

A Very Brief Introduction to Neural Networks and Deep Learning

Ashley Lee, Isabel Restrepo, & Paul Stey

December 8, 2017

Table of Contents

1 Background

- Neural Network Basics
- History
- Mechanics of Neural Networks

2 CNNs

- Convolutional Neural Network (CNNs)
- Real-World Constraints
- Pooling Layer
- Normalization Layer

3 RNNs

- RNNs and LSTMs

4 Autoencoders

- Autoencoder Basics
- Stacked and Denoising Autoencoders

5 Summary

- When to use Neural Networks

What is a neural network?

- ① A species of directed acyclic graphs (usually)
 - ② “Universal function approximator”
 - ③ “An engineering solution to a statistics problem”

What do neural networks do?

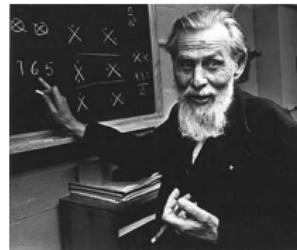
Like many other statistical or machine learning models (e.g., GLM, random forests, boosting), neural networks:

- ① Attempt to approximate a data-generating mechanism
 - ② Can be used for classification problems
 - ③ Can be used for regression problems
 - ④ Can also be used for dimension reduction like principal components analysis (PCA)

History of Neural Networks

The history of neural networks is long and tumultuous.

- ① McCulloch and Pitts (1943) “*A Logical Calculus of Ideas Immanent in Nervous Activity*”
 - ② Rosenblatt (1958) “*The Perceptron: A Probabilistic Model For Information Storage And Organization In The Brain*”

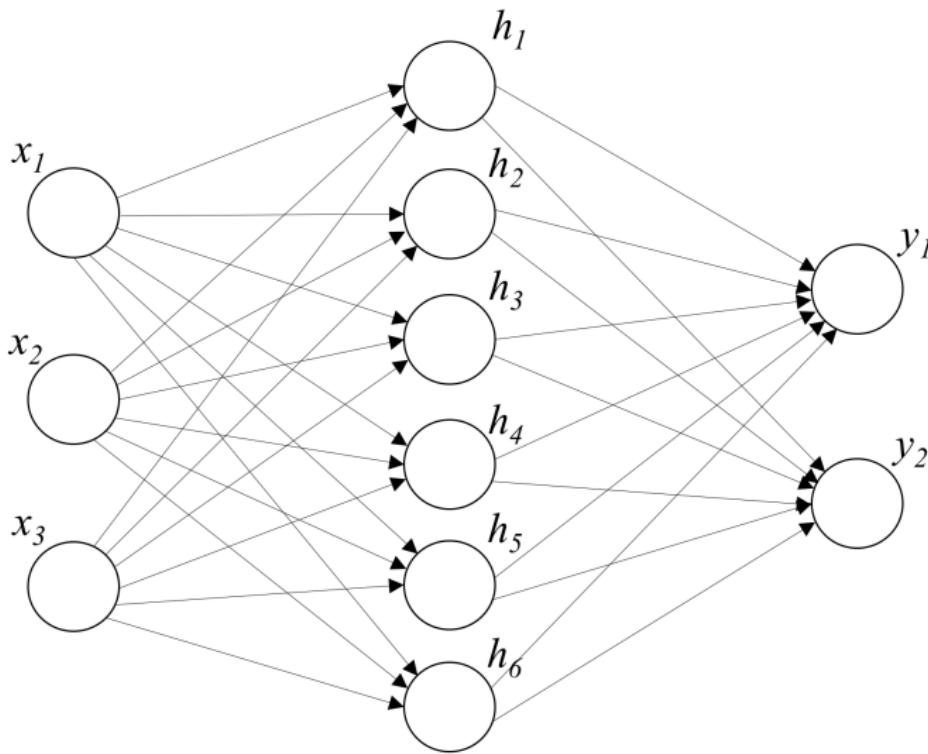


More Recently

Neural networks are experiencing a major resurgence. There are at least two reasons.

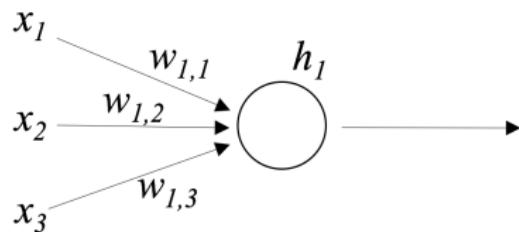
- ➊ Better algorithms for back-propagation
 - ➋ GPUs are well suited to building neural networks
 - Matrix multiplies can be made embarrassingly parallel
 - GPUs have much better memory bandwidth
 - ➌ More labeled data

The Multilayer Perceptron



Single Neuron

A single neuron takes inputs, x_j , and applies the weights, w_j to the input by computing the dot product of the vectors x and w . The result is the input to the “activation” function.



Activation Functions

The notion of an activation function comes again from the theoretical relationship to neurons in the brain.

Activation functions are analogous to “link” functions in generalized linear models (GLMs).

In fact, one common activation function is the sigmoid function, which is just our old friend the logistic function which you are using when you fit logistic regression models.

Purpose of Activation Functions

There are a few reasons we use activation functions.

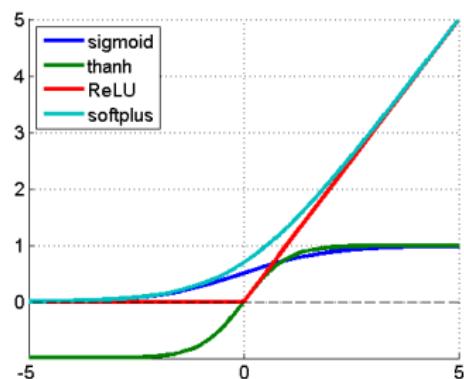
The most basic reason is that—like link functions in GLMs—we want to take some linear predictor and transform it so that it is bounded appropriate. For instance, the value of logistic function is in the range $(0, 1)$.

And the second key reason is that this allows us to introduce non-linearities. Recall that a neural network (like many statistical or machine learning models) is trying to approximate a data-generating mechanism. So we are trying to approximate a function that might be very complex and include many non-linearities.

Common Activation Functions

Some common activation functions include the following:

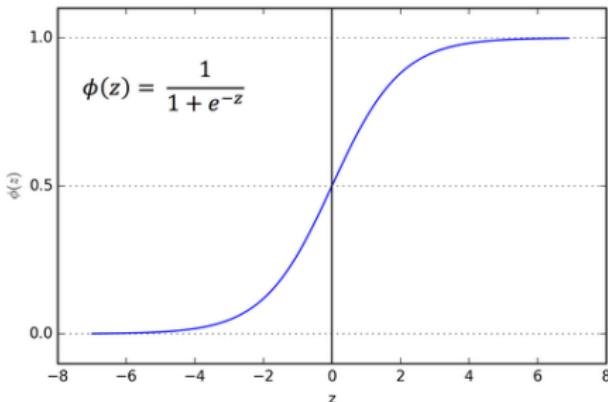
- ① Sigmoid (i.e., logistic)
 - ② Hyperbolic tangent: \tanh
 - ③ Rectified linear unit (ReLU)
 - ④ softplus



Sigmoid Functions

Sigmoid function: $\phi(z) = \frac{1}{1+e^{-z}}$

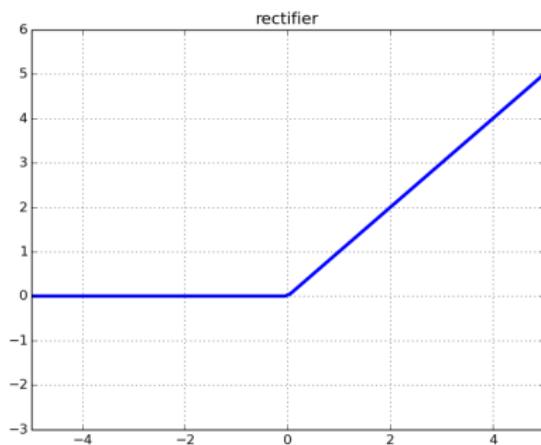
- ① Range between 0 and 1
 - ② Sometimes interpreted as probability
 - ③ Special case of the softmax:
 $\psi(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$, which is used for vector-value outcome



Rectified Linear Unit (ReLU)

ReLU: $f(z) = \max(0, z)$

- ① Range between 0 and inf
 - ② Biological justification



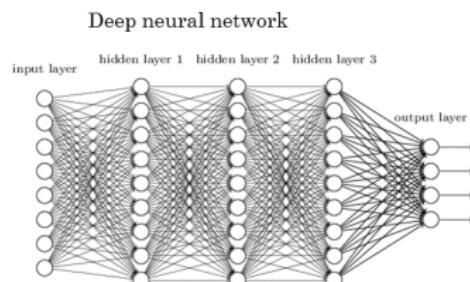
Activation Functions

For the hidden layers of neural networks, ReLUs—or some variation thereof—tend to be the most common.

Sigmoid (and probably $tanh$) are not frequently used in the hidden layers more because of the vanishing gradient problem (discussed later).

Convolutional Neural Networks

Why do we need ConvNets?

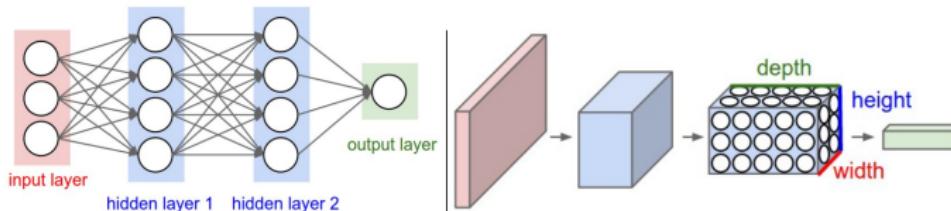


- ① Regular neural nets don't scale well to images
 - For images of size $32 \times 32 \times 3$, a *single* fully-connected neuron in the first layer would have 3072 weights.
 - Images of size $200 \times 200 \times 3$, a *single* gives 120000 weights.
- ② Full connectivity is wasteful and the huge number of parameters would quickly lead to overfitting.

CNNs

What are ConvNets?

- ① ConvNets are very similar to neural networks discussed thus far. Dot product, followed by non-linearity, and loss function at the end.
- ② Explicit assumption that input are images.
- ③ Layers have neurons arranged in 3 dimensions (width, height, depth) to form an **activation volume**



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

CNN Architecture

Types of layers used to build ConvNets

① Convolutional Layer

- Input: 3-d volume
- Output: 3-d volume
- Convolve “filters” with small regions in the image
- Output depth, depends on the number of filters

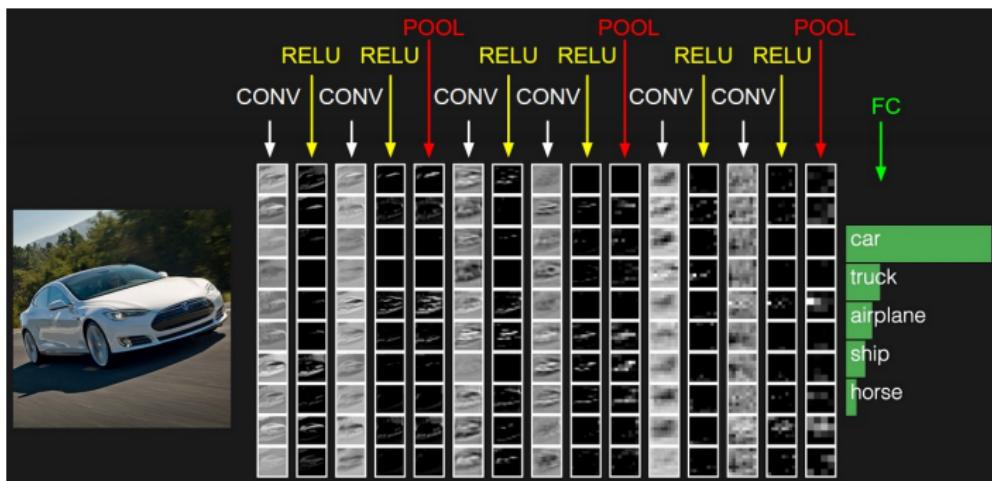
② Pooling Layer

- Downsampling along spatial dimensions (width, height)

③ Fully-Connected Layer (what we've seen so far)

- Compute class score. Dimensions are transformed to $1 \times 1 \times k$, where k is number of classes

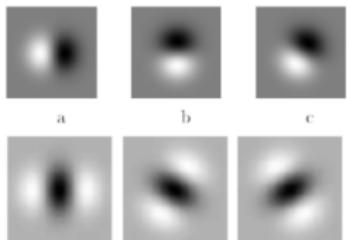
CNN Architecture



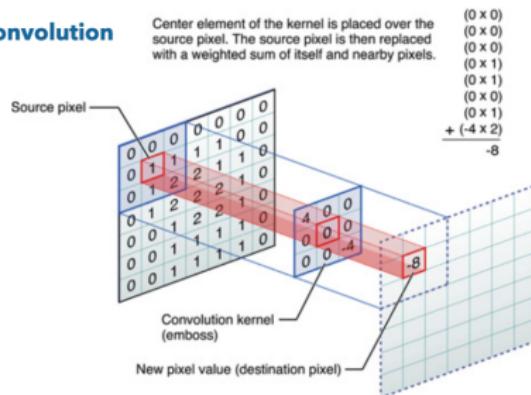
- ① CONV/FC layers perform transformations that are a function of not only the activations in the input volume, but also of the parameters (weights and biases of neurons).
- ② RELU/POOL layers will implement a fixed function (e.g., $\text{ReLU} = \max(0, x)$)

Classical Computer Vision

1. Filter Design. E.g., Steerable filters - Gaussian and it's derivatives



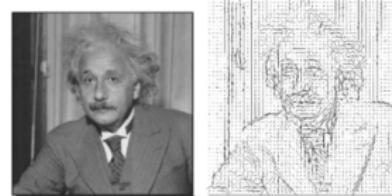
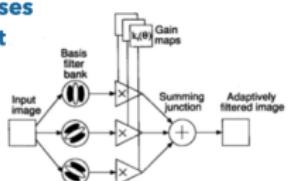
2. Convolution



3. Individual responses (apply a, b, c)



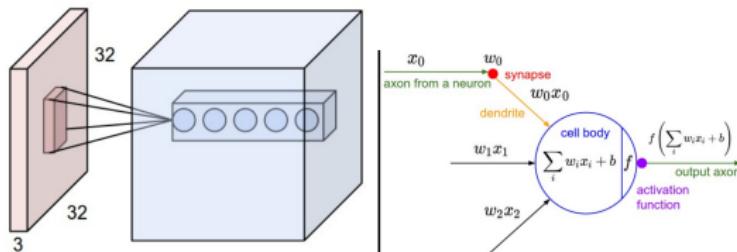
4. Compose responses (optional) - gradient orientation



Convolutional Layer - Local Connectivity

Connecting each Neuron to the Input Volume

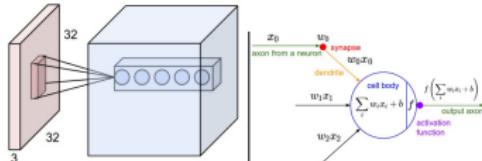
- ① **Receptive Field (hyper-parameter):** Spatial extend of each neuron = filter size
- ② Connections are local in space (along width and height), but always full along the entire depth of the input volume.
 - Example: Suppose an input volume had size $16 \times 16 \times 20$. Then using an example receptive field size $3 \times 3 \times 20 = 180$ connections to the input volume. Notice that, again, the connectivity is local in space (i.e., 3×3), but full along the input depth (20)

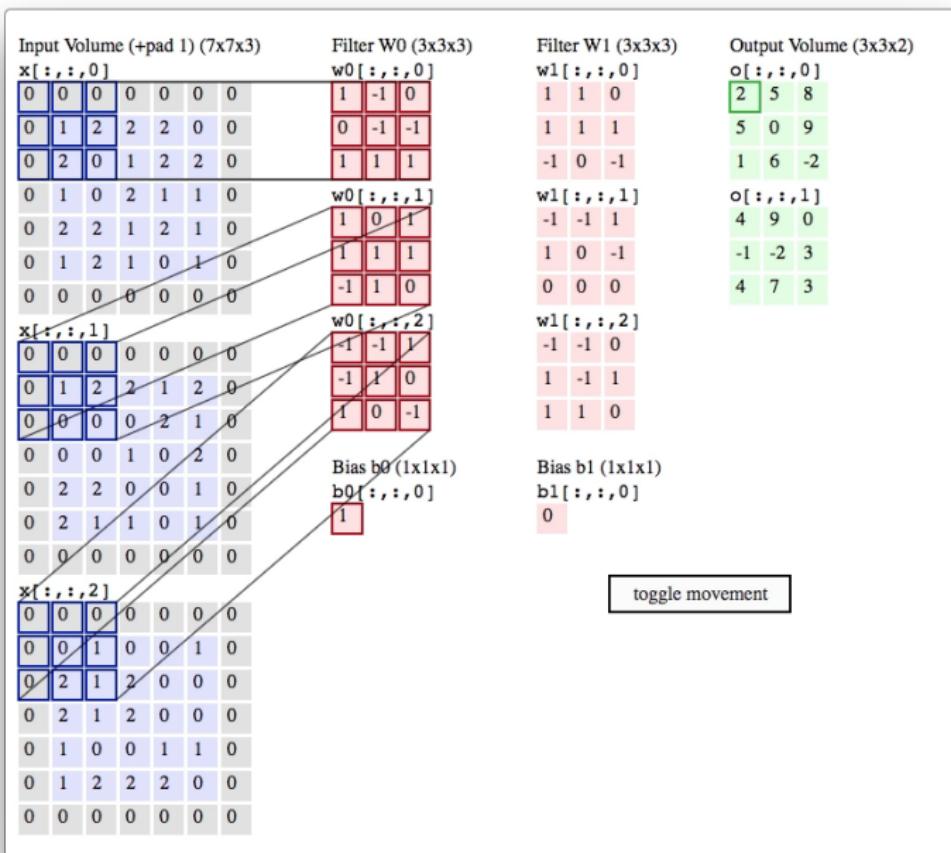


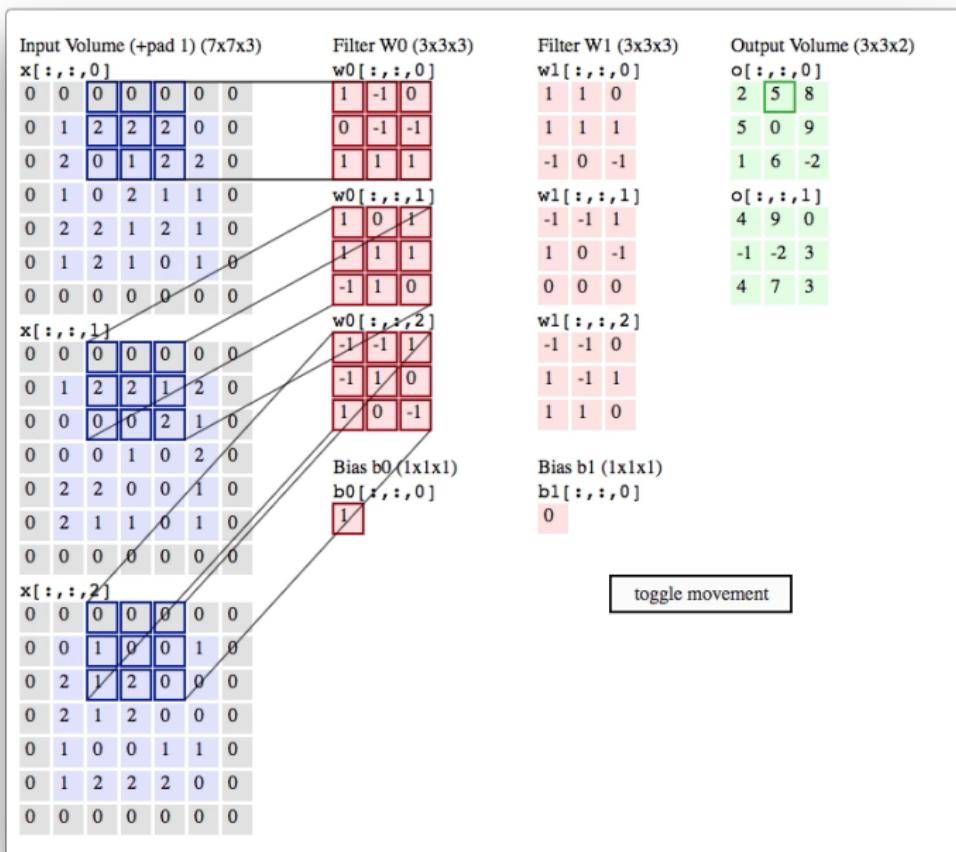
Convolutional Layer - Spatial Arrangement

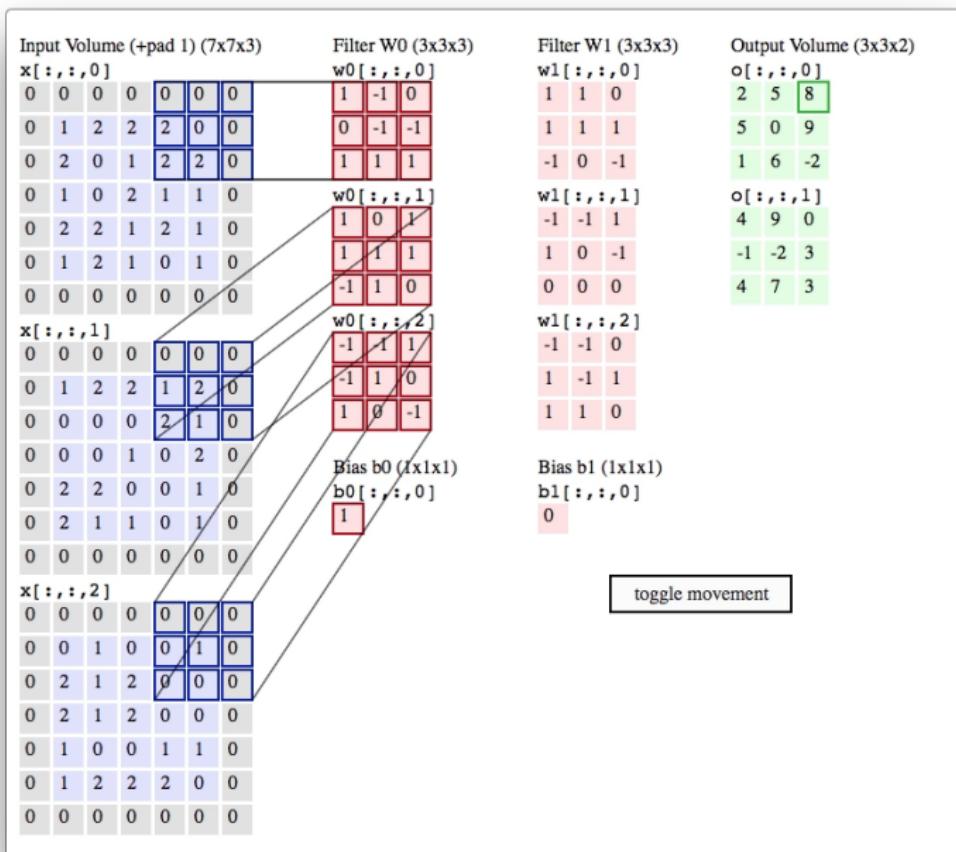
Number of Neurons in the Output Volume

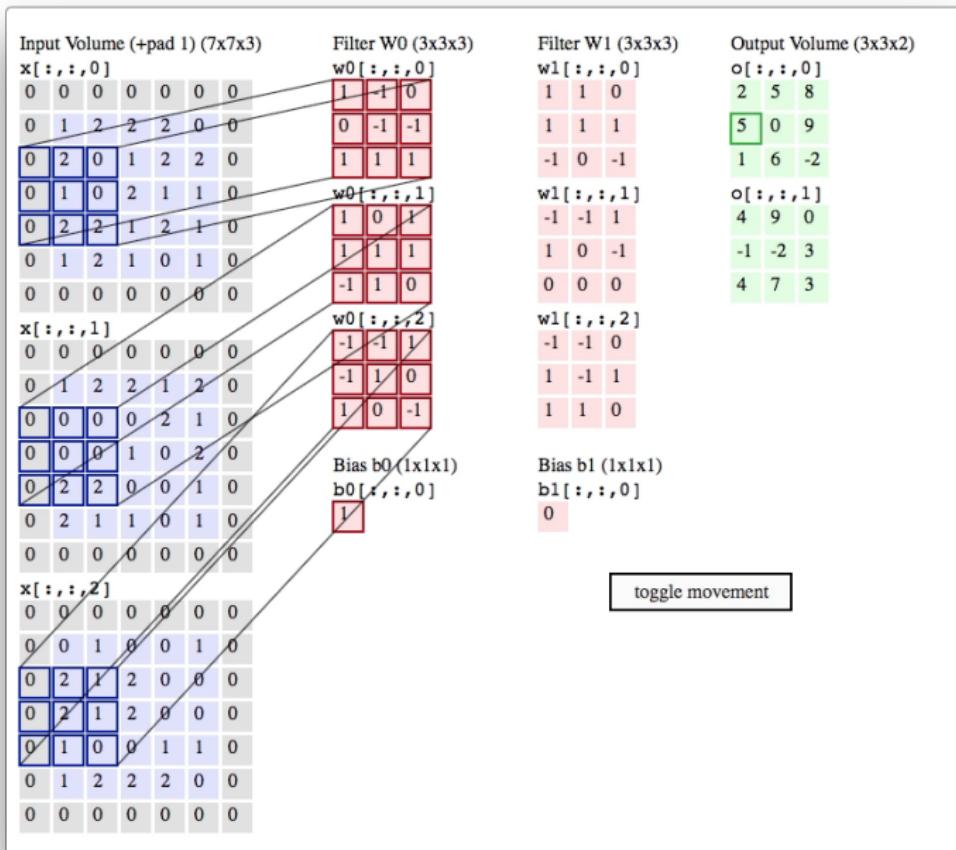
- ➊ **Depth (hyper-parameter):** Number of filters: Number of patterns to look for in input. Think back to steerable filters. The set of neurons looking at the same region may be referred to as **depth column** or **fiber**.
- ➋ **Stride (hyper-parameter):** Number of pixels used when sliding the filter. This controls the overlap between neurons. If stride is 1, we move one pixel at a time (common, but depends on the width of filter).
- ➌ **Padding (hyper-parameter):** Number of zeros to add to the border of the volume. Convenient for convolution.

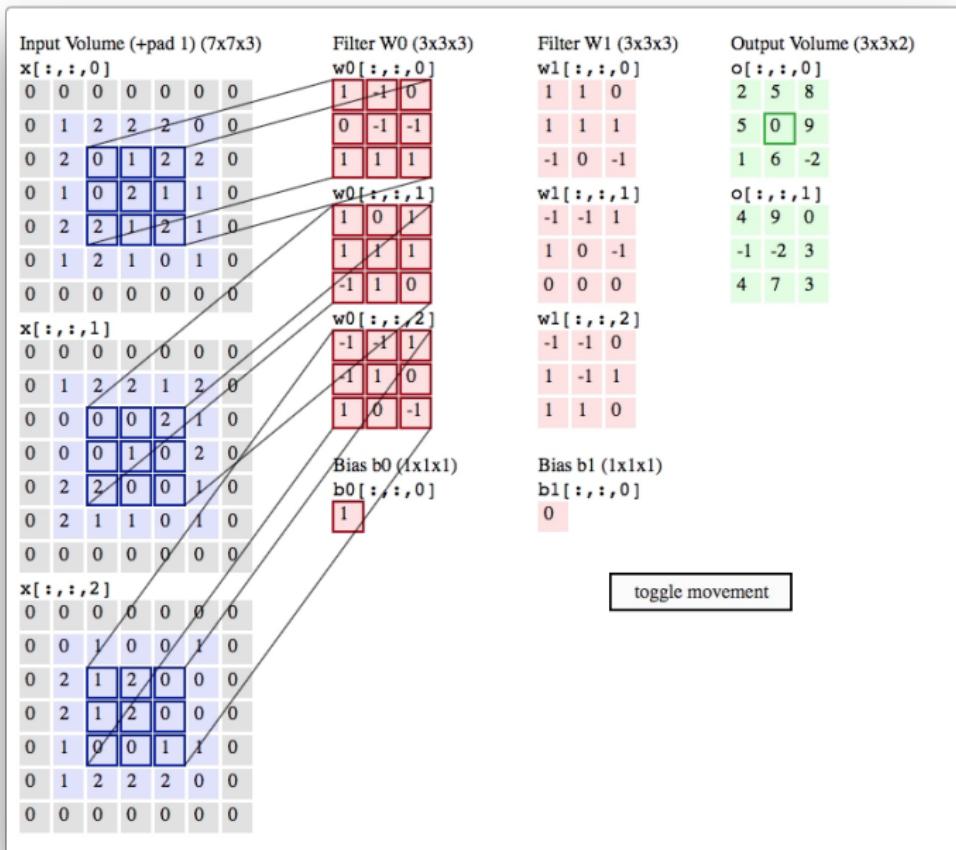


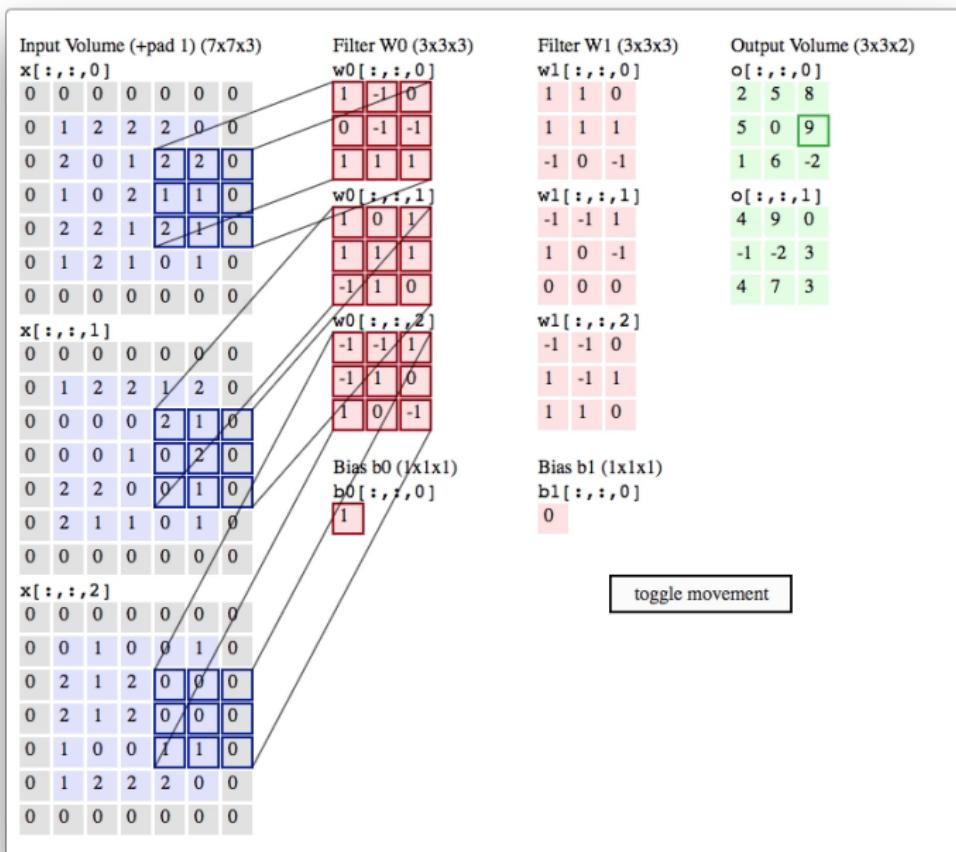


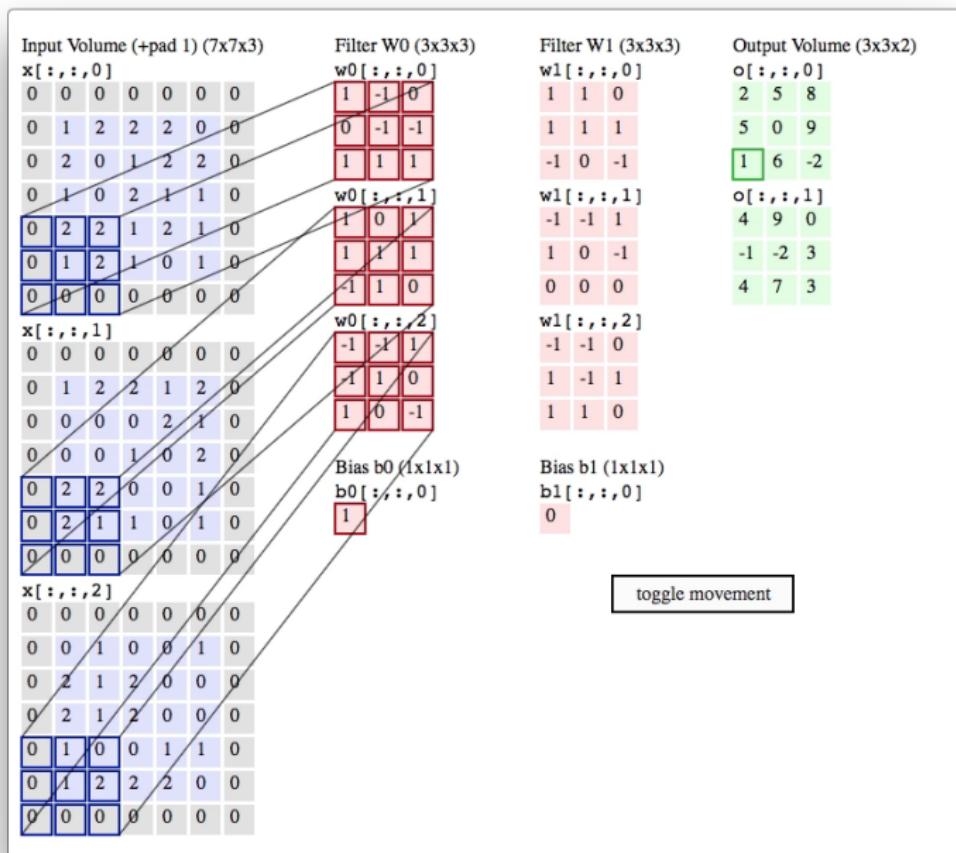


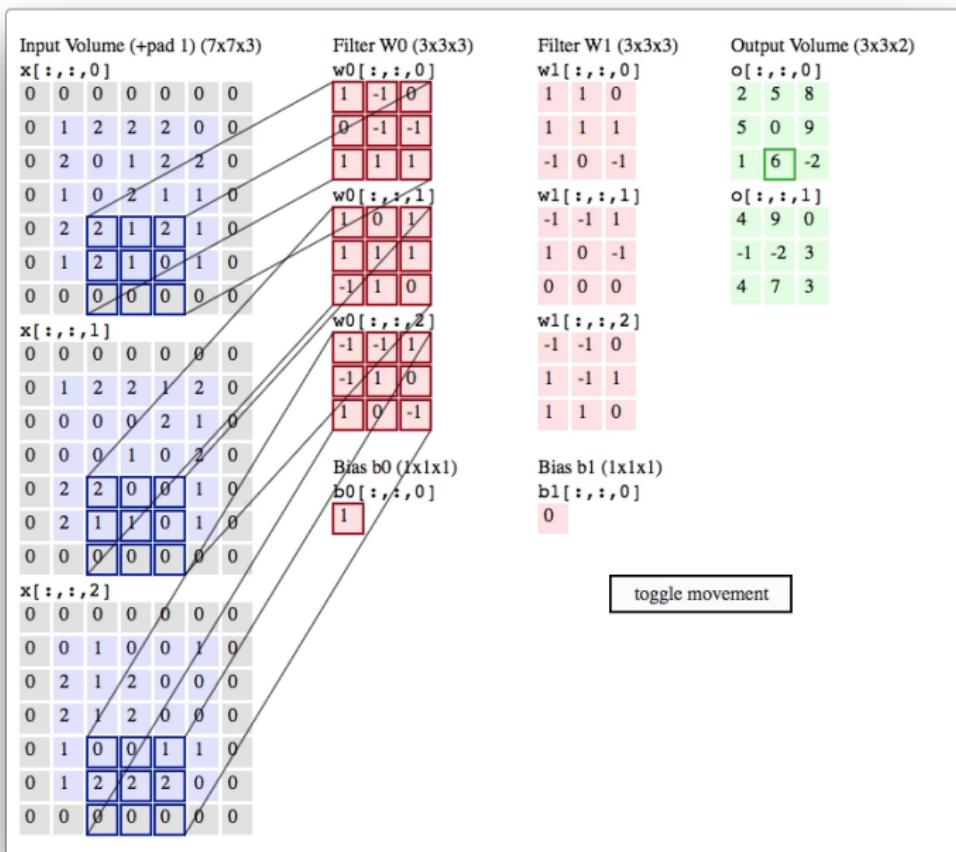


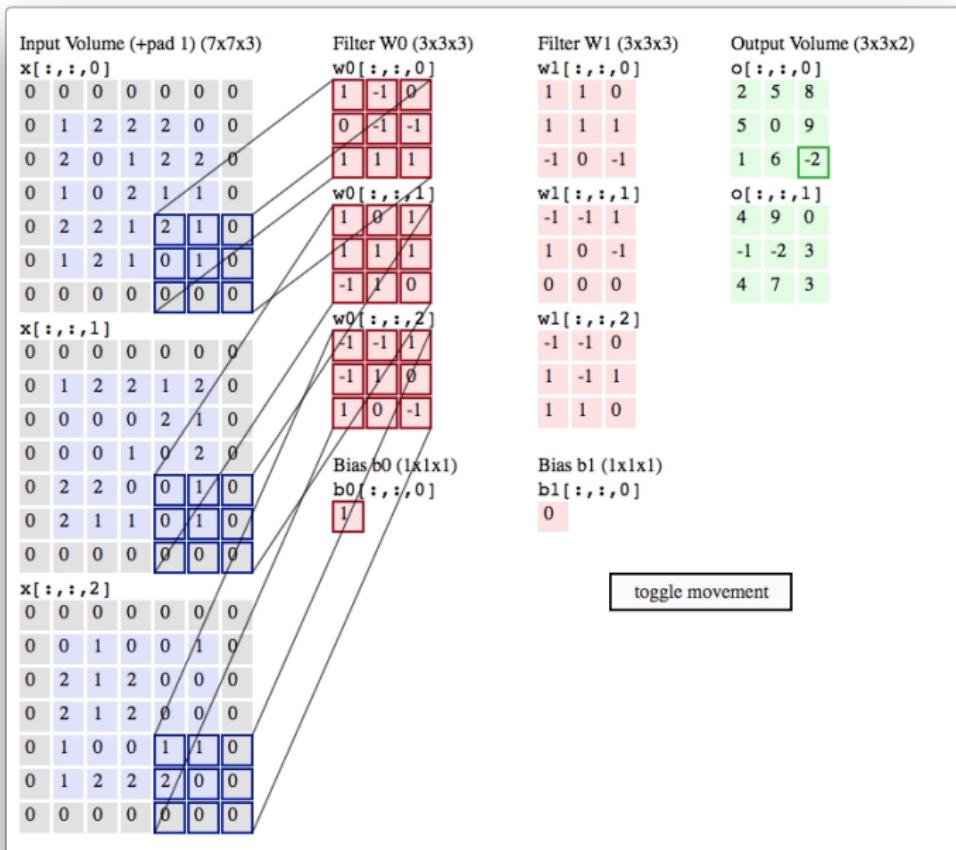


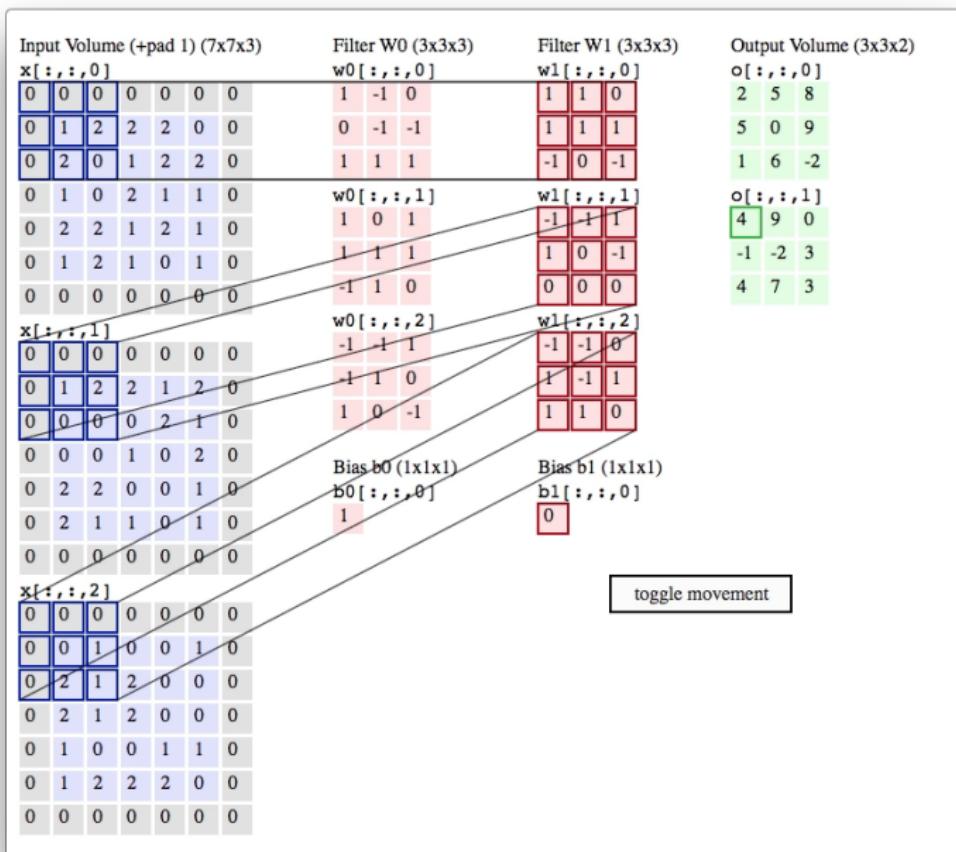


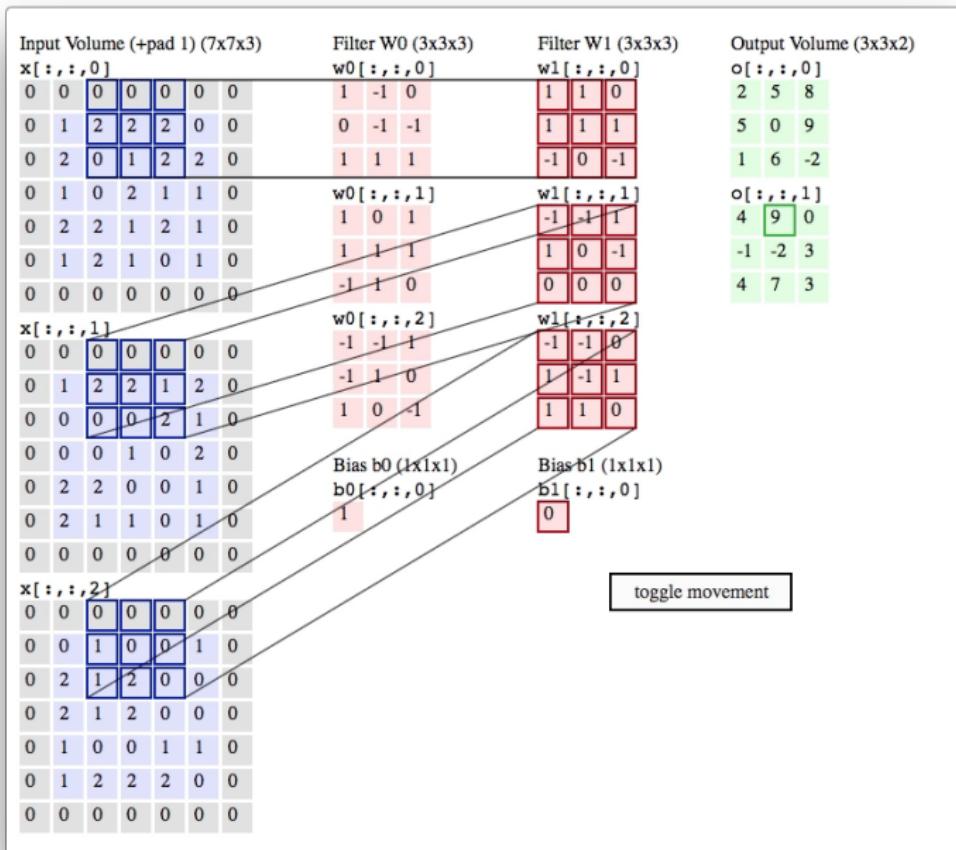


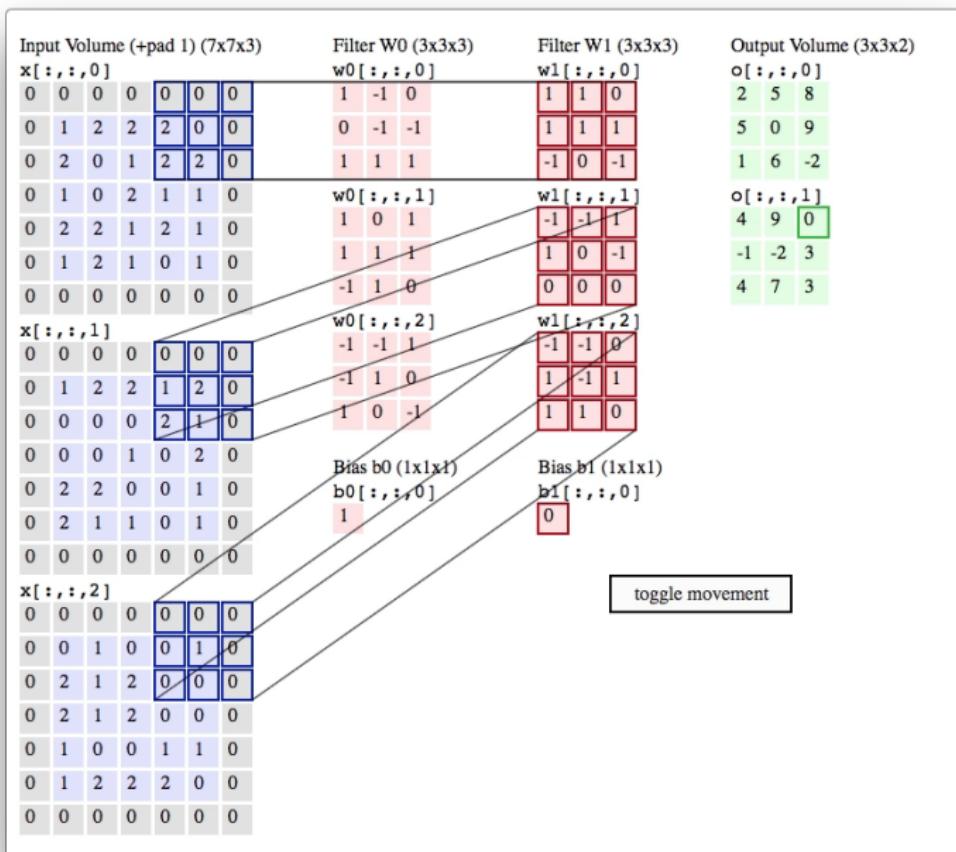


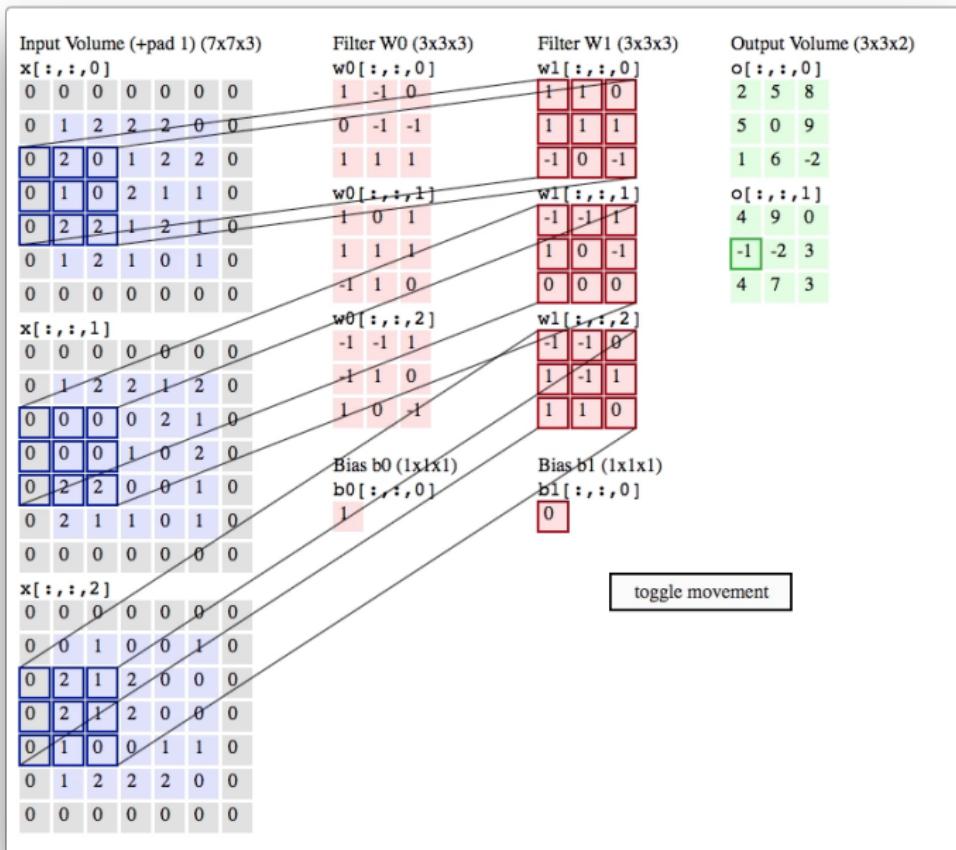


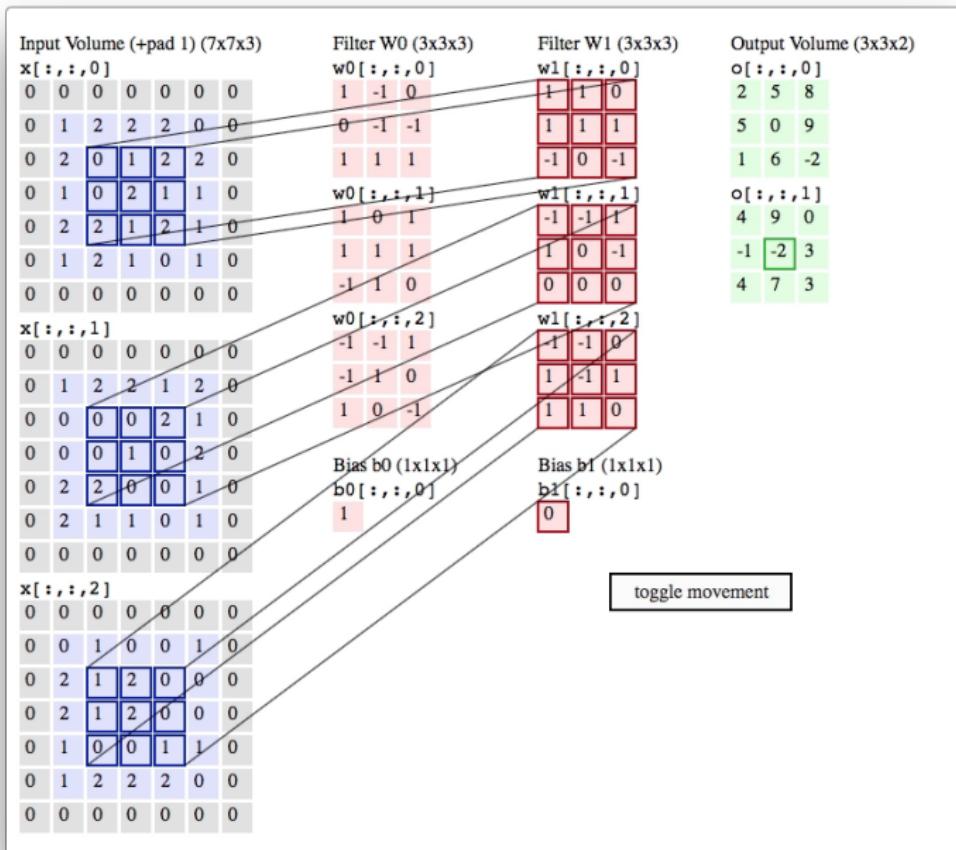


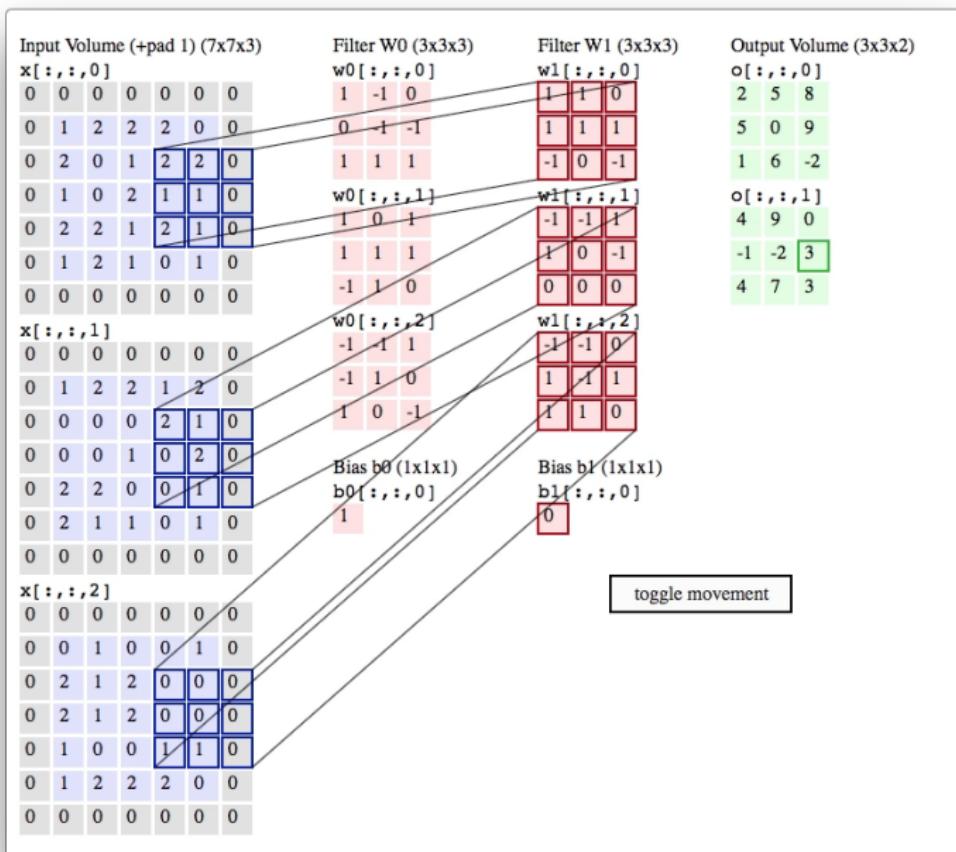


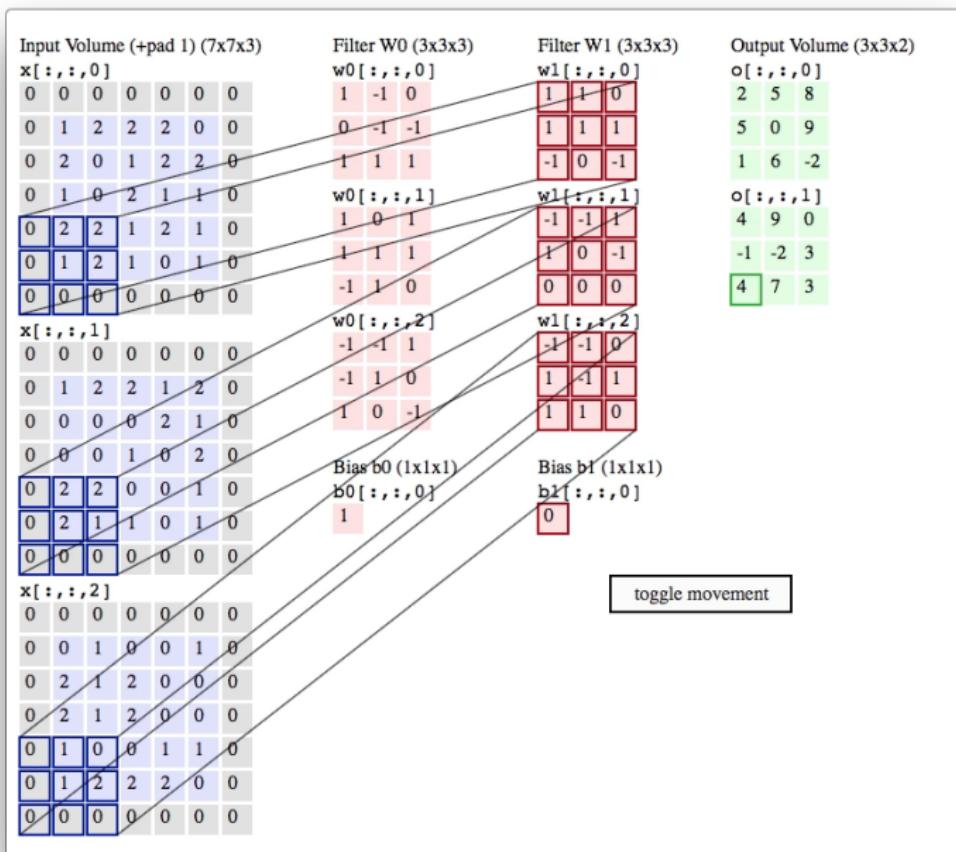


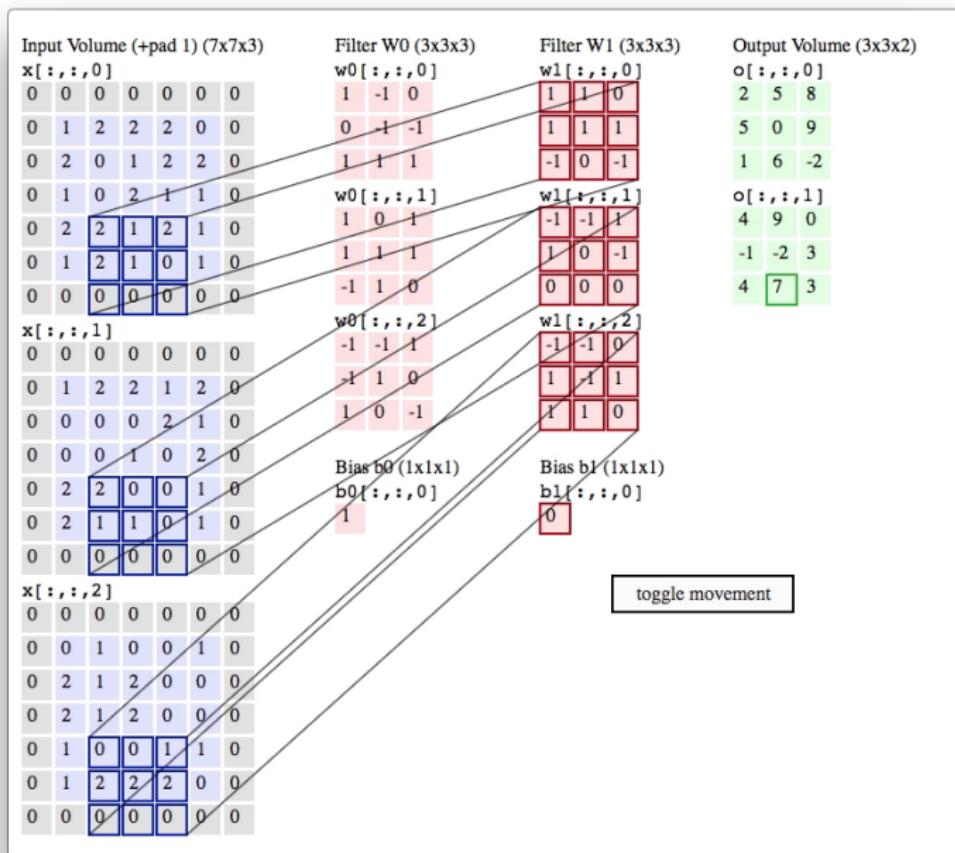


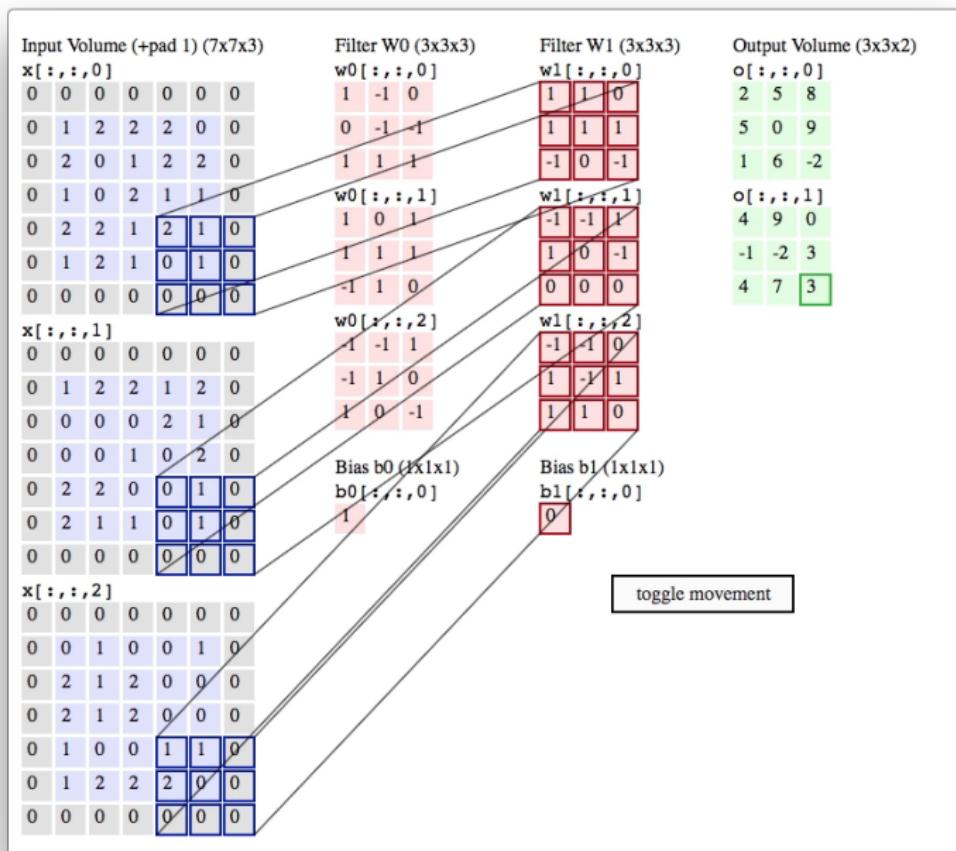












Mutual Constraints of Hyper-Parameters

- ➊ If stride is one, $S = 1$, it's common to set padding, $P = \frac{F-1}{2}$, to ensure that input and output volumes have the same size.
- ➋ The number of neurons that fit is $\frac{W-F+2P}{S} + 1$. Must set hyper-parameters so there number of neurons is an integer!
- ➌ **Real-World Example:** (Winner of ImageNet 2012) Images $227 \times 227 \times 3$, $F = 11$, $S = 4$, $P = 0$, $K = 96$.

- Number of neurons in *first layer*

$$\frac{227-11}{4} + 1 = 55 \text{ (per width)}$$

$$55 \times 55 \times 96 = 290400 \text{ (width} \times \text{height} \times \text{depth})$$

Each neuron has $11 \times 11 \times 3 = 363$ weights plus 1 bias

$290400 \times 364 = 105705600$ parameters *in the first layer!*

Parameter Sharing

If one feature (filter/kernel/depth slice) is useful at pixel (x_1, y_1) , then it is also useful at pixel (x_2, y_2)

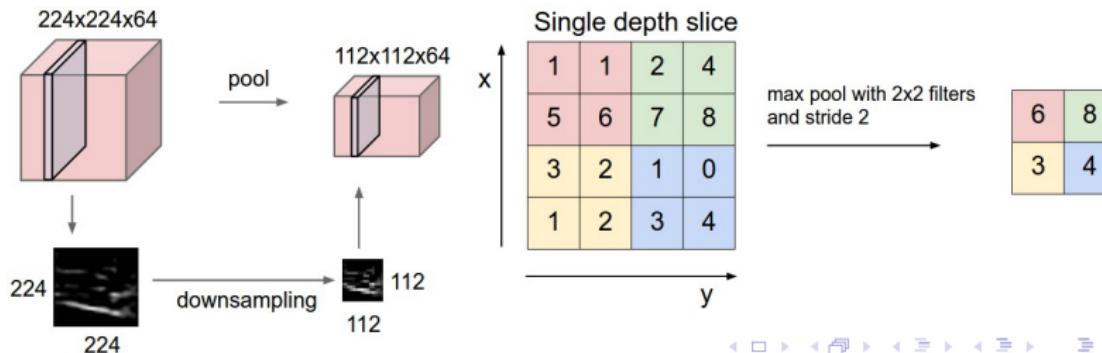
- ① All neurons in the same depth slice have same weights and bias.
- ② Example before $96 \times 11 \times 11 \times 3 = 34848$ unique weights + 96 biases = 34944 parameters
- ③ Then forward pass can be seen as the convolution of that slice neuron's weights with the input volume.
- ④ Weights = filter/kernel that is convolved with the input



Example filters learned by Krizhevsky et al. Each of the 96 filters shown here is of size $11 \times 11 \times 3$

Pooling Layer

- ① Inserted in-between convolutional layers
- ② Goal: To reduce spatial size, amount of parameters, computation requirements, and overfitting
- ③ Operates independently on every depth slice and resizes spatially using $\max(\cdot)$ (or could be $\text{mean}(\cdot)$ or $L2$)
- ④ Introduces no parameters, since it is a fixed operation
- ⑤ Hyper-parameters: F (spatial size) and S (stride)



Normalization Layer

Future - Hope of Getting Rid of Pooling

- ① Striving for Simplicity: Newer works propose discarding pooling. Suggest using larger strides
- ② Variational auto-encoders
- ③ Generative Adversarial Networks

Normalization Layer

- ① Have shown no real advantages

Practice

Hyper-Parameters and Layer Sizing

- ➊ Input layer should be divisible by 2 many times
- ➋ Conv layer use small filters (3×3), (5×5), (notice odd) and stride $S = 1$ (reduce sizing headaches - downsampling done by pooling)
- ➌ Use padding that doesn't alter dimension
- ➍ Pool layers (2×2), $F = 2$, and $S = 2$. Discards 75% of activations
- ➎ Compromise based on memory constraints (GPU). May have to sacrifice at the 1st layer.

Popular Networks

- ➊ LeNet: LeCun 90's - read zip code (first successful ConvNet)
- ➋ AlexNet: Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton. Popularize ConvNets for computer vision. ImageNet. Deeper, bigger than LeNet
- ➌ ZF Net: Matthew Zeiler and Rob Fergus. ILSVRC 2013 Winner. Improved by expanding the size of the middle convolutional layers making the stride and filter size on the first layer smaller
- ➍ GoogLeNet: Szegedy et al. from Google. ILSVRC 2014 winner. Reduced the number of parameters in the network from 60M to 4M (inception model)
- ➎ ResNet: Residual Network. Kaiming He et al. ILSVRC 2015 winner. It features special skip connections and heavy use of batch normalization. State of the art as of 2016

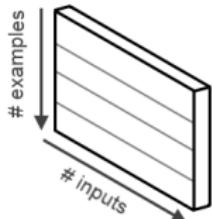
Recurrent Neural Networks (RNNs)

RNNs are a type of neural network that share parameters across time. They are used frequently for natural language processing (NLP) and time series data.

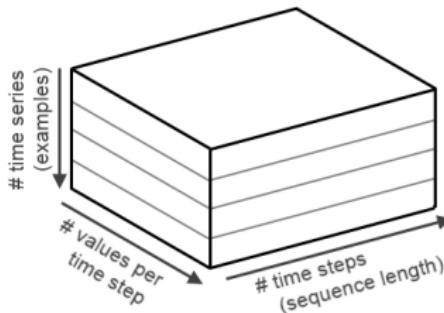
RNNs

- 1 Unlike standard neural networks or CNNs, RNNs introduce a temporal or sequential component to the system.
- 2 While CNNs introduce idea of spatial dependence, RNNs rely on the notion that *the current input relies on previous inputs*

Feed Forward Network Data

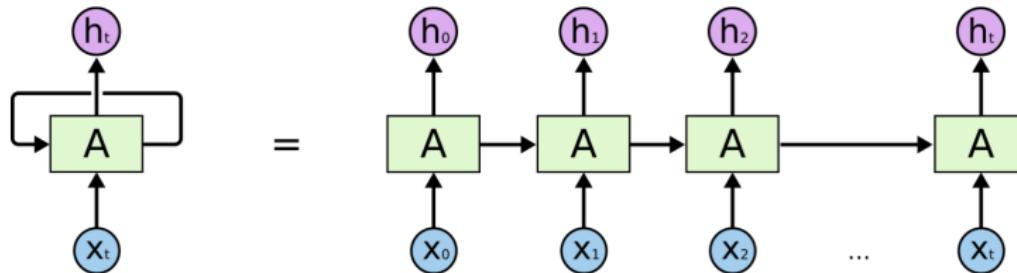


Recurrent Network Data



RNN Architecture and Notation

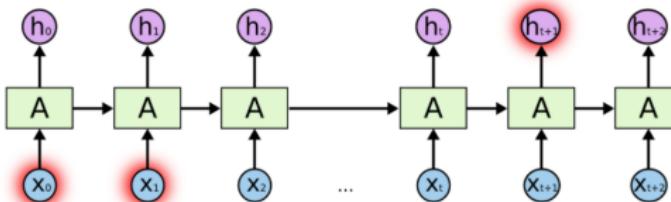
- ① Current input: x_t
- ② Current hidden state: h_t
- ③ Non-linear activation function σ or $tanh$
- ④ Predicted output: \hat{y}_t



Long-Term Dependencies

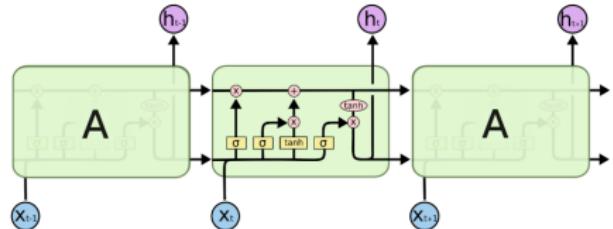
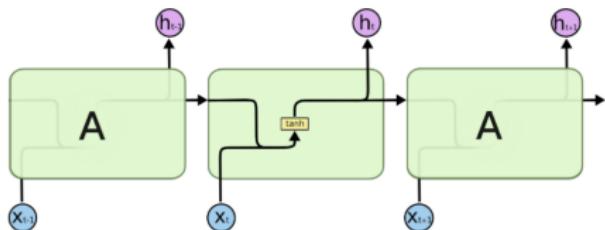
- ① Suppose we want a model that predicts the next word in this sentence:
 - “*The clouds are in the <WORD>*”

- ② Now suppose we want to model the next word in this sentence:
 - “*I grew up in France and lived there most of my life. As a result I am a native speaker of <WORD>*”



Long Short Term Memory Networks (LSTM)

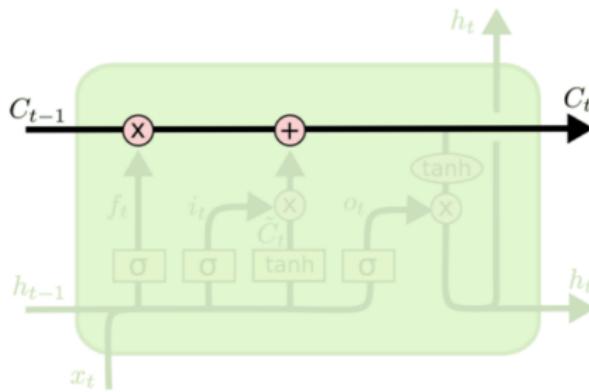
LSTMs are a species of RNN that is designed for learning long-term dependencies. While standard RNNs are able to model temporal dependencies, LSTMs have a much more elaborate structure.



Cell State

The key idea of LSTMs is that of cell state to propagate information.

- ① Cell state is subject to linear operations
- ② Information in cell state can be “forgotten” using a system of gates

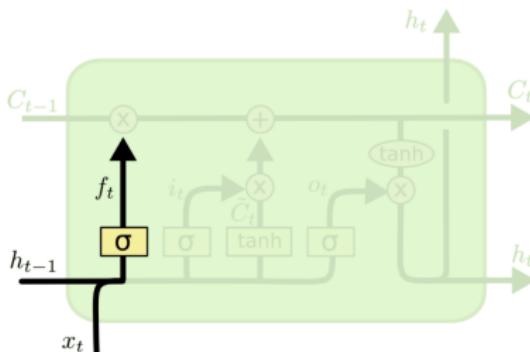


LSTM Forget Gate Layer

Forget gate layer is a sigmoid layer that decides the values to be updated with

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

- ➊ This layer looks at h_{t-1} and x_t and updates the cell state
- ➋ f_t has range $(0, 1)$
- ➌ 0 indicates “complete get rid of this”, and 1 indicates “completely keep this”.



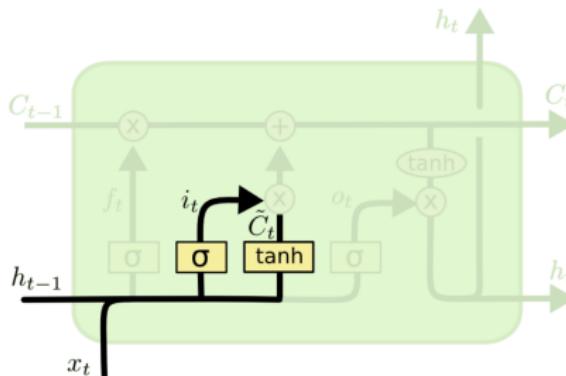
Input Gate Layer

Input gate layer is a combination of a sigmoid layer and $tanh$ layer that decide the new information to keep with

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

- ① Sigmoid layer is our “input gate layer”
- ② $tanh$ layer creates candidate values

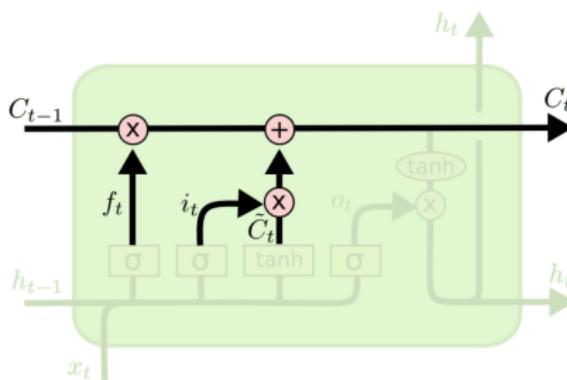


Update the Cell State

Now we update the cell state using

$$C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t$$

- ➊ Multiply old state by f_t , forgetting what we had decided
- ➋ And add $i_t \cdot \tilde{C}_t$, our new candidate values



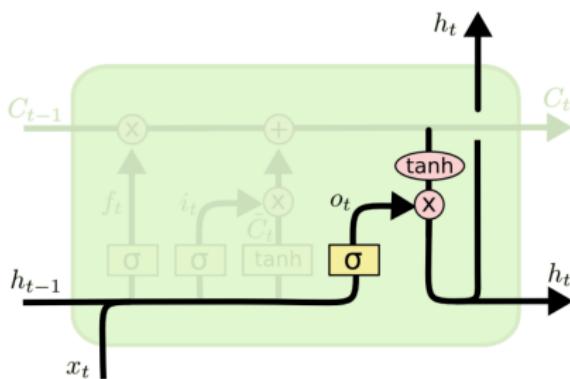
Output Layer

Now we decide what to output with

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \cdot \tanh(C_t)$$

- ① Take updated cell state to help compute output
- ② Apply sigmoid layer and \tanh layer



Autoencoder Basics

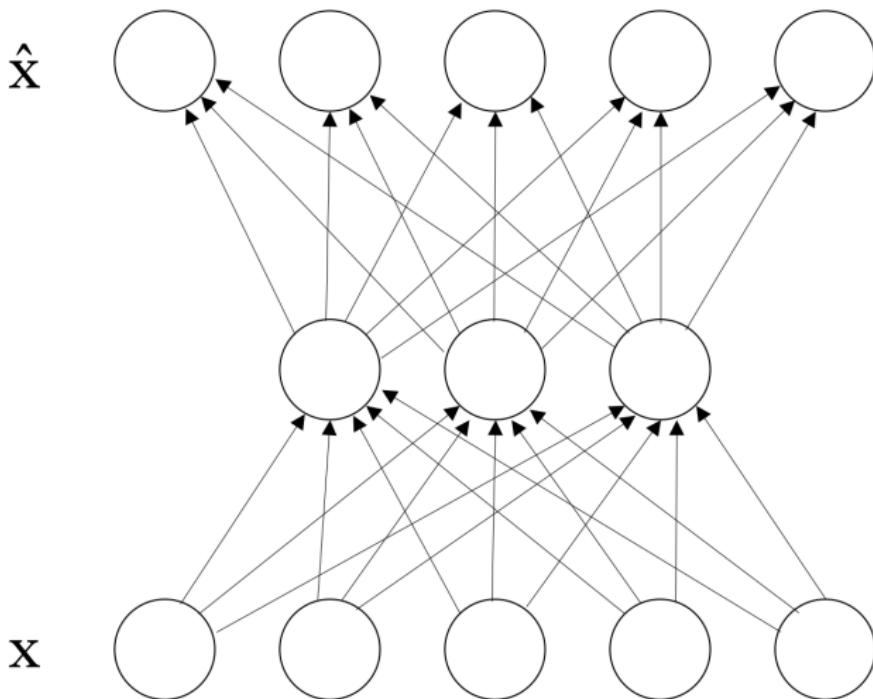
- ① Feed-forward neural network
- ② Used for unsupervised learning
- ③ Trained to reproduce input layer in the output layer
- ④ Training uses gradient descent

Autoencoder Basics cont.

- ➊ Similar to PCA*, but more flexible
- ➋ Autoencoders can accommodate non-linear transformations
 - 1 Consequence of flexible activation functions
 - 2 Turns out to be hugely useful

* Actually, with squared error loss and no sigmoid transformation, it is equivalent to PCA (Baldi & Hornik, 1989)

Simple Autoencoder



Simple Autoencoder cont.

From input layer to hidden layer is “encoder”:

$$f(x) = \sigma(Wx + b)$$

Hidden layer to the output is the “decoder”:

$$g(x) = \sigma(W'x + b')$$

where σ is the activation function (often sigmoid) and W' is often W^\top (known as having “tied” weights).

Simple Autoencoder cont.

- ① Input layer → hidden layer (fewer units) → output layer
- ② Hidden layers values can be thought of as lossy compression of input layer
- ③ Hidden layer with fewer units than input is often called “bottleneck” or under-complete layer

Question

In a sense, autoencoders want to learn an approximation to the identity function.

Why would this be useful?

Answer

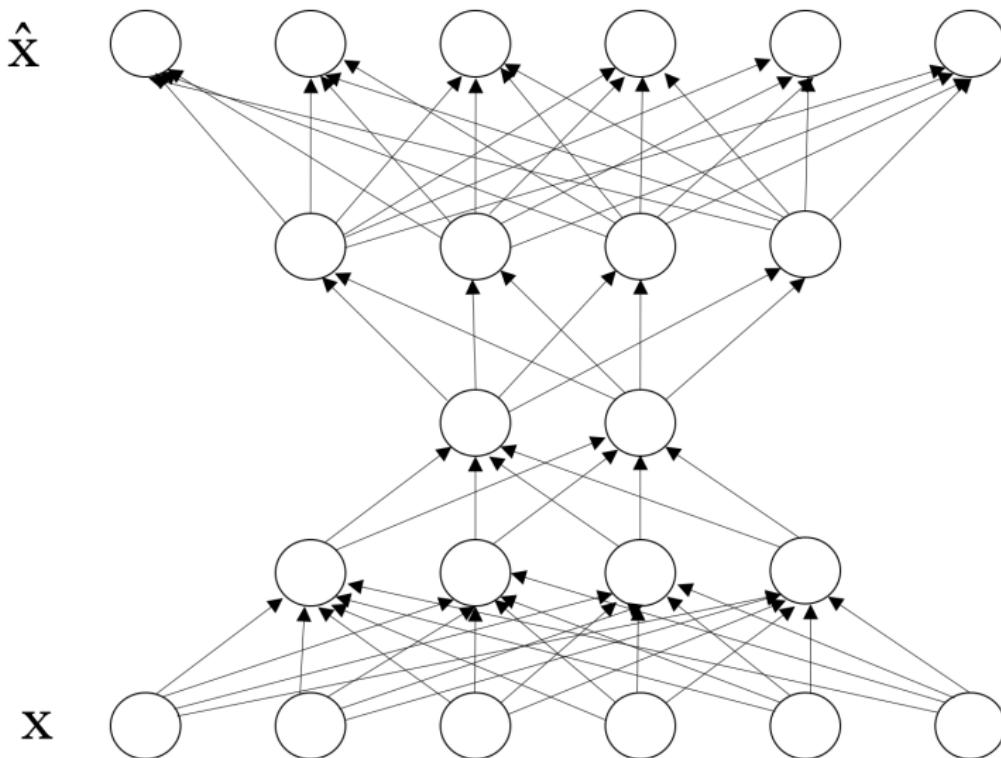
This turns out to have several practical uses.

- ➊ Autoencoders learn a kind of latent representation of the input
- ➋ Can be used for initializing weights and biases prior to fitting neural net
- ➌ Very well suited to dimensionality reduction in NLP tasks

Stacked Autoencoder

- ➊ Sometimes called deep autoencoder
- ➋ Feed-forward neural network
- ➌ Has additional layers
- ➍ Still “unsupervised” in the sense of no labels or real-valued outcome variable
- ➎ Output layer is still \hat{x}
- ➏ Trained in stages

Stacked Autoencoder



Stacked Autoencoder cont.

- ① Beneficial for data with hierarchical organization
- ② Can be difficult to train using back propagation

Denoising Autoencoder

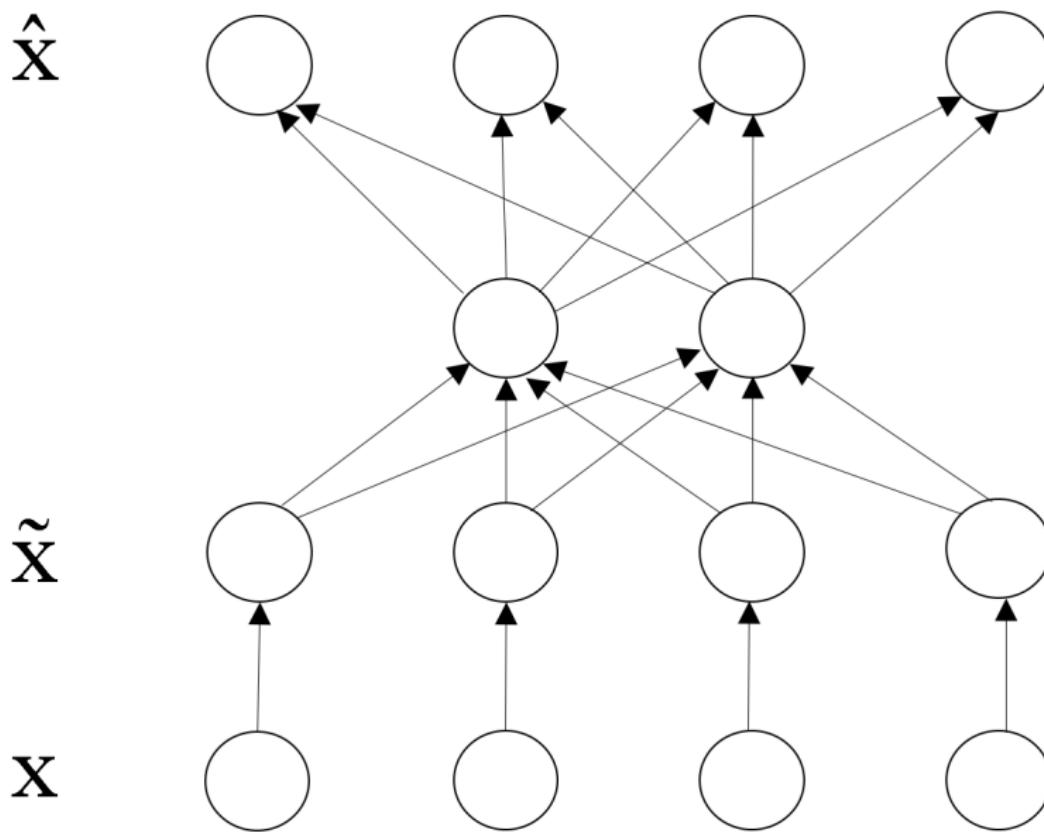
- ➊ Stochastic version of autoencoder
- ➋ Random noise injected into the input layer
 - 1 *Masking noise*: some fraction of elements* in x set to 0
 - 2 *Gaussian additive noise*: $\tilde{x}|x \sim \mathcal{N}(x, \sigma^2 I)$
 - 3 *Salt-and-pepper noise*: fraction of elements* in x set to either min or max possible value (often 0 or 1) according to coin toss

* chosen at random for each example

Denoising Autoencoder

- ➊ Output layer is evaluated against original (uncorrupted) data
- ➋ Goal is slightly different from simple autoencoder
 - 1 Not attempting to learn identity function
 - 2 We want robustness of representation
 - 3 Learned representation is insensitive to perturbations in the input

Denoising Autoencoder



Denoising Autoencoder cont.

- ① Training proceeds fairly similarly to standard autoencoder
- ② Crucial difference is the loss function is calculated using *uncorrupted* input layer

Extending Denoising Autoencoder

- ➊ Some noise types only corrupt subset of the input's components
 - 1 Masking noise
 - 2 Salt-and-pepper noise
- ➋ For these, we can extend denoising by adding emphasis on corrupted components.
- ➌ Use hyperparameters α and β to control emphasis; for squared error loss, this yields

$$L(x, \hat{x}) = \alpha \left(\sum_{i \in \mathcal{I}(\tilde{x})} (x_i - \hat{x}_i) \right) + \beta \left(\sum_{i \notin \mathcal{I}(\tilde{x})} (x_i - \hat{x}_i) \right)$$

where $\mathcal{I}(\tilde{x})$ denotes indexes components of x that were corrupted.

When to use neural networks

Conditions under which you might consider using neural networks:

- ① Have a huge amount of labeled training data
- ② Image classification (with huge amount of labeled images)
- ③ Certain NLP tasks
- ④ Some signal processing problems

When *not* to use neural networks

Probably should **not** use neural networks when:

- ① You have specific hypotheses you want to test
 - E.g., "*Drug X improves condition Y*".
- ② Interested in estimating the effect of some variable(s) on some outcome variable
- ③ You have highly structured data and/or few features

In the case of (1.) and (2.), a traditional statistical model is better. In the case of (3.), using some ensemble-of-trees method will give as-good or better results with minimal tuning.

Sources

- 1 <http://cs231n.github.io/convolutional-networks/>
- 2 <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- 3 Dave Berenbaum "Activation Functions" (2016)
https://github.com/brown-data-science/deep_learning_group/blob/master/slides/activation_functions.pdf
- 4 Yinong Wang "Recurrent Neural Networks" (2016)
https://github.com/brown-data-science/deep_learning_group/blob/master/slides/recurrent_neural_networks.pdf
- 5 Vincent P., Larochelle, H., Bengio Y., Manzagol P.A. (2008)
Extracting and Composing Robust Features with Denoising Autoencoders.
- 6 Vincent P., Larochelle, H., LaJoie I., Bengio Y., Manzagol P.A. (2010) *Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion*
- 7 Makhzani A., Frey B. (2015) *Winner-Take-All Autoencoders*
- 8 Ng A. *Sparse Autoencoders CS294A Lecture notes*