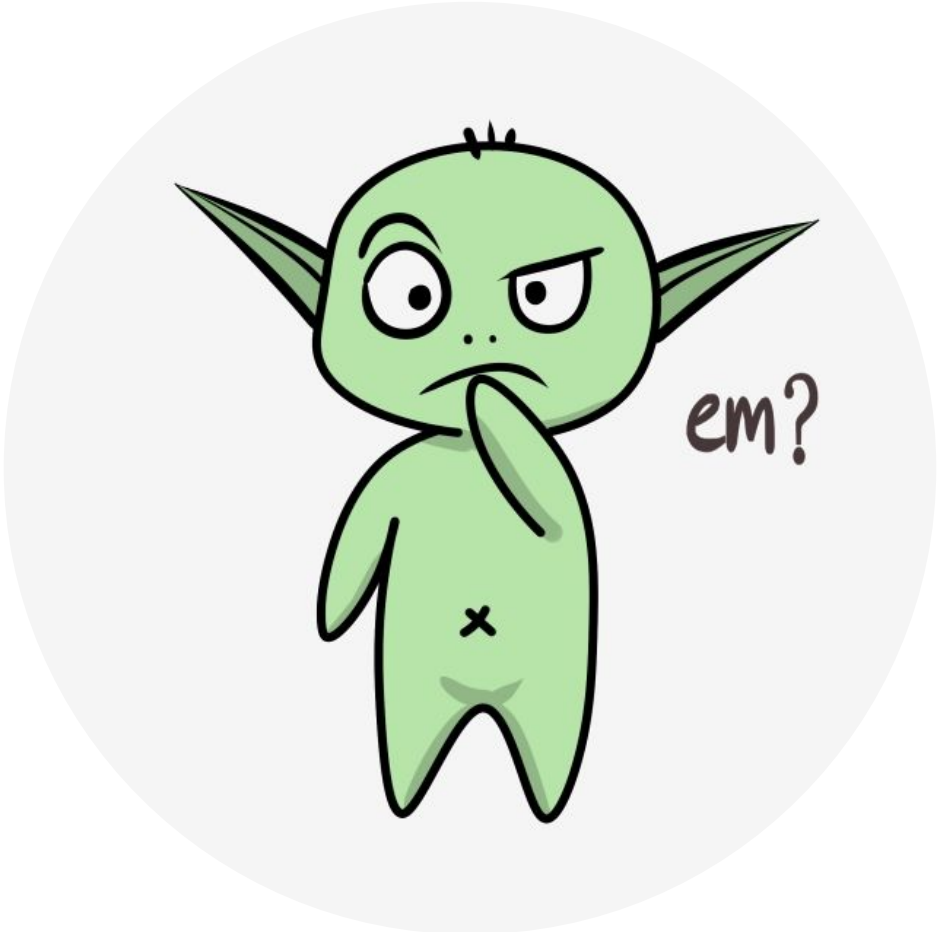
- Instead of normalizing individual features within a batch, we normalize the activations of each layer across the entire dataset.
 - Instead of computing the mean and std. dev., μ_B and σ_B , within a batch, we do so across all features (i.e. units), getting μ_L and σ_L .
 - Obviously, it has no dependency on the batch size. Thus, more amenable to sequences.
 - Similar to batch normalization, it can be applied before or after the activation function.

$$\textcircled{1} \quad \hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mu_L^{(k)}}{\sqrt{(\sigma_L^{(k)})^2 + \epsilon}}$$

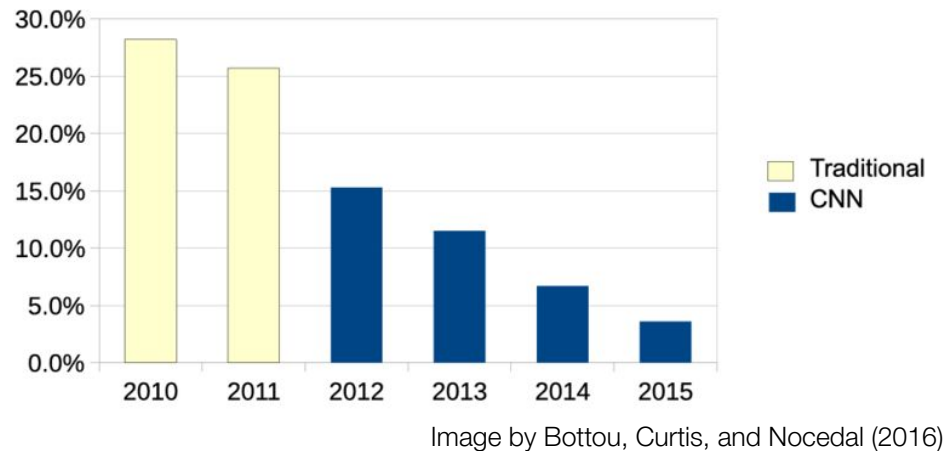
$$\textcircled{2} \quad y_i^{(k)} = \gamma^{(k)} \hat{x}_i^{(k)} + \beta^{(k)}$$

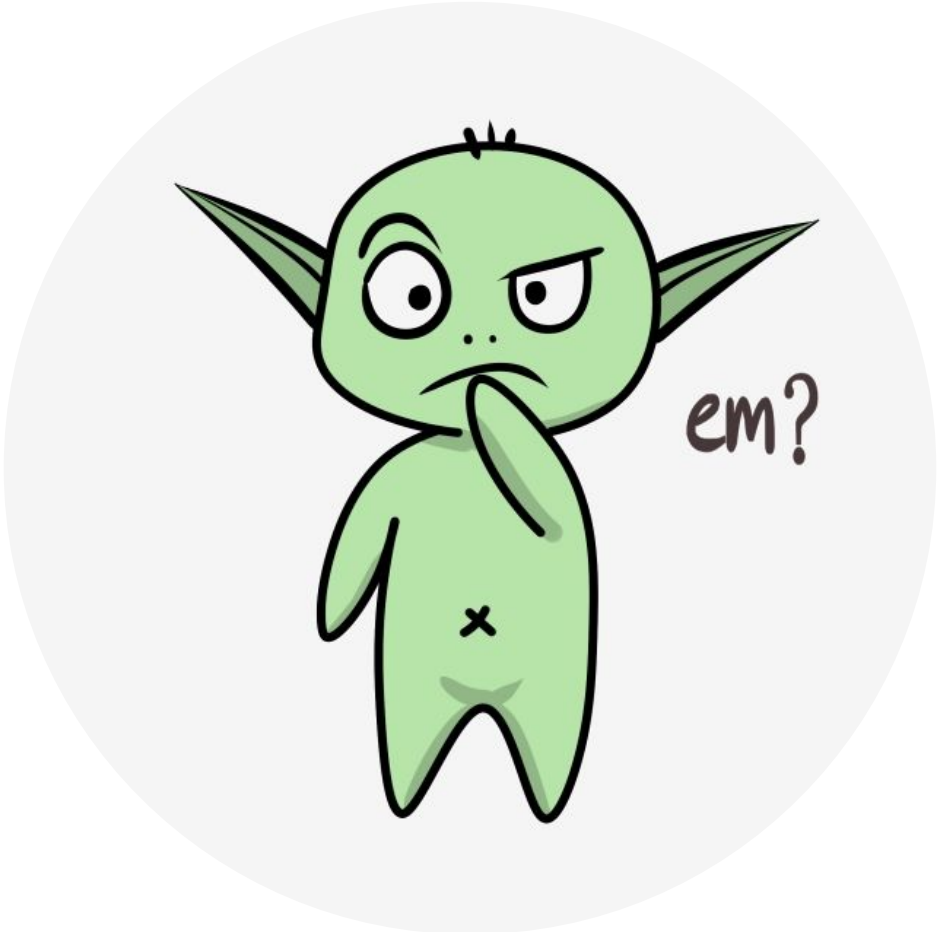
**learnable
parameters**

$\textcircled{3}$ No need for special adjustments for inference.



Are we improving anything? From features to architecture...





Convolutional Neural Networks

Deep Convolutional Network

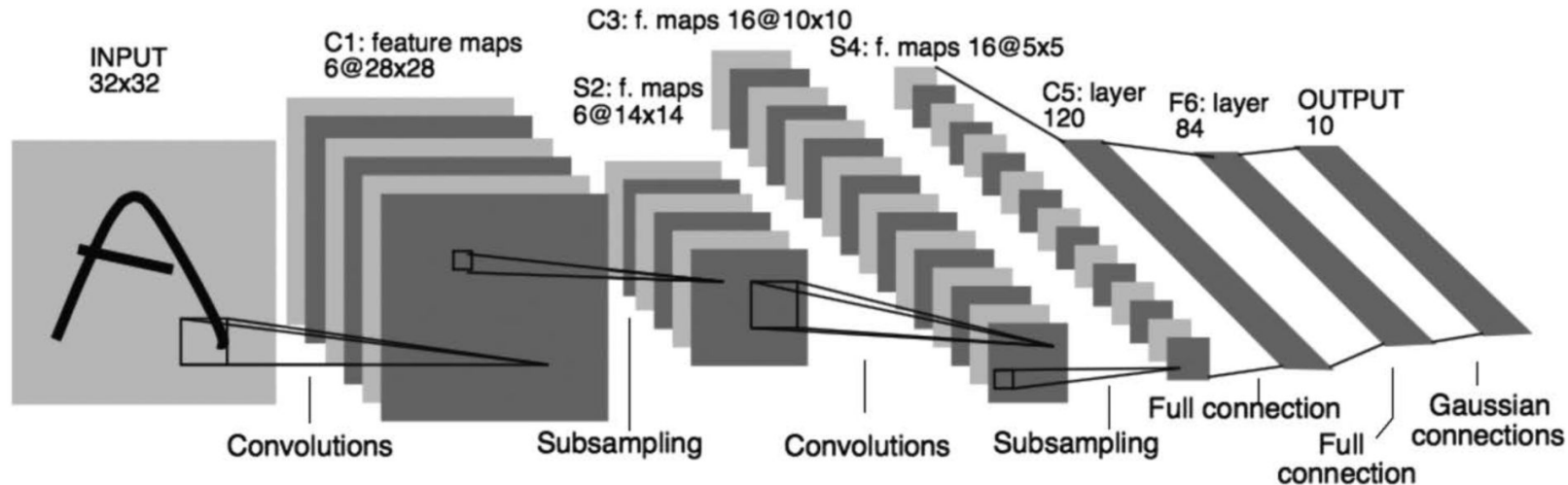
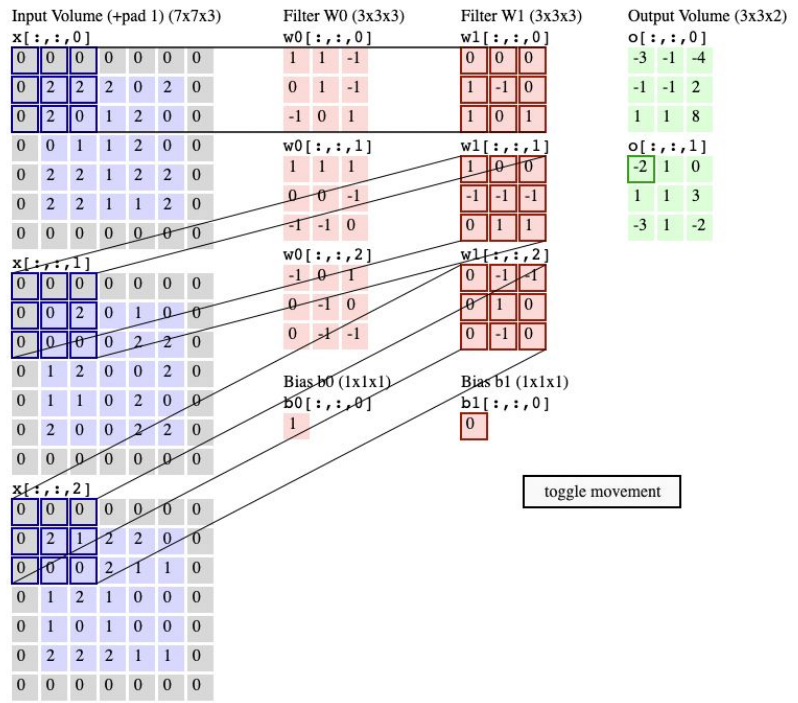


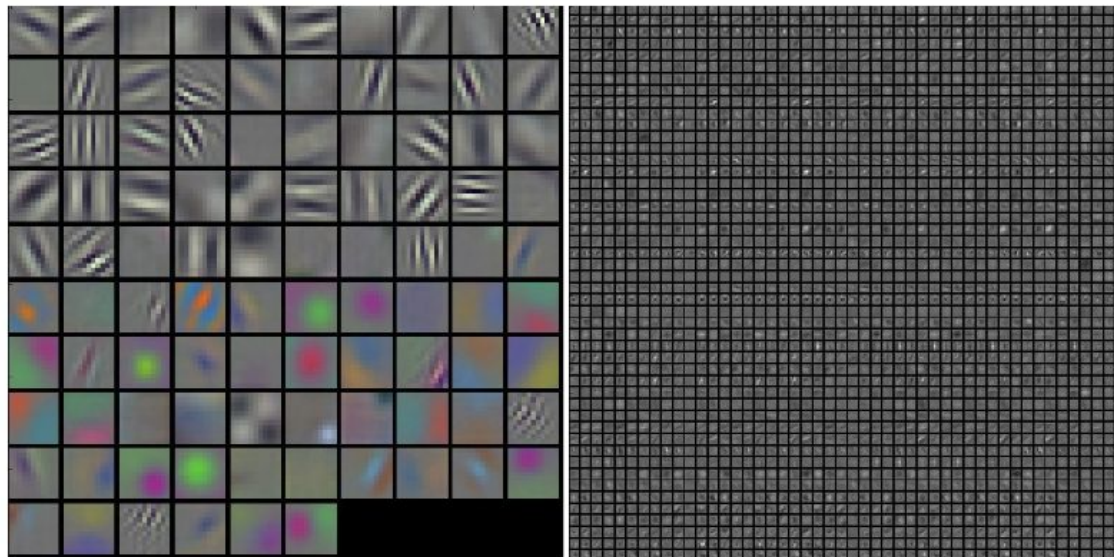
Figure 9.15: Deep Convolutional Network. Republished with permission of Proceedings of the IEEE, from Gradient-based learning applied to document recognition, LeCun, Bottou, Bengio, and Haffner, volume 86, 1998; permission conveyed through Copyright Clearance Center, Inc.

Deep Convolutional Network

Notice weights are shared

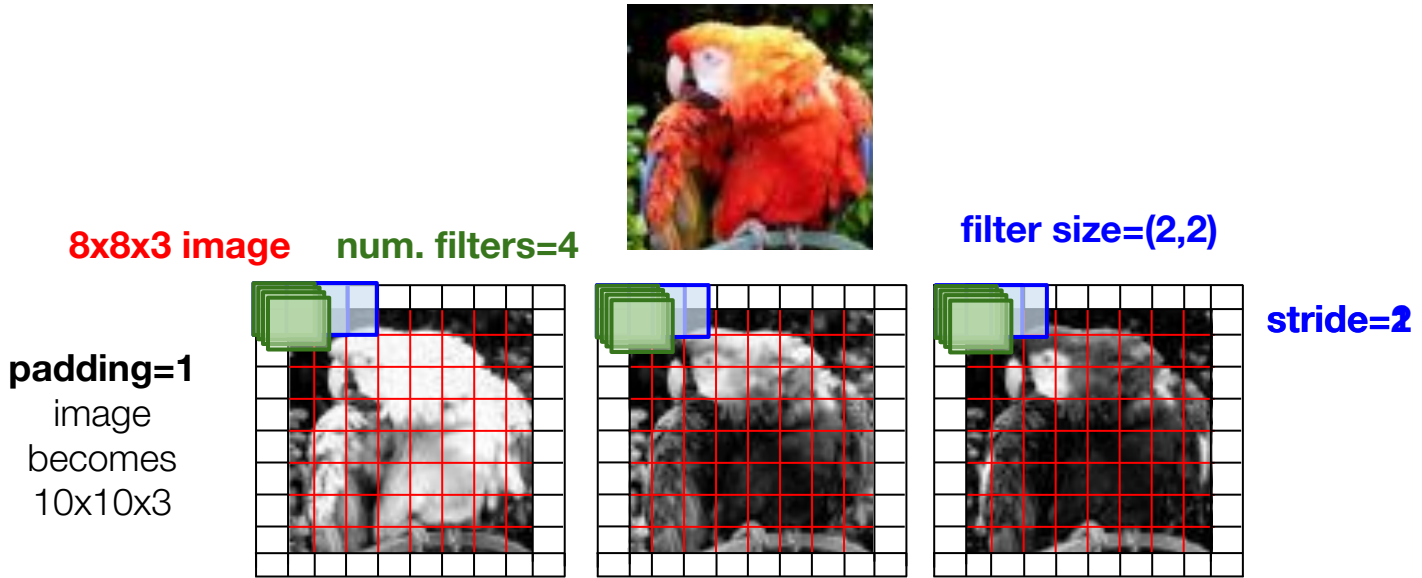


Learned Representations



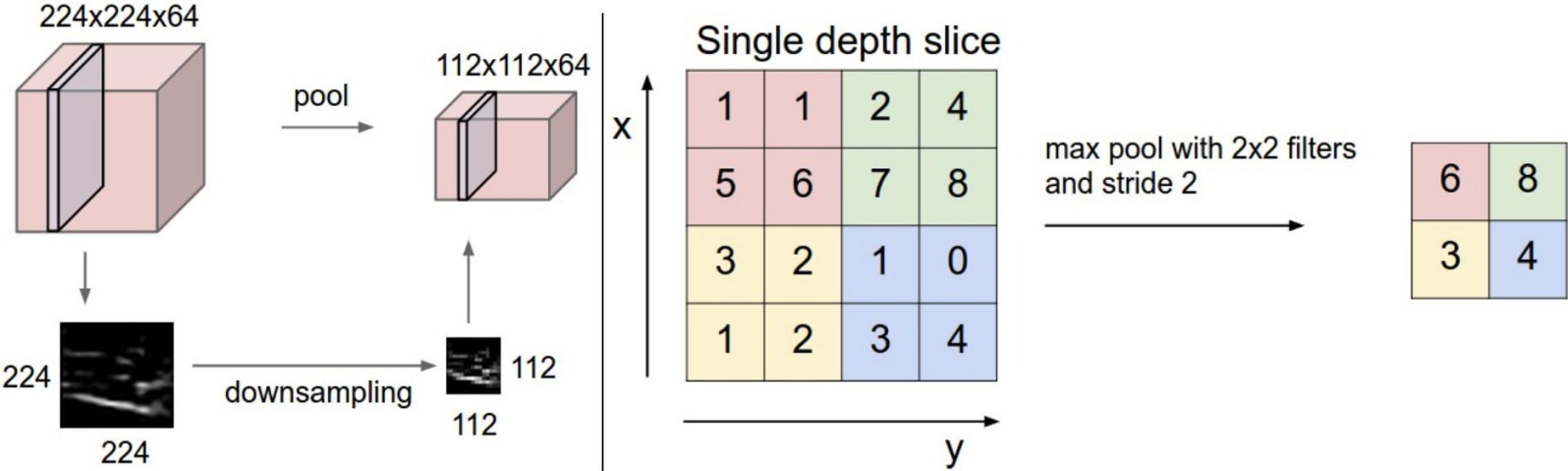
Typical-looking filters on the first CONV layer (left), and the 2nd CONV layer (right) of a trained AlexNet. Notice that the first-layer weights are very nice and smooth, indicating nicely converged network. The color/grayscale features are clustered because the AlexNet contains two separate streams of processing, and an apparent consequence of this architecture is that one stream develops high-frequency grayscale features and the other low-frequency color features. The 2nd CONV layer weights are not as interpretable, but it is apparent that they are still smooth, well-formed, and absent of noisy patterns.

The hyperparameters of a convolutional neural network



Size of output: $(W - F + 2P)/S + 1$

Pooling



Pooling layer downsamples the volume spatially, independently in each depth slice of the input volume. Left: In this example, the input volume of size [224x224x64] is pooled with filter size 2, stride 2 into output volume of size [112x112x64]. Notice that the volume depth is preserved. Right: The most common downsampling operation is max, giving rise to max pooling, here shown with a stride of 2. That is, each max is taken over 4 numbers (little 2x2 square).

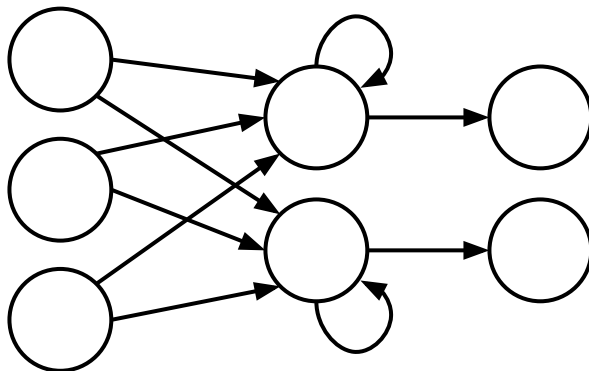


Recurrent Neural Networks

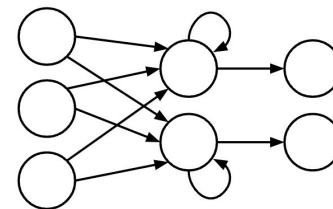
Recurrent Neural Networks (RNNs)

- RNNs are generally used to process *sequential* data.
- Similarly to how convolutional networks share parameters across patches of an image, RNNs share parameters across positions in a sequence.

$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta})$$



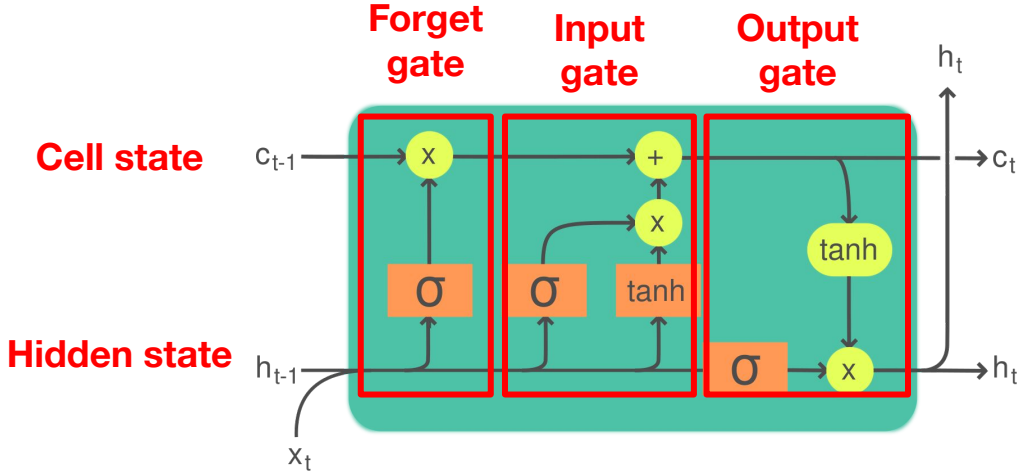
Recurrent Neural Networks (RNNs)



- RNNs are often trained using backpropagation through time, which like regular backpropagation but applied to an *unrolled* computational graph.
- They need to be trained in sequence.
- Their recurrence potentially gives us a lossy *memory*.
- In reinforcement learning, they are a natural instantiation of an agent state (not to be confused with the environment state nor the observation).

$$S_{t+1} \doteq u(S_t, A_t, O_{t+1})$$

Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber, 1997]



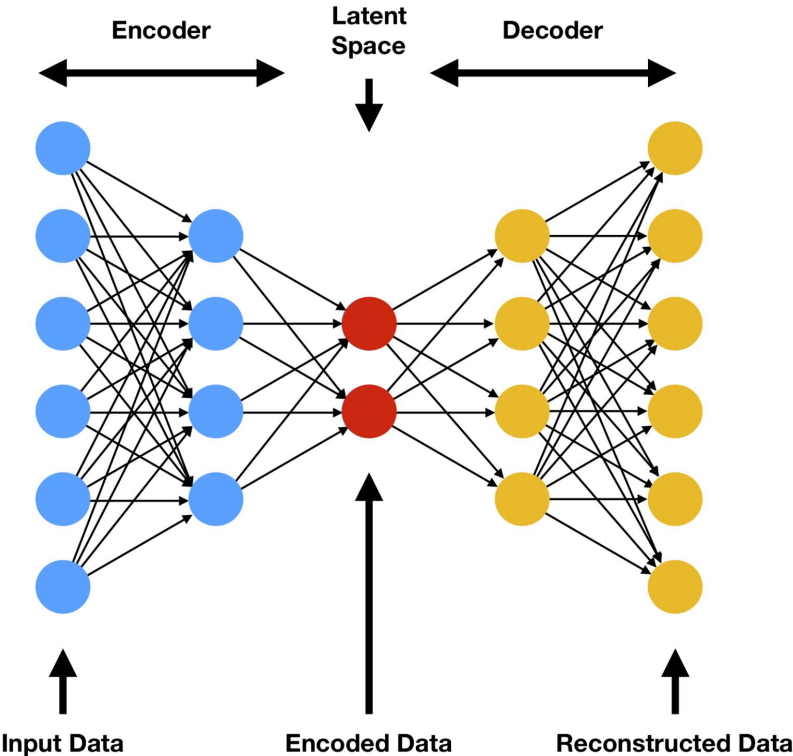
Legend:

Layer	Componentwise	Copy	Concatenate

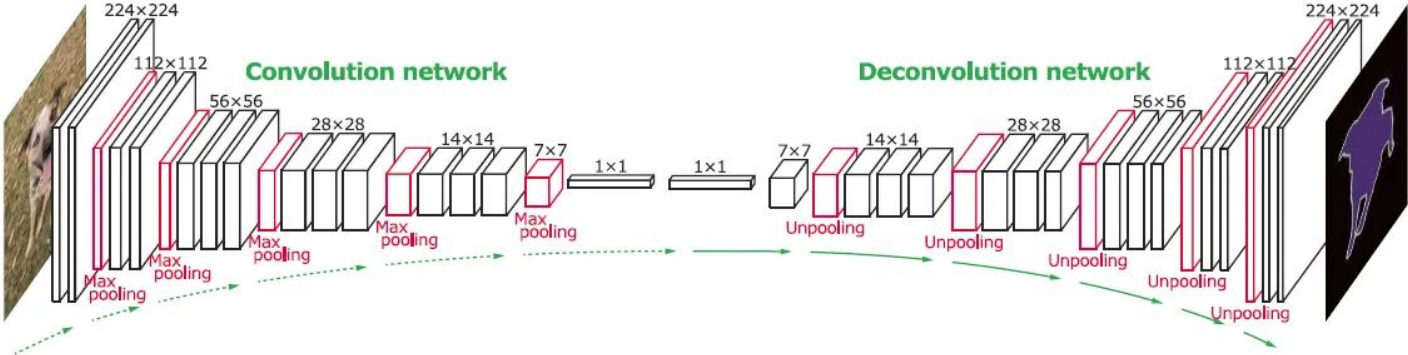


Other Neural Network Architectures

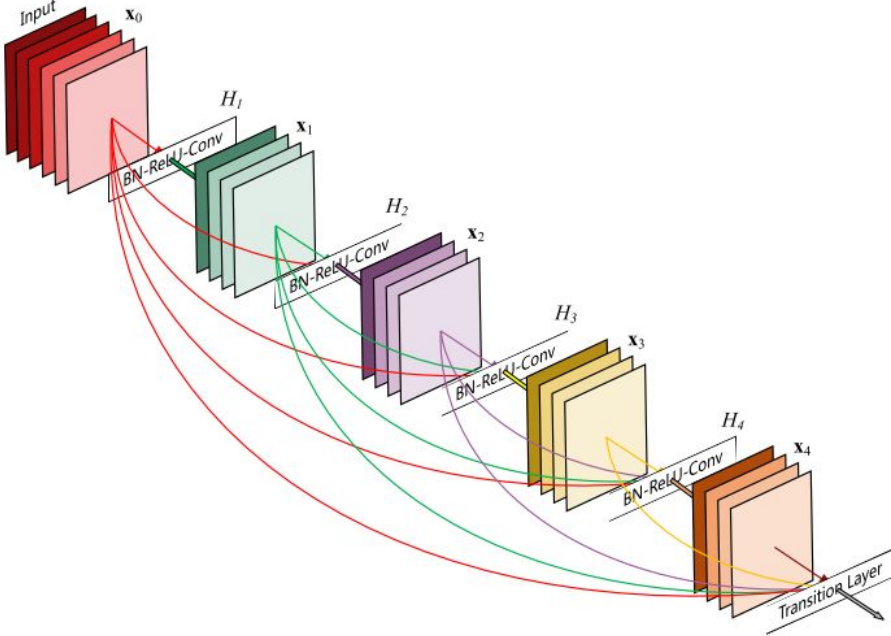
Autoencoders



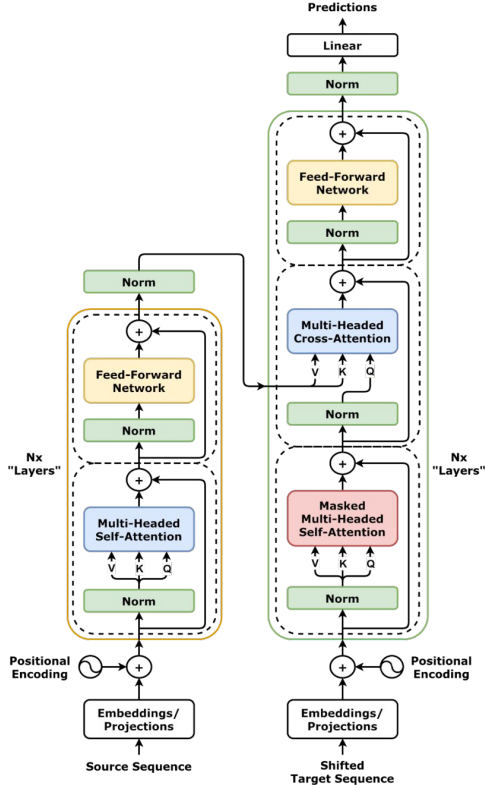
Inverse Convolutional Network



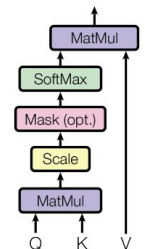
Residual Networks



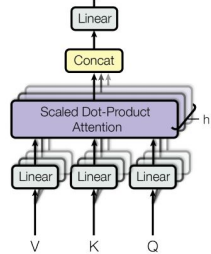
Transformers [Vaswani et al., 2017]



Scaled Dot-Product Attention



Multi-Head Attention





Next class

- What I plan to do:
 - Finish overview of neural networks / deep learning.
 - Start talking about Deep RL; more specifically, DQN.
 - I'll try to make more lecture notes available.

- What I recommend YOU to do for next class:
 - *Brush-up* on the basics of deep learning if you don't remember.
Specifically, Goodfellow, Bengio & Courville (2016)'s chapters 6–10.
 - Read the DQN paper: *V. Mnih et al.: Human-level control through deep reinforcement learning. Nature 518(7540): 529-533 (2015).*