

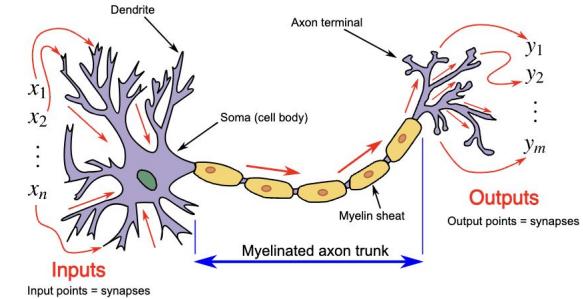
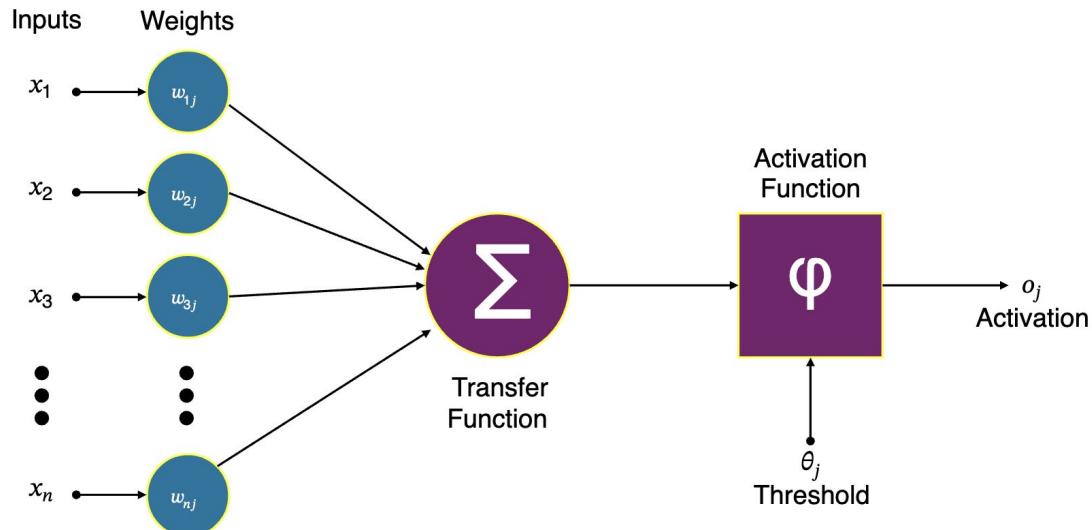
Week 8

Deep Learning

Content

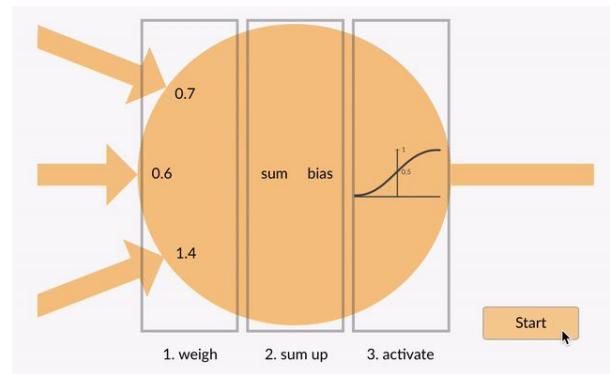
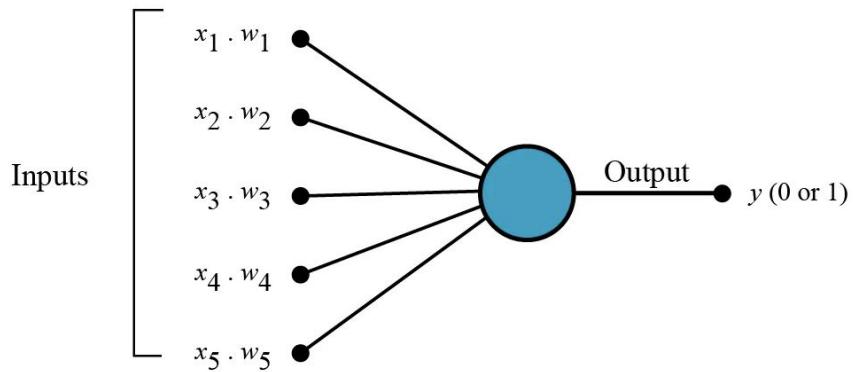
- Feedforward Neural Networks (FFNN)
- Back-Propagation
- Convolutional Neural Network (CNN)
- Training ConvNets
- Example: CNN in Python

Recap: Artificial neuron structure



Biological neuron model

Recap: A single Neuron from an Artificial Neural Network (ANN)

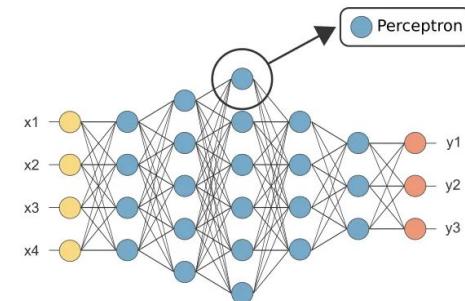
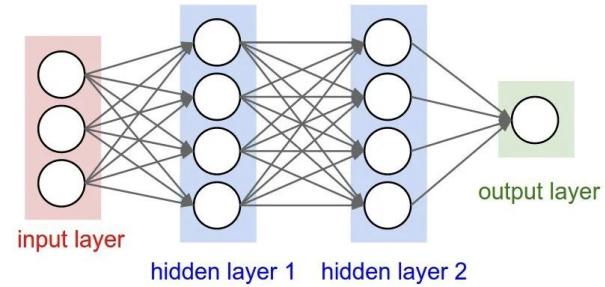


One of the first ANNs was known as the perceptron, and it consisted of only a single neuron.

Feedforward Neural Networks (FFNN)

A fully-connected feed-forward neural network (FFNN) — aka A multi-layered perceptron (MLP)

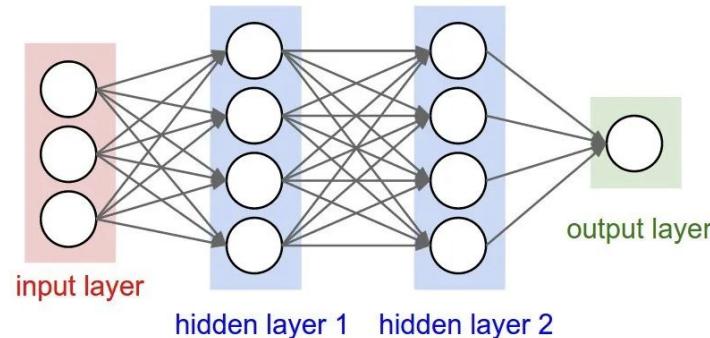
- A type of artificial neural network where connections between nodes do not form a cycle. They are called “feedforward” because data flows in one direction—from input to output.
- Consists of an input layer, one or more hidden layers, and an output layer.
- Key Components
 - The **Perceptron**: one of the first Neural Networks, is made of just a single Neuron
 - **Neurons**: Basic units of the network, performing computations.
 - **Layers**:
 - **Input Layer**: Receives the initial data.
 - **Hidden Layers**: Intermediate layers that perform computations and feature extraction.
 - **Output Layer**: Produces the final output.



Feedforward Neural Networks (FFNN)

A network of artificial neurons with multiple layers

- A type of artificial neural network where connections between nodes do not form a cycle. They are called “feedforward” because data flows in one direction—from input to output.
- Consists of an input layer, one or more hidden layers, and an output layer.
- Key Components
 - Neurons: Basic units of the network, performing computations.
 - Layers:
 - Input Layer: Receives the initial data.
 - Hidden Layers: Intermediate layers that perform computations and feature extraction.
 - Output Layer: Produces the final output.
- Activation Functions
 - Purpose: Introduce non-linearity into the model, enabling complex functions.
 - Common Types:
 - Sigmoid
 - ReLU (Rectified Linear Unit)
 - Tanh



Feedforward Neural Networks (FFNN)

Example: Learning XOR

- XOR Function (“Exclusive or”)

x_1	x_2	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

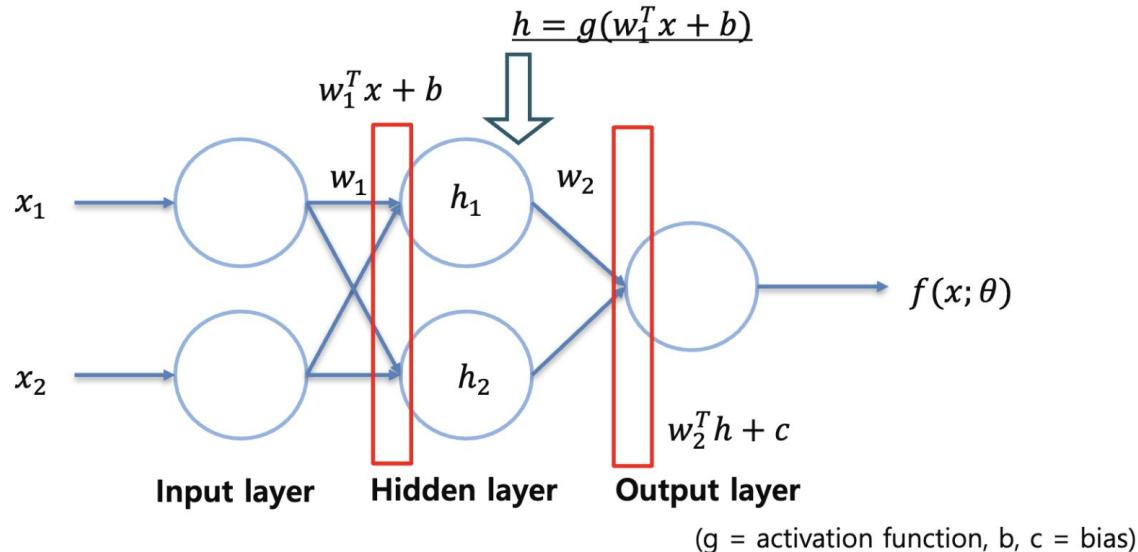
- $X = \{[0, 0]^T, [0, 1]^T, [1, 0]^T, [1, 1]^T\}$
- Mean squared error (MSE) loss function

$$J(\boldsymbol{\theta}) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \boldsymbol{\theta}))^2$$

- $f^*(x)$: correct answer
- $f(x; \theta)$: calculated result by neural network

Feedforward Neural Networks (FFNN)

Example: Learning XOR



$$f(x; \theta) = f(x; w_1, b, w_2, c) = w_2^T g(w_1^T x + b) + c$$

Feedforward Neural Networks (FFNN)

Gradient-based Learning

- **Loss function of neural network is non-convex**
 - Trained by using iterative, gradient-based optimizers
- **Cost function**
 - Learning conditional distribution
 - Learning conditional statistics
- **Output layers**
 - Linear units for Guassian output distributions
 - Sigmoid units
 - Softmax units

Feedforward Neural Networks (FFNN)

Learning Conditional Distribution

- Negative log-likelihood as cross-entropy between training data and model distribution

$$J(\boldsymbol{\theta}) = H(\hat{p}_{data}, p_{model}(y|x)) = - \sum_{x,y \sim \hat{p}_{data}} \hat{p}_{data} \log p_{model}(y|x)$$

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{x,y \sim \hat{p}_{data}} \log p_{model}(y|x)$$

- \hat{p}_{data} : *training data – generating distribution*
- $p_{model}(y|x)$: *probability distribution estimating* \hat{p}_{data}
- $H(\hat{p}_{data}, p_{model}(y|x))$: *cross entropy between* \hat{p}_{data} *and* p_{model}
- \log function undoes the \exp of some output units

Feedforward Neural Networks (FFNN)

Learning Conditional Distribution

- Mean Square Error (MSE)

$$f^* = \arg \min_f \mathbb{E}_{x,y \sim p_{data}} \|y - f(x)\|^2$$

- Mean Absolute Error (MAE)

$$f^* = \arg \min_f \mathbb{E}_{x,y \sim p_{data}} \|y - f(x)\|_1$$

- MSE, MAE often lead to poor results when used with gradient-based learning

Feedforward Neural Networks (FFNN)

Linear Units for Gaussian Output Distributions

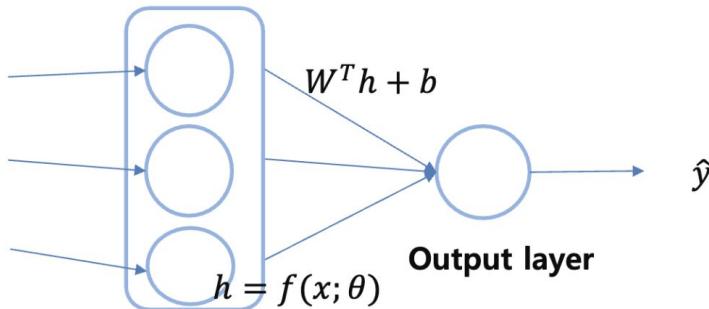
- Feature h , produce a vector \hat{y}

$$\hat{y} = W^T h + b$$

- They produce the mean of a conditional Gaussian distribution

$$p(y|x) = \mathcal{N}(y; \hat{y}, I)$$

- No saturation \rightarrow little difficult to gradient-based optimization



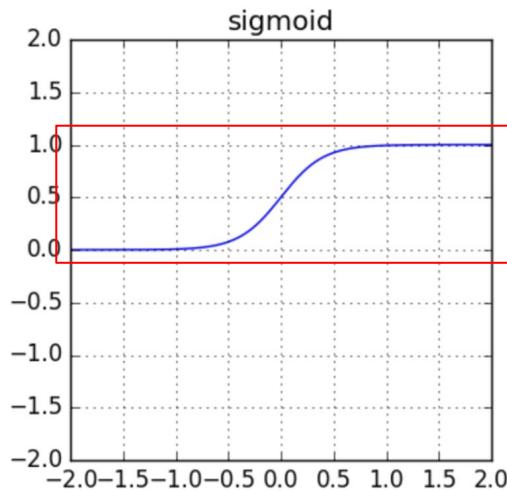
Feedforward Neural Networks (FFNN)

Sigmoid Units (Two plots, “soft” meanings)

- Binary cl $\hat{y} = \sigma(w^T h + b)$
- Output

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

- Saturate to 0 and 1



Feedforward Neural Networks (FFNN)

Softmax Units

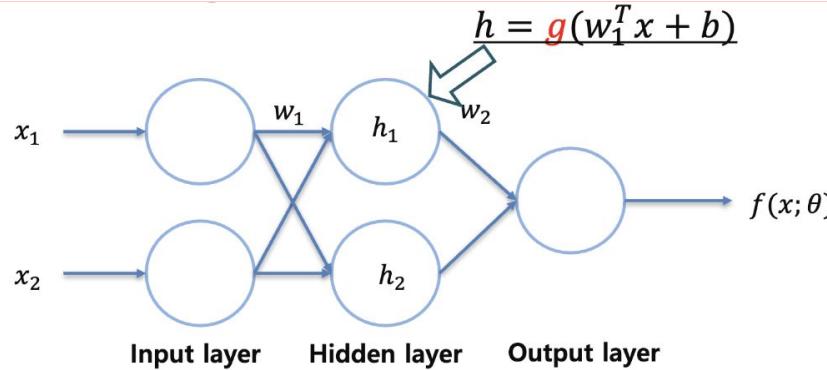
- Multiclass classification
- To generalize to the case of a discrete variable with n values, vector \hat{y} , with $\hat{y}_i = P(y = i|x)$

$$\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b} \quad (z_i = \log \tilde{P}(y = i|x))$$
$$softmax(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

Feedforward Neural Networks (FFNN)

Hidden Units (Which active functions always use)

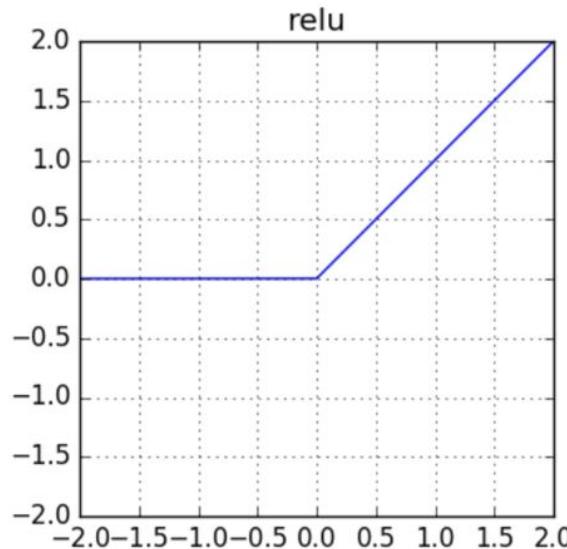
- How to choose the type of hidden unit to use in the hidden layers of the model
- Input: $\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$
- Activation function $\mathbf{g}(\mathbf{z})$
 - Rectified Linear Units (ReLU)
 - Logistic Sigmoid and Hyperbolic Tangent



Feedforward Neural Networks (FFNN)

Rectified Linear Units (ReLU)

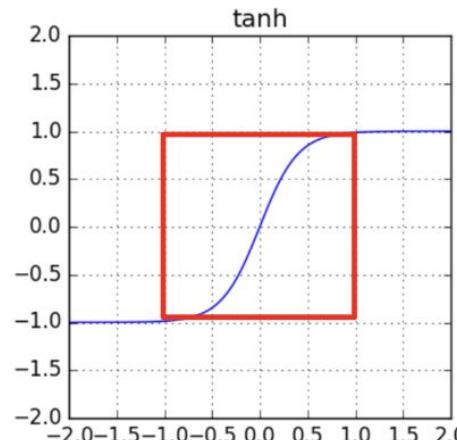
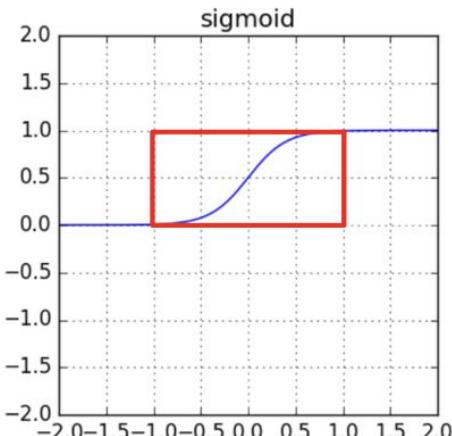
- Activation function $g(z) = \max(0, z)$
- If model's behavior is closer to linear, models are easier to optimize



Feedforward Neural Networks (FFNN)

Logistic Sigmoid and Hyperbolic Tangent

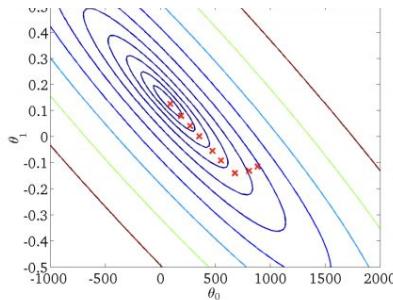
- Activation function
 - $g(z) = \sigma(z)$
 - $g(z) - \tanh(z) = 2\sigma(2z) - 1$
- Saturation of sigmoidal units make gradient-based learning difficult



Back-Propagation

Recap: Gradient Descent

- Recall: gradient descent moves opposite the gradient (the direction of steepest descent)



- Weight space for a multilayer neural net: one coordinate for each weight or bias of the network, in all the layers
- Conceptually, not any different from what we've seen so far — just higher dimensional and harder to visualize!
- We want to compute the cost gradient dE/dw , which is the vector of partial derivatives.
 - This is the average of dL/dw over all the training examples, so in this lecture we focus on computing dL/dw .

Back-Propagation

Recap: Univariate Chain Rule (Add figure of “Movement”)

- We've already been using the univariate Chain Rule.
- Recall: if $f(x)$ and $x(t)$ are univariate functions, then

$$\frac{d}{dt} f(x(t)) = \frac{df}{dx} \frac{dx}{dt}$$

Back-Propagation

Recap: Univariate Chain Rule

- Recall: Univariate logistic least squares model

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

- Let's compute the loss derivatives.

Back-Propagation

Recap: Univariate Chain Rule

$$\begin{aligned}\mathcal{L} &= \frac{1}{2}(\sigma(wx + b) - t)^2 \\ \frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial}{\partial w} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial w} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) x\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial}{\partial b} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial b} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial b} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b)\end{aligned}$$

What are the disadvantages of this approach?

Back-Propagation

Recap: Univariate Chain Rule

- A more structured way to do it

Computing the loss

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

Computing the derivatives

$$\frac{d\mathcal{L}}{dy} = y - t$$

$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{dy} \sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{d\mathcal{L}}{dz} \times$$

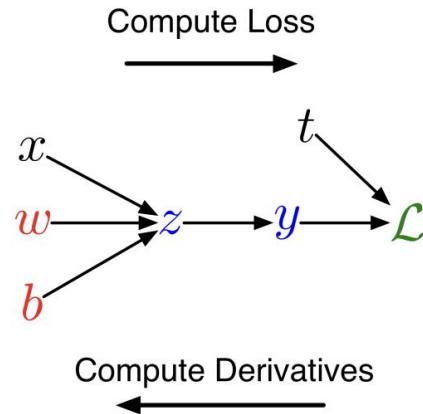
$$\frac{\partial \mathcal{L}}{\partial b} = \frac{d\mathcal{L}}{dz}$$

Remember, the goal isn't to obtain closed-form solutions, but to be able to write a program that efficiently computes the derivatives.

Back-Propagation

Recap: Univariate Chain Rule

- We can diagram out the computations using a computation graph.
- The nodes represent all the inputs and computed quantities, and the edges represent which nodes are computed directly as a function of which other nodes.



Back-Propagation

Recap: Univariate Chain Rule

- A slightly more convenient notation:
 - Use \bar{y} to denote the derivative $d\mathcal{L}/dy$, sometimes called the error signal.
 - This emphasizes that the error signals are just values our program is computing (rather than a mathematical operation).
 - This is not a standard notation, but I couldn't find another one that I liked.

Computing the loss

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

Computing the derivatives

$$\bar{y} = y - t$$

$$\bar{z} = \bar{y} \sigma'(z)$$

$$\bar{w} = \bar{z} x$$

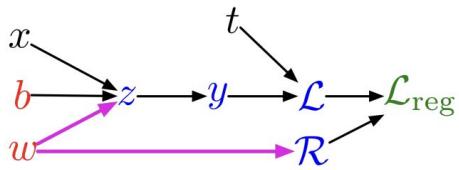
$$\bar{b} = \bar{z}$$

Back-Propagation

Recap: Multivariate Chain Rule

- Problem: what if the computation graph has fan-out > 1? This requires the multivariate Chain Rule!

L2-Regularized regression



$$z = wx + b$$

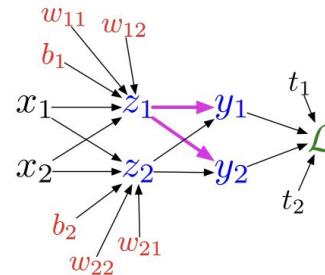
$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

$$\mathcal{R} = \frac{1}{2}w^2$$

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda\mathcal{R}$$

Multiclass logistic regression



$$z_\ell = \sum_j w_{\ell j} x_j + b_\ell$$

$$y_k = \frac{e^{z_k}}{\sum_\ell e^{z_\ell}}$$

$$\mathcal{L} = - \sum_k t_k \log y_k$$

Back-Propagation

Recap: Multivariate Chain Rule

- Suppose we have a function $f(x, y)$ and functions $x(t)$ and $y(t)$. (All the variables here are scalar-valued.) Then

$$\frac{d}{dt}f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

- Example:

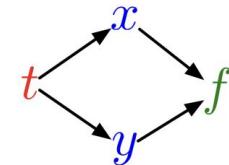
$$f(x, y) = y + e^{xy}$$

$$x(t) = \cos t$$

$$y(t) = t^2$$

- Plug in to Chain Rule:

$$\begin{aligned}\frac{df}{dt} &= \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} \\ &= (ye^{xy}) \cdot (-\sin t) + (1 + xe^{xy}) \cdot 2t\end{aligned}$$



Back-Propagation

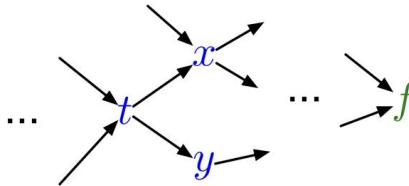
Recap: Multivariate Chain Rule

- In the context of backpropagation:

Mathematical expressions
to be evaluated

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Values already computed
by our program



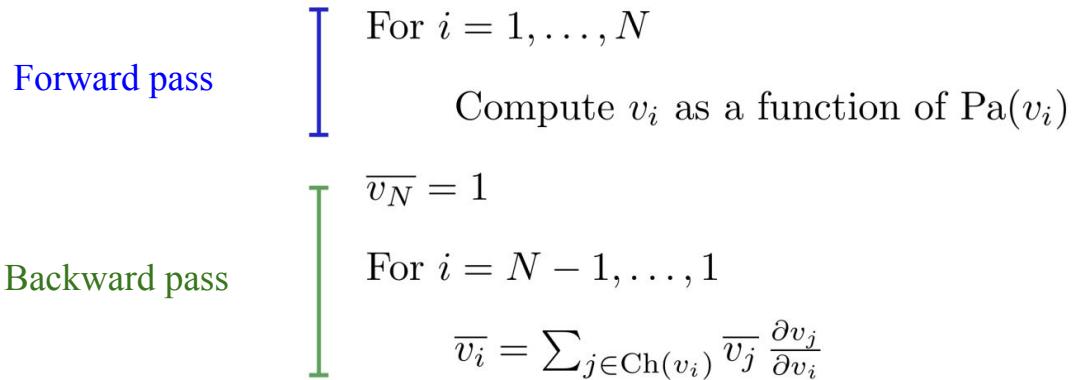
- In our notation:

$$\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$$

Back-Propagation

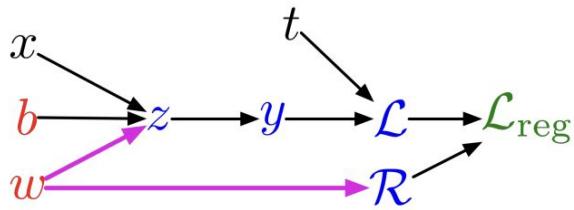
Full backpropagation algorithm

- Let v_1, \dots, v_N be a topological ordering of the computation graph (i.e. parents come before children.) v_N denotes the variable we're trying to compute derivatives of (e.g. loss).



Back-Propagation

Example: univariate logistic least squares regression



Forward pass:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

$$\mathcal{R} = \frac{1}{2}w^2$$

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda \mathcal{R}$$

Backward pass:

$$\overline{\mathcal{L}_{\text{reg}}} = 1$$

$$\begin{aligned}\overline{\mathcal{R}} &= \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{R}} \\ &= \overline{\mathcal{L}_{\text{reg}}} \lambda\end{aligned}$$

$$\begin{aligned}\overline{\mathcal{L}} &= \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{L}} \\ &= \overline{\mathcal{L}_{\text{reg}}}\end{aligned}$$

$$\begin{aligned}\overline{y} &= \overline{\mathcal{L}} \frac{d\mathcal{L}}{dy} \\ &= \overline{\mathcal{L}}(y - t)\end{aligned}$$

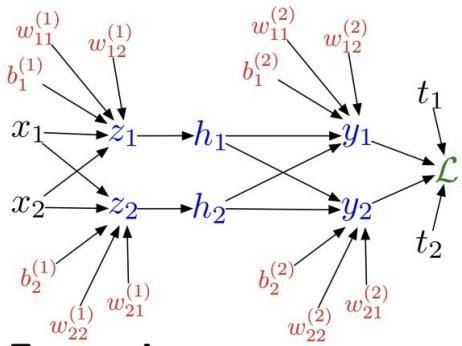
$$\begin{aligned}\overline{z} &= \overline{y} \frac{dy}{dz} \\ &= \overline{y} \sigma'(z)\end{aligned}$$

$$\begin{aligned}\overline{w} &= \overline{z} \frac{\partial z}{\partial w} + \overline{\mathcal{R}} \frac{d\mathcal{R}}{dw} \\ &= \overline{z}x + \overline{\mathcal{R}}w\end{aligned}$$

$$\begin{aligned}\overline{b} &= \overline{z} \frac{\partial z}{\partial b} \\ &= \overline{z}\end{aligned}$$

Back-Propagation

Multilayer Perceptron (multiple outputs):



Forward pass:

$$z_i = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

$$h_i = \sigma(z_i)$$

$$y_k = \sum_i w_{ki}^{(2)} h_i + b_k^{(2)}$$

$$\mathcal{L} = \frac{1}{2} \sum_k (y_k - t_k)^2$$

Backward pass:

$$\overline{\mathcal{L}} = 1$$

$$\overline{y_k} = \overline{\mathcal{L}} (y_k - t_k)$$

$$\overline{w_{ki}^{(2)}} = \overline{y_k} h_i$$

$$\overline{b_k^{(2)}} = \overline{y_k}$$

$$\overline{h_i} = \sum_k \overline{y_k} w_{ki}^{(2)}$$

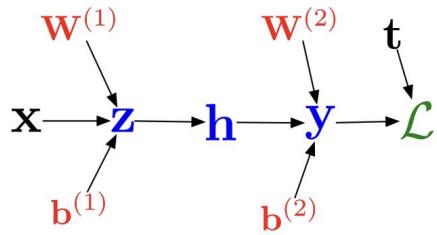
$$\overline{z_i} = \overline{h_i} \sigma'(z_i)$$

$$\overline{w_{ij}^{(1)}} = \overline{z_i} x_j$$

$$\overline{b_i^{(1)}} = \overline{z_i}$$

Back-Propagation

In vectorized form:



Forward pass:

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$

$$\mathbf{h} = \sigma(\mathbf{z})$$

$$\mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$\mathcal{L} = \frac{1}{2}\|\mathbf{t} - \mathbf{y}\|^2$$

Backward pass:

$$\bar{\mathcal{L}} = 1$$

$$\bar{\mathbf{y}} = \bar{\mathcal{L}}(\mathbf{y} - \mathbf{t})$$

$$\overline{\mathbf{W}^{(2)}} = \bar{\mathbf{y}}\mathbf{h}^\top$$

$$\overline{\mathbf{b}^{(2)}} = \bar{\mathbf{y}}$$

$$\bar{\mathbf{h}} = \mathbf{W}^{(2)\top}\bar{\mathbf{y}}$$

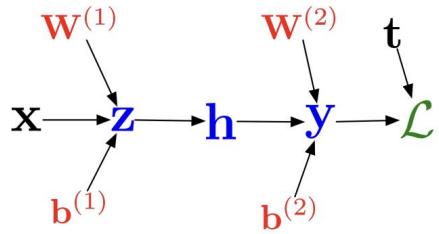
$$\bar{\mathbf{z}} = \bar{\mathbf{h}} \circ \sigma'(\mathbf{z})$$

$$\overline{\mathbf{W}^{(1)}} = \bar{\mathbf{z}}\mathbf{x}^\top$$

$$\overline{\mathbf{b}^{(1)}} = \bar{\mathbf{z}}$$

Back-Propagation

In vectorized form: (Add an Example, simple vector to calculate, explain it)



Forward pass:

$$z = W^{(1)}x + b^{(1)}$$

$$h = \sigma(z)$$

$$y = W^{(2)}h + b^{(2)}$$

$$\mathcal{L} = \frac{1}{2} \|t - y\|^2$$

Backward pass:

$$\bar{\mathcal{L}} = 1$$

$$\bar{y} = \bar{\mathcal{L}}(y - t)$$

$$\overline{W^{(2)}} = \bar{y}h^\top$$

$$\overline{b^{(2)}} = \bar{y}$$

$$\bar{h} = W^{(2)\top}\bar{y}$$

$$\bar{z} = \bar{h} \circ \sigma'(z)$$

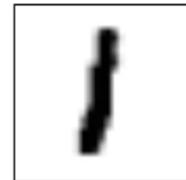
$$\overline{W^{(1)}} = \bar{z}x^\top$$

$$\overline{b^{(1)}} = \bar{z}$$

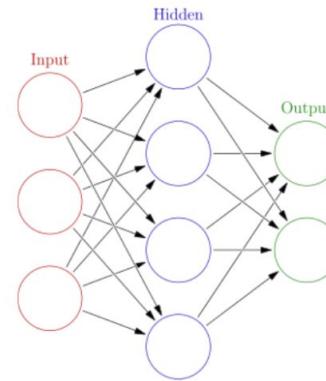
Convolutional Neural Network (CNN)

Introduction

- Multi-Layer Perceptron (MLP) for image recognition task
 - Use intensity of each pixel as 1D input vector
 - There several problems such as:
 - Required extremely large number of parameter
 - No invariance to shifting, scaling
 - Unable to deal with variable size input data



28x28 image



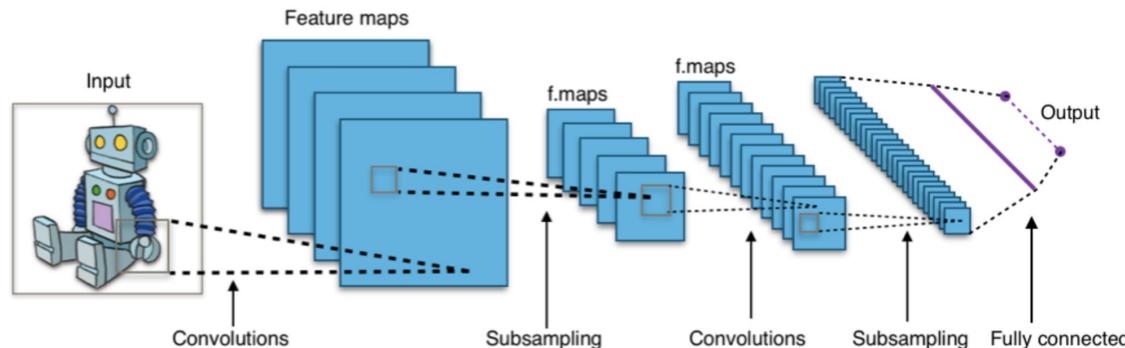
784 input neurons and large number of parameters

[https://en.wikipedia.org/wiki/Neural_network_\(machine_learning\)](https://en.wikipedia.org/wiki/Neural_network_(machine_learning))

Convolutional Neural Network (CNN)

Introduction

- The image contains information about Convolutional Neural Networks (CNNs). Here are the key points:
 - **CNN consists of:**
 - Convolutional layer
 - Pooling layer
 - Fully connected layer
 - **CNN learns the convolutional layer's kernel.**
 - To make computation easier, CNN uses cross-correlation instead of convolution.



Convolutional Neural Network (CNN)

ImageNet

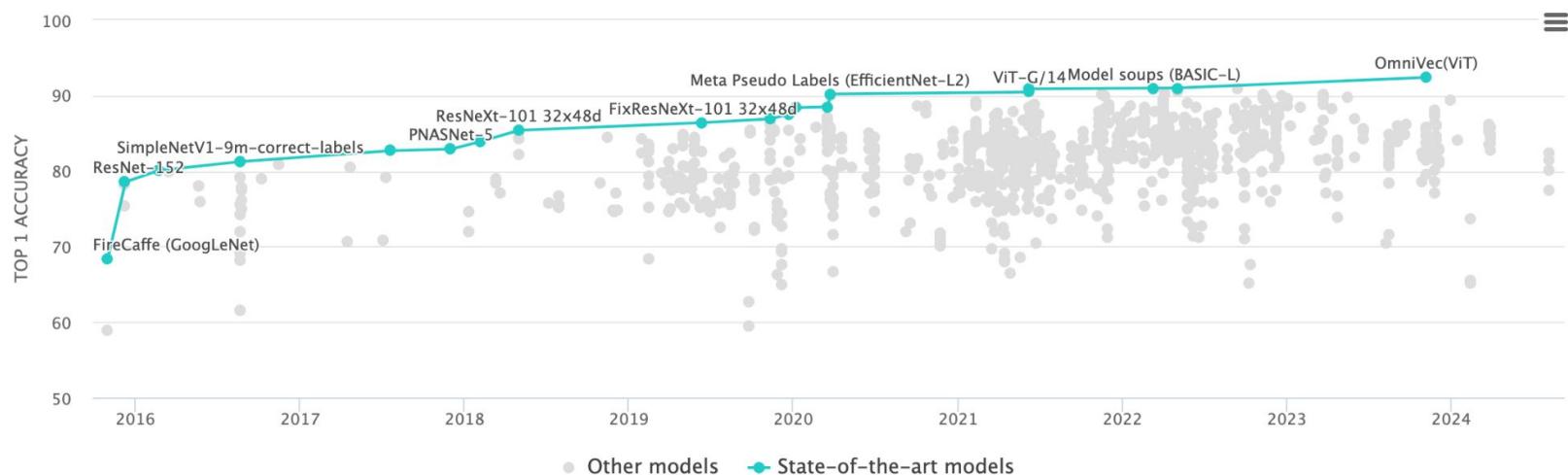
- ImageNet is a **large-scale** visual database designed for use in visual object recognition software research.
- It contains over **14 million images** that have been manually annotated to indicate the objects they depict.
- Organized based on the WordNet hierarchy, a lexical database of English.
- Synsets (synonym sets): Groupings of words with the same or similar meaning.
- Developed by Fei-Fei Li and her team at Stanford University in 2009.



Convolutional Neural Network (CNN)

ImageNet

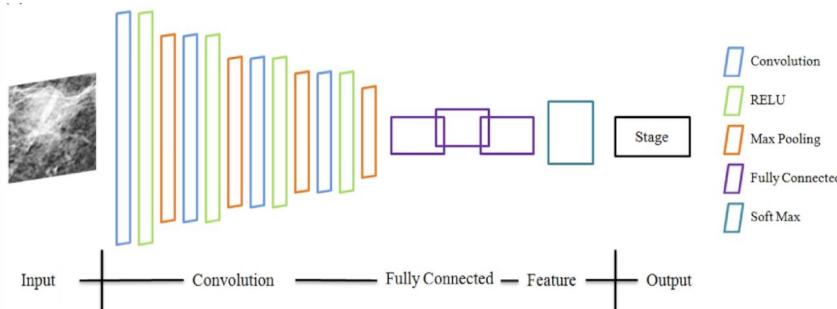
- Image Classification on ImageNet



Convolutional Neural Network (CNN)

CNN Architecture Overview

- **Image: Diagram of CNN layers**
 - CNNs consist of several **layer types**:
 - **Convolutional Layer**: Extracts features using filters/kernels.
 - **Pooling Layer**: Reduces the spatial dimensions of the feature maps.
 - **Fully Connected Layer**: Makes final predictions by connecting all neurons.
 - **Activation Functions**: Typically use **ReLU** to introduce non-linearity.



Convolutional Neural Network (CNN)

Convolutional Layer

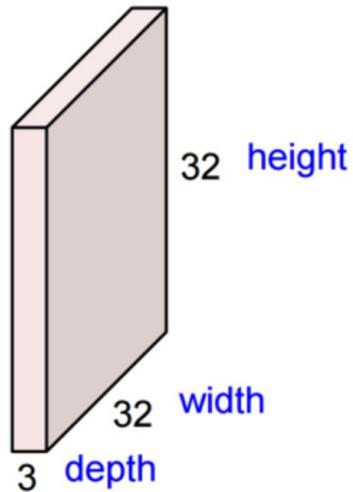
- **Image: Convolution process, showcasing filter moving across the image**
 - CNNs consist of several layer types:
 - **Convolution Operation:** Applies filters (small matrices) over the input image to detect features like edges, textures, and shapes.
 - Filters scan the image in sections (strides) and output **feature maps**.
 - Helps detect **low-level features** (e.g., edges) in early layers and **high-level features** (e.g., objects) in later layers.

Convolutional Neural Network (CNN)

Convolution Layer

- Example

32x32x3 image

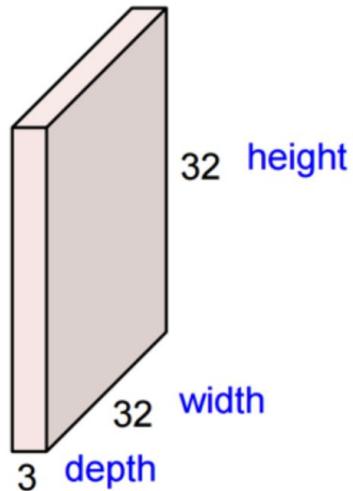


Convolutional Neural Network (CNN)

Convolution Layer

- Example

32x32x3 image



5x5x3 filter

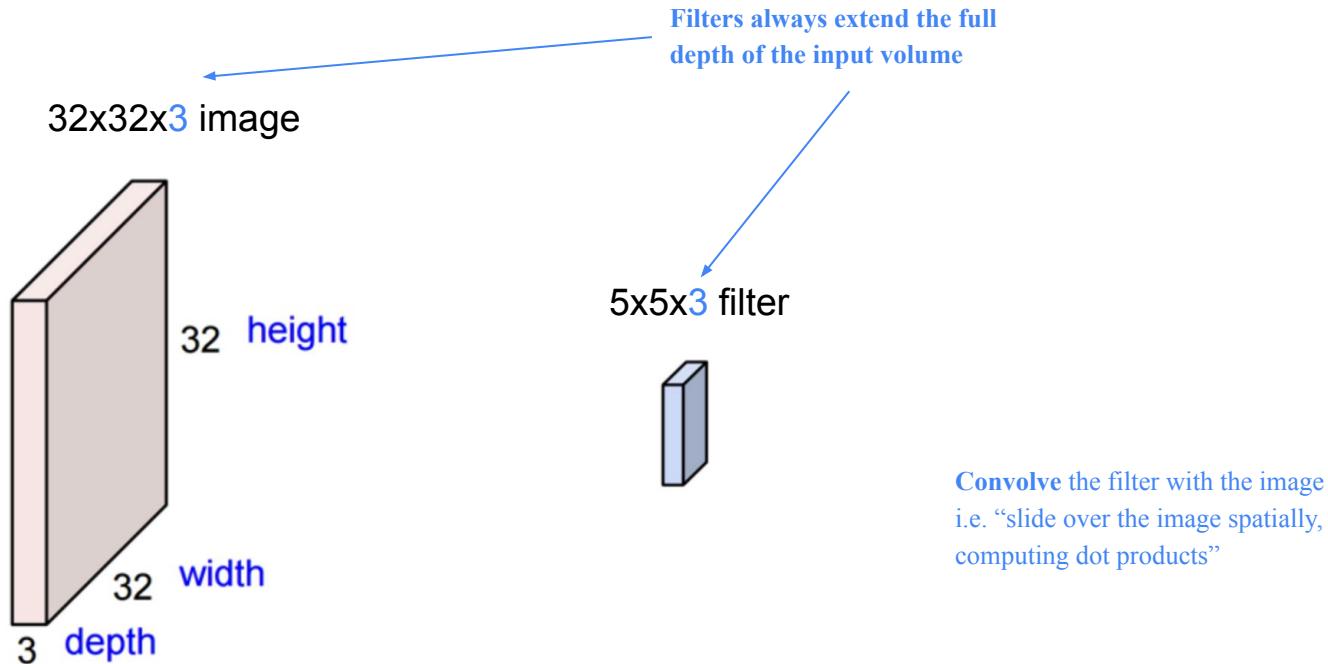


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolutional Neural Network (CNN)

Convolution Layer

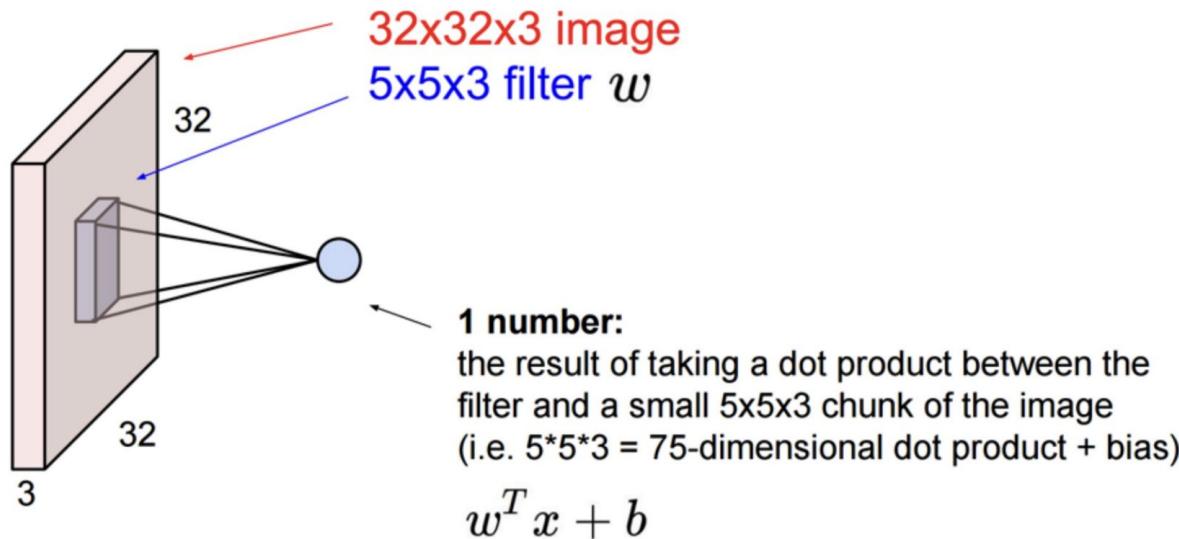
- Example



Convolutional Neural Network (CNN)

Convolution Layer

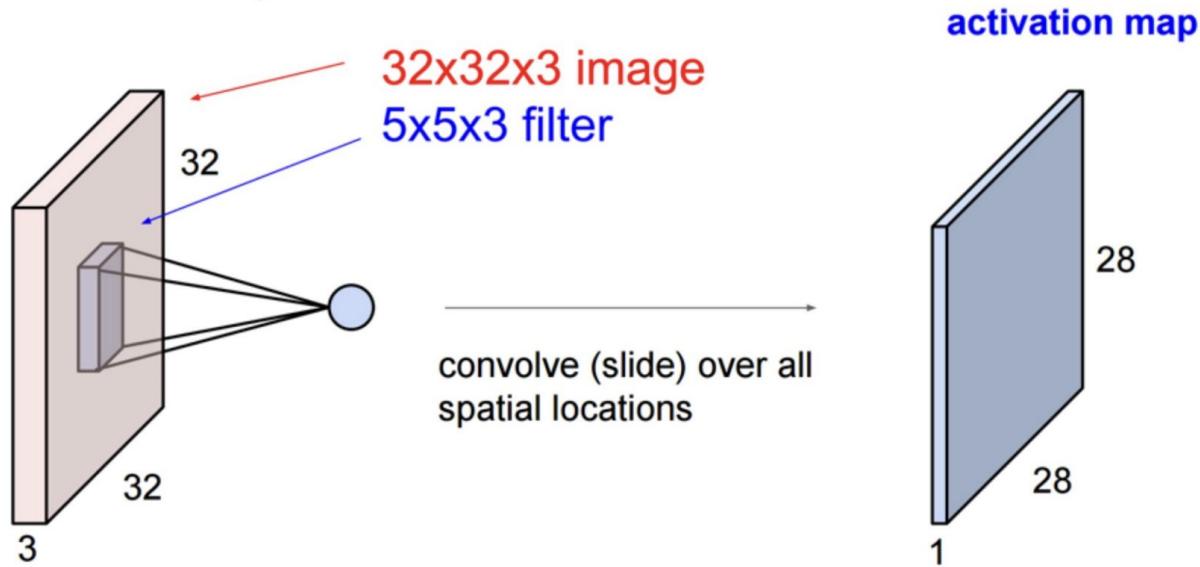
- Example



Convolutional Neural Network (CNN)

Convolution Layer

- Example

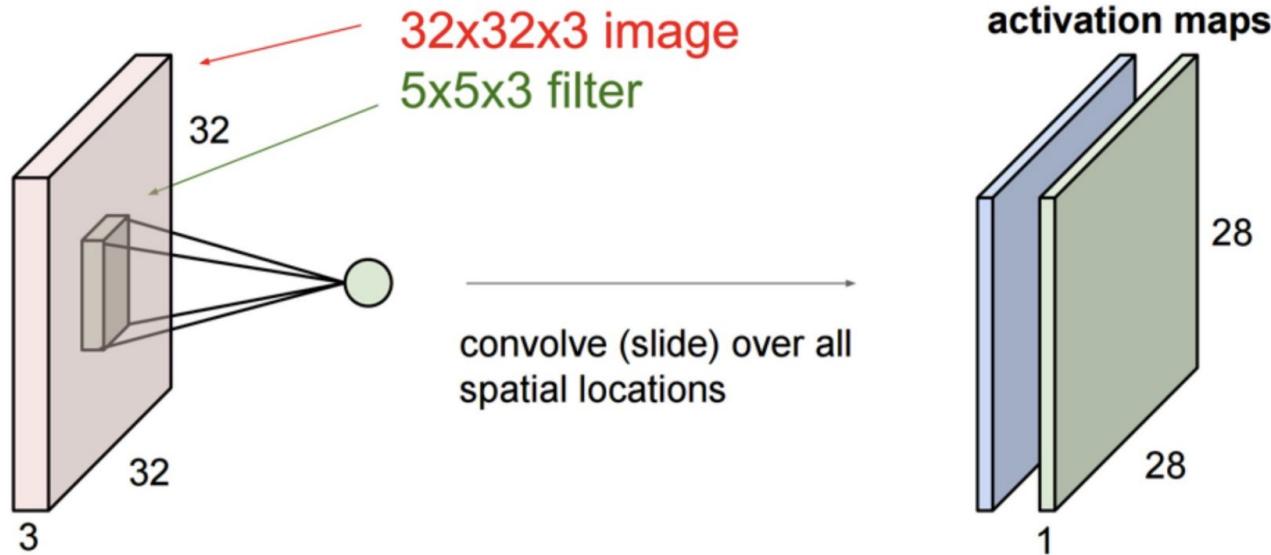


Convolutional Neural Network (CNN)

Convolution Layer

- Example

Consider a second, green filter

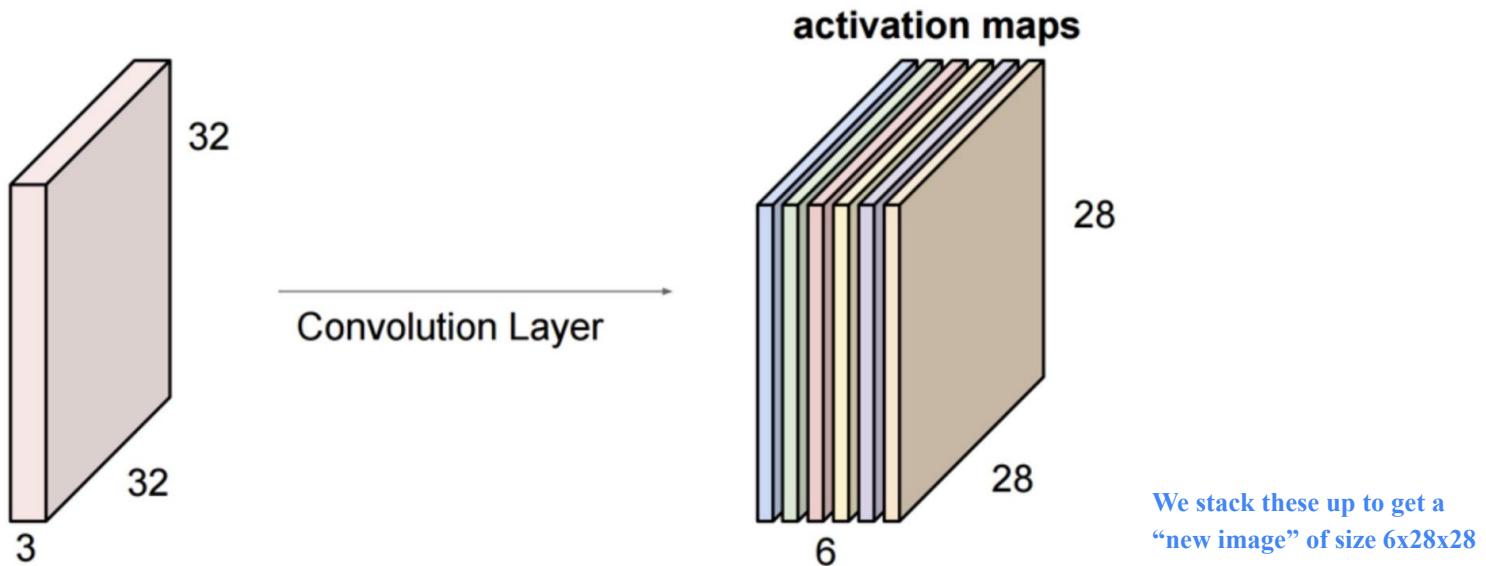


Convolutional Neural Network (CNN)

Convolution Layer

- Example

If we had 6 5x5 filters, we will get 6 separate activation maps



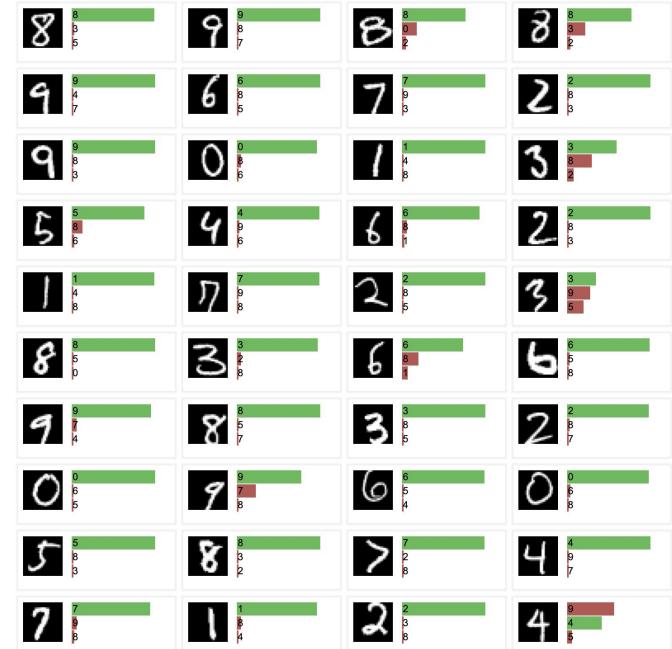
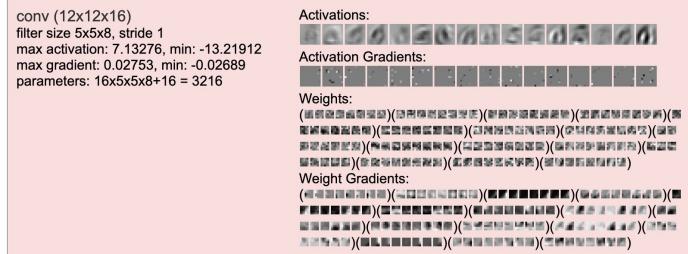
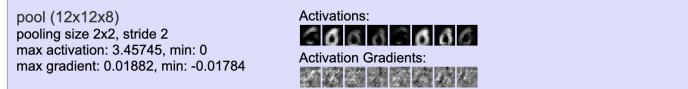
Convolutional Neural Network (CNN)

Convolution Layer

- Demos

<http://cs231n.stanford.edu/>

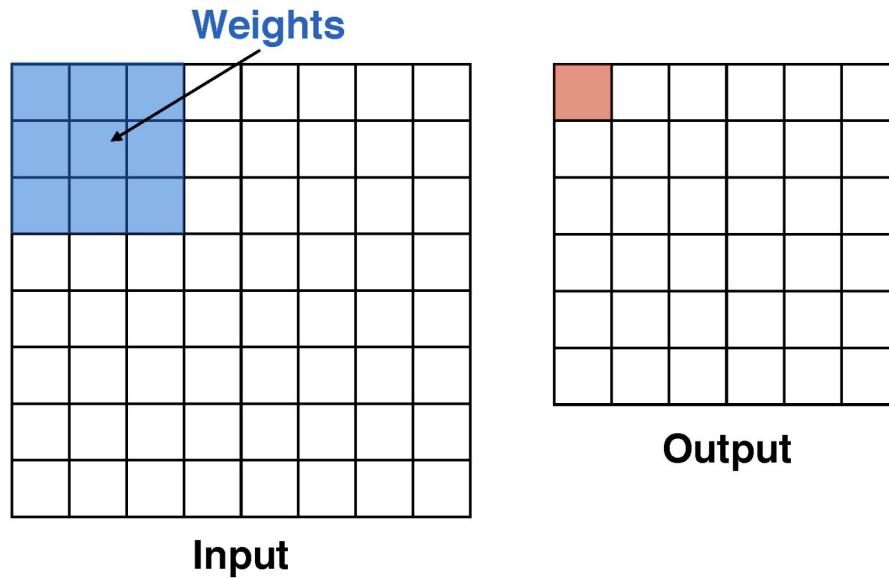
<http://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>



Convolutional Neural Network (CNN)

Convolution Layer

- Stride

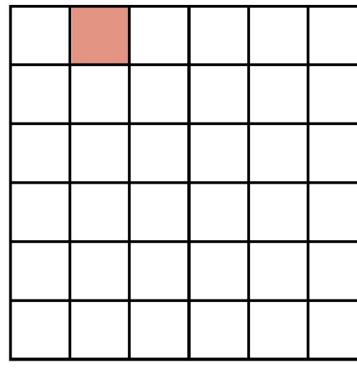
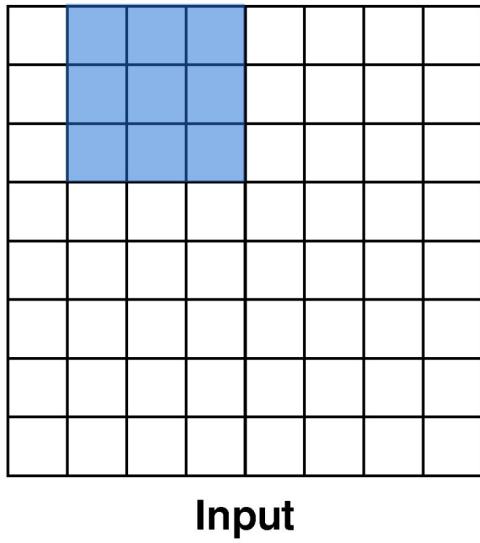


During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride



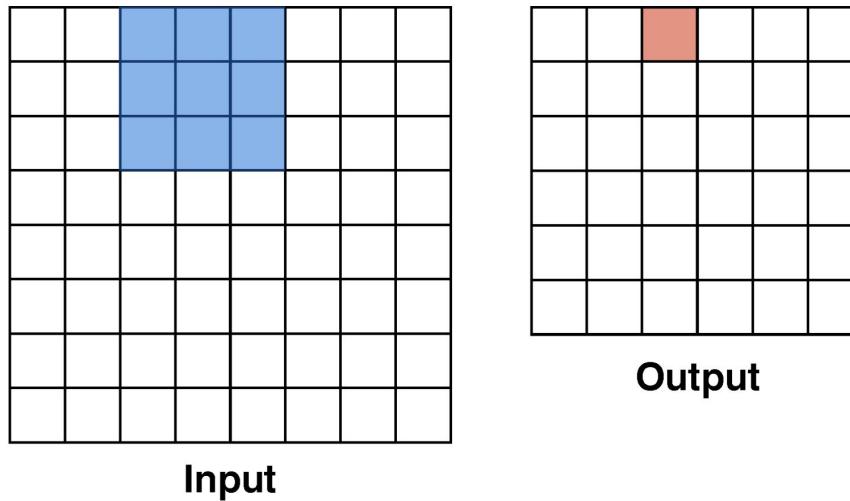
Output

During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride

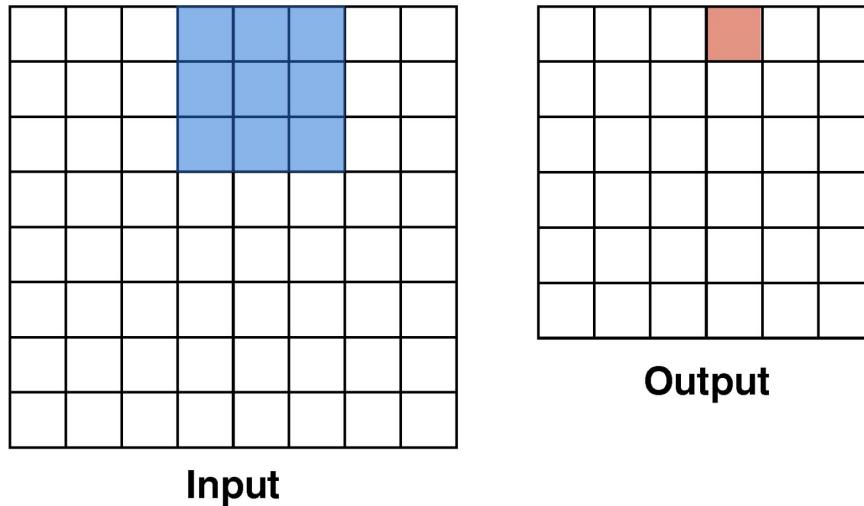


During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride

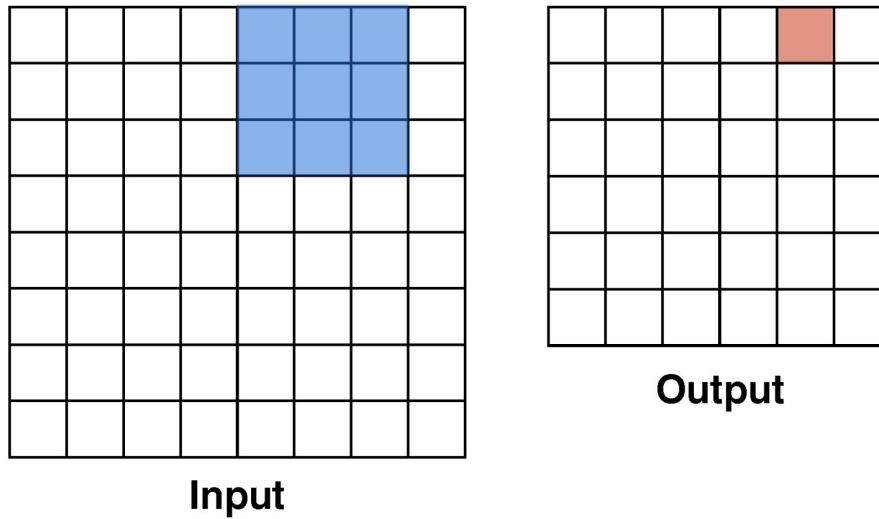


During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride

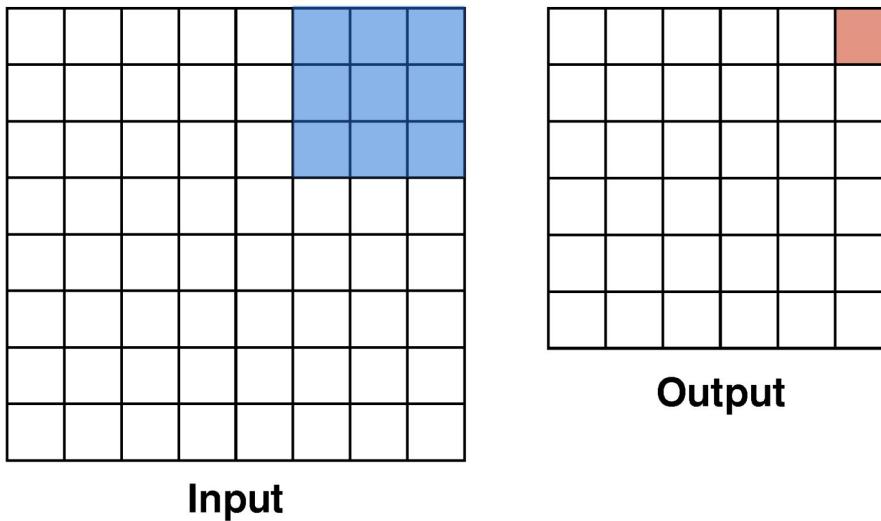


During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride

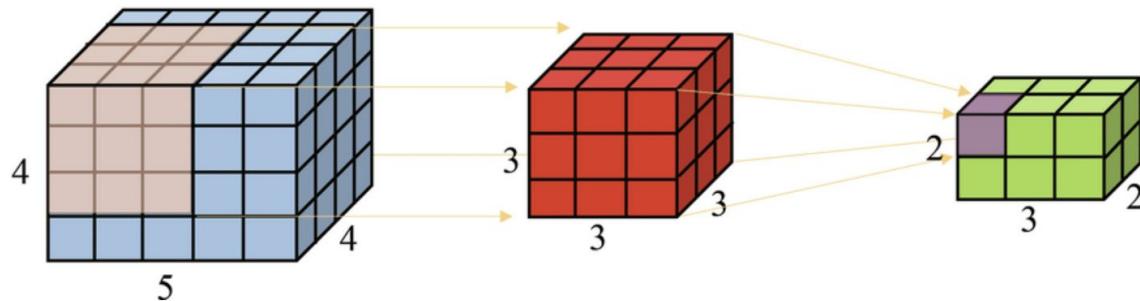


During convolution,
the weights “slide”
along the input to
generate each output

Convolutional Neural Network (CNN)

Convolution Layer

- Stride
 - Example of a 3D convolution operation with a stride of 1



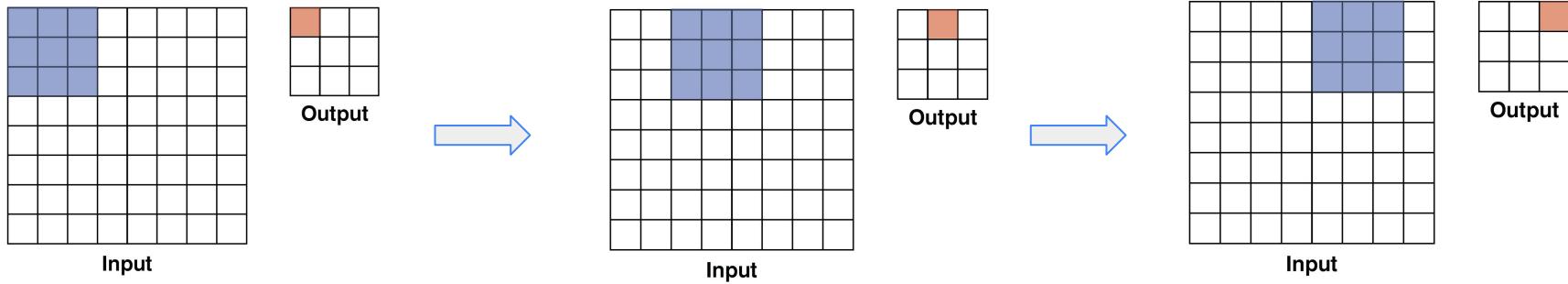
Note: each position we are doing a 3D sum: channel, row, column

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

Convolutional Neural Network (CNN)

Convolution Layer

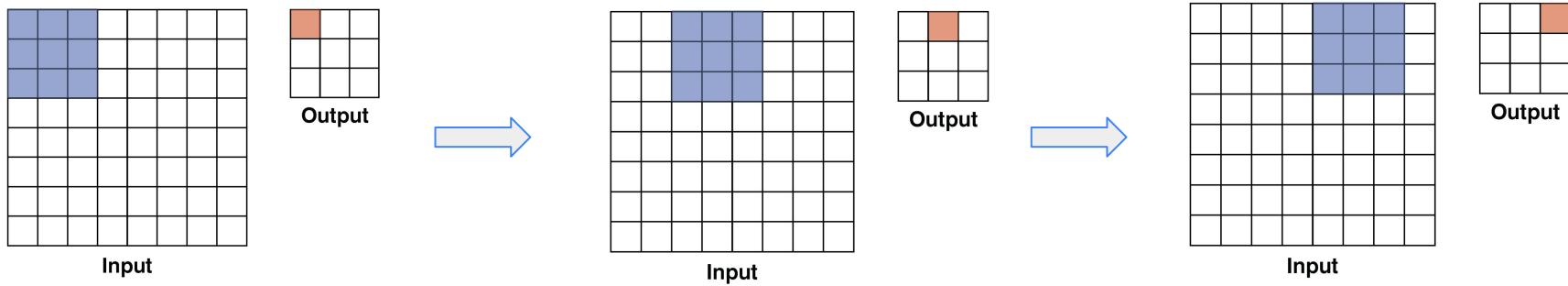
- Stride
 - Example: stride = 2



Convolutional Neural Network (CNN)

Convolution Layer

- Stride
 - Example: stride = 2

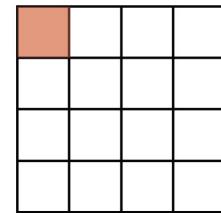
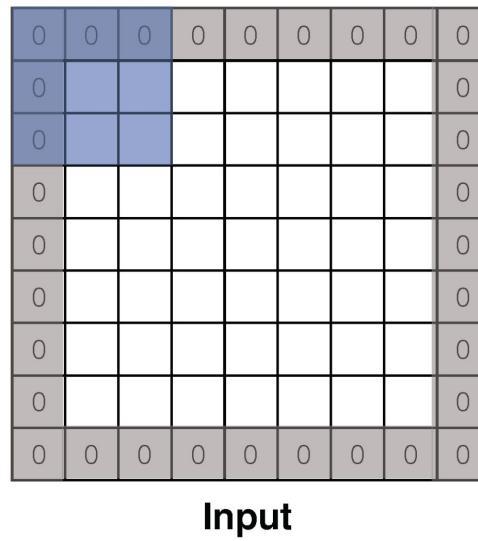


Notice that with certain strides, we may not be able to see all the input.
The output is also half the size of the input.

Convolutional Neural Network (CNN)

Convolution Layer

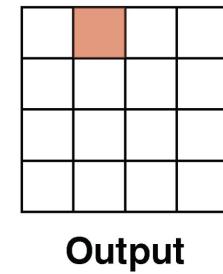
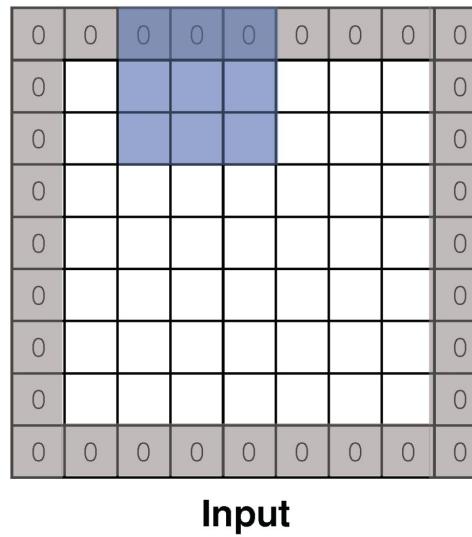
- Padding
 - We can also pad the input with zeros.
 - Example, pad = 1, stride = 2.



Convolutional Neural Network (CNN)

Convolution Layer

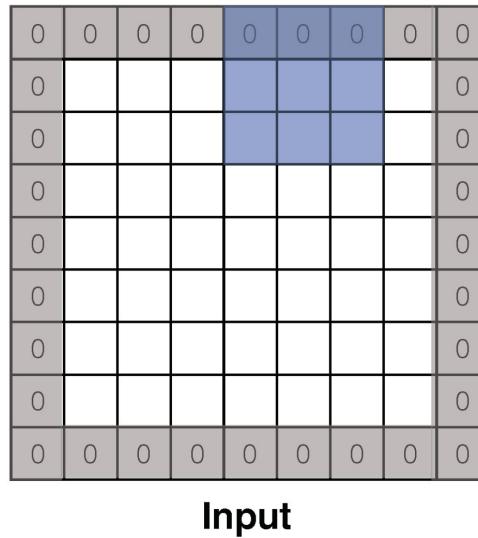
- Padding
 - We can also pad the input with zeros.
 - Example, pad = 1, stride = 2.



Convolutional Neural Network (CNN)

Convolution Layer

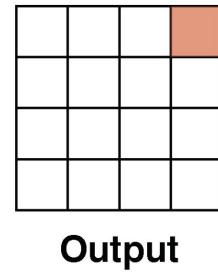
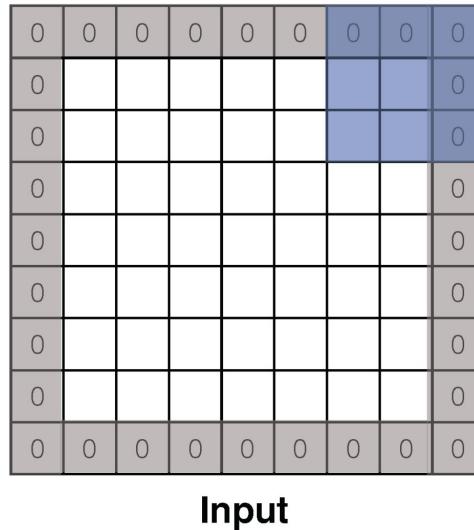
- Padding
 - We can also pad the input with zeros.
 - Example, pad = 1, stride = 2.



Convolutional Neural Network (CNN)

Convolution Layer

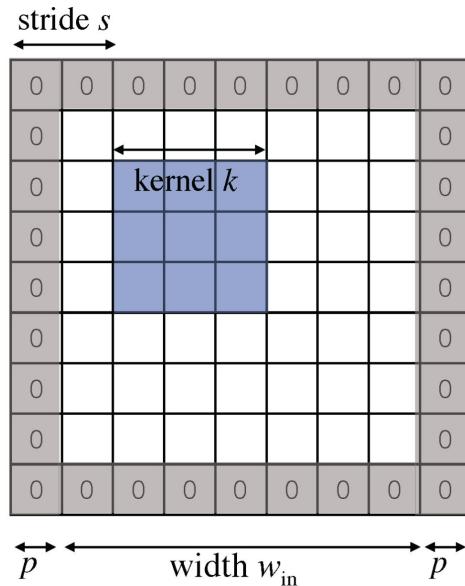
- Padding
 - We can also pad the input with zeros.
 - Example, pad = 1, stride = 2.



Convolutional Neural Network (CNN)

Convolution Layer

- Q: How big is the output?



In general, the output has size:

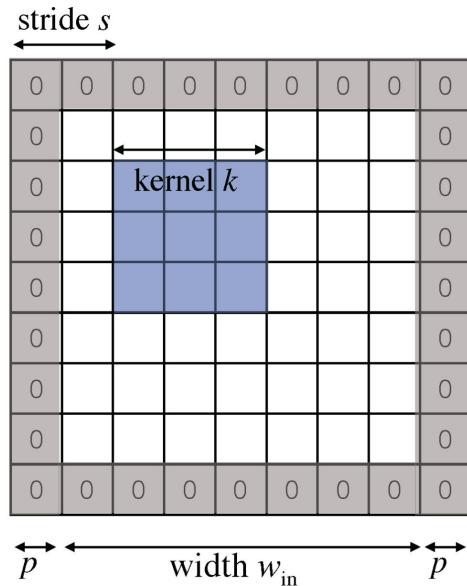
$$w_{out} = \left\lfloor \frac{w_{in} + 2p - k}{s} \right\rfloor + 1$$

Convolutional Neural Network (CNN)

Convolution Layer

- Q: How big is the output?

Example: $k = 3$, $s = 1$, $p = 1$



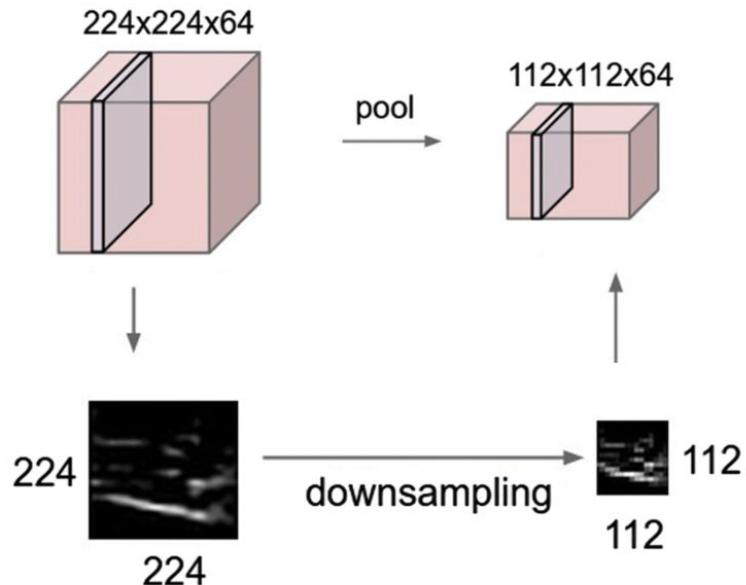
$$\begin{aligned}w_{\text{out}} &= \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1 \\&= \left\lfloor \frac{w_{\text{in}} + 2 - 3}{1} \right\rfloor + 1 \\&= w_{\text{in}}\end{aligned}$$

VGGNet [Simonyan 2014]
used filters of this shape

Convolutional Neural Network (CNN)

Pooling Layer

- For most ConvNets, convolution is often followed by **pooling**
 - Creates a smaller representation while retaining the most important information
 - The “**Max**” operation is the most common
 - Max-pooling is a **downsampling** operation often used in between convolution layers to reduce the spatial dimensions of the feature maps.
 - It helps in simplifying the network and extracting important features.

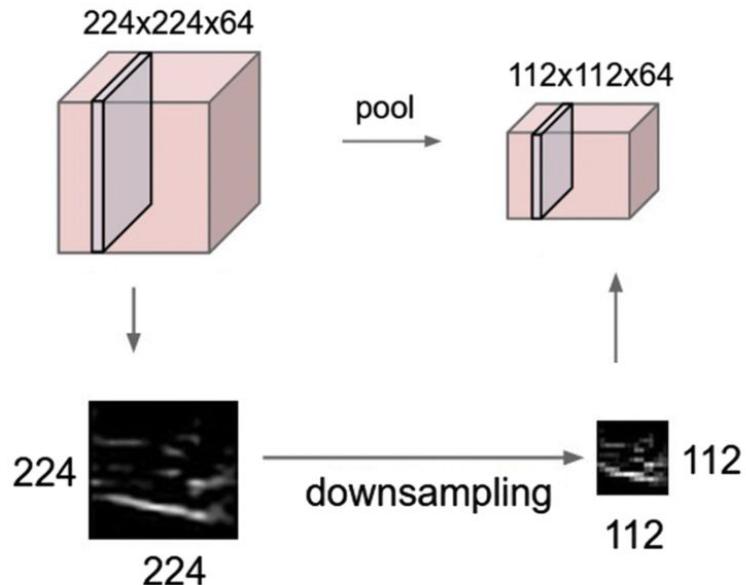


Q: Why might “Avg” be a poor choice?

Convolutional Neural Network (CNN)

Pooling Layer

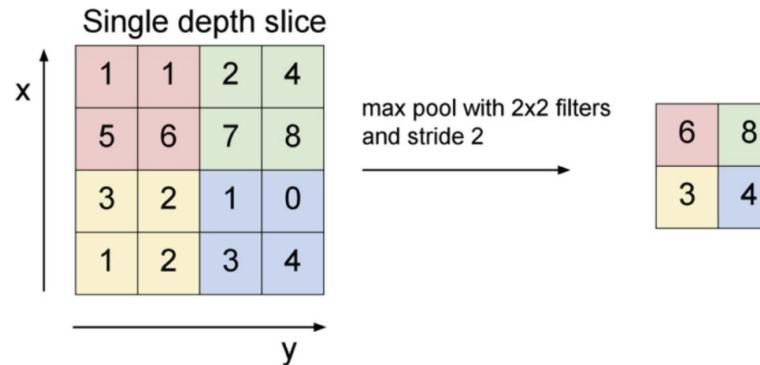
- Q: Why might “Avg” be a poor choice?
 - Makes the representations smaller and more manageable
 - Operates over each activation map independently



Convolutional Neural Network (CNN)

Pooling Layer

- Max Pooling
 - What's the backprop rule for max pooling?
 - In the forward pass, store the index that took the max
 - The backprop gradient is the input gradient at the index

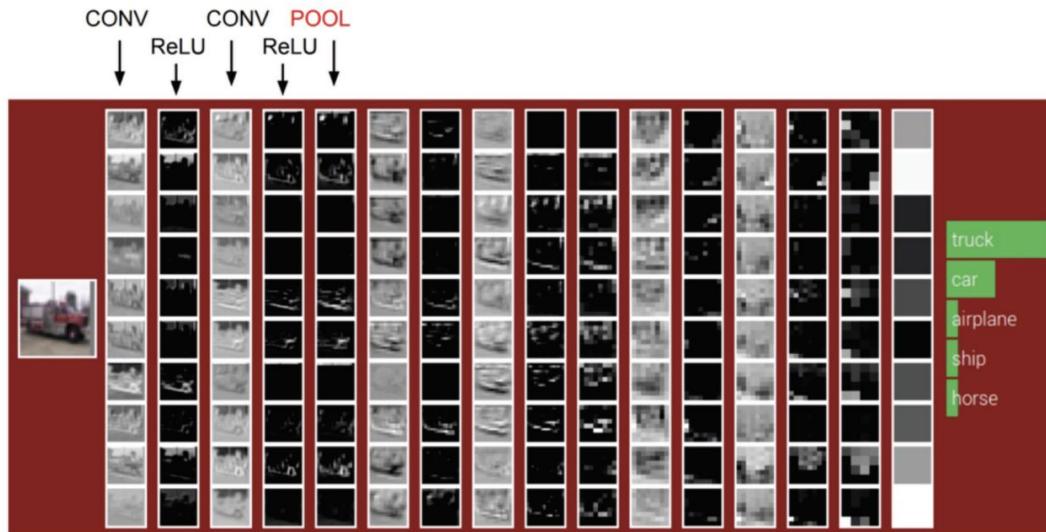


Interpretation of intelligence in CNN-pooling processes: a methodological survey - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/Down-sampling-using-pooling-5_fig1_334208543 [accessed 16 Sept 2024]

Convolutional Neural Network (CNN)

Pooling Layer

- Example ConvNet

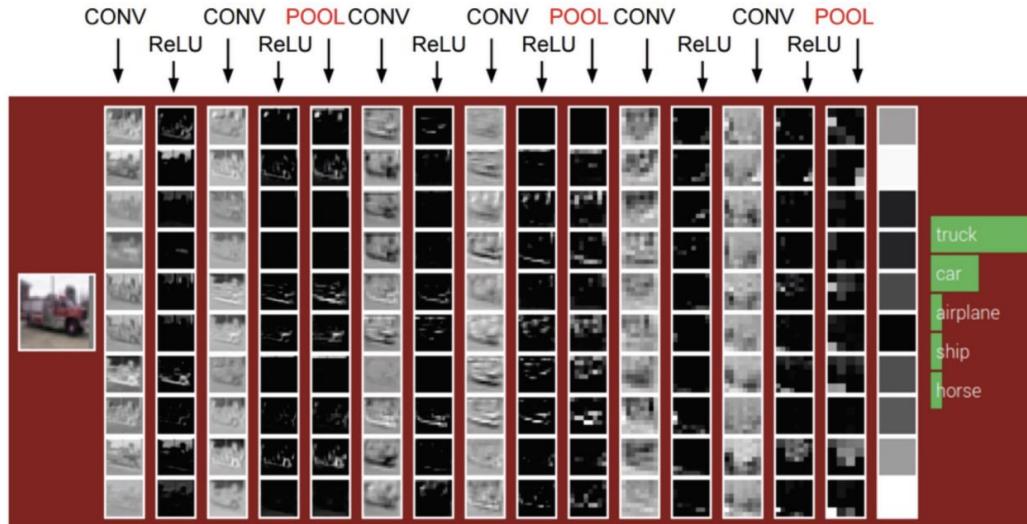


*Interpretation of intelligence in CNN-pooling processes: a methodological survey - Scientific Figure on ResearchGate. Available from:
https://www.researchgate.net/figure/Down-sampling-using-pooling-5_fig1_334208543 [accessed 16 Sept 2024]*

Convolutional Neural Network (CNN)

Pooling Layer

- Example ConvNet

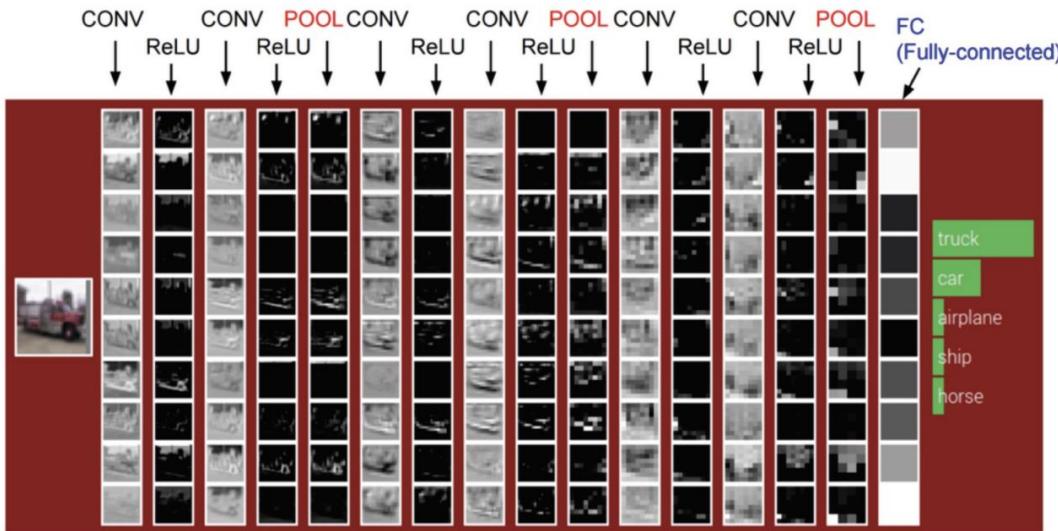


Interpretation of intelligence in CNN-pooling processes: a methodological survey - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/Down-sampling-using-pooling-5_fig1_334208543 [accessed 16 Sept 2024]

Convolutional Neural Network (CNN)

Pooling Layer

- Example ConvNet

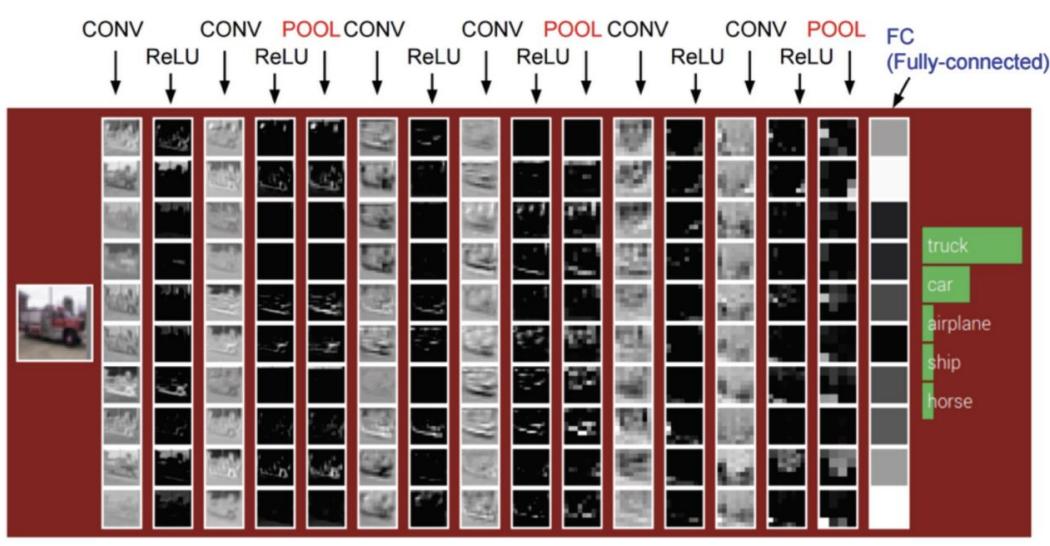


Interpretation of intelligence in CNN-pooling processes: a methodological survey - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/Down-sampling-using-pooling-5_fig1_334208543 [accessed 16 Sept 2024]

Convolutional Neural Network (CNN)

Pooling Layer

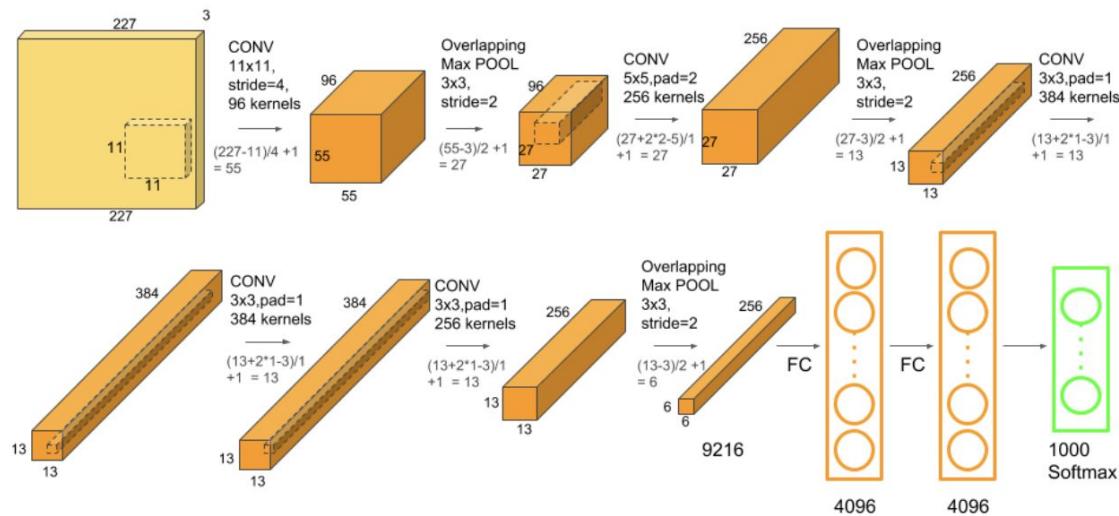
- Example ConvNet



10x3x3 conv filters,
stride=1, pad=1
2x2 pool filters,
stride=2

Convolutional Neural Network (CNN)

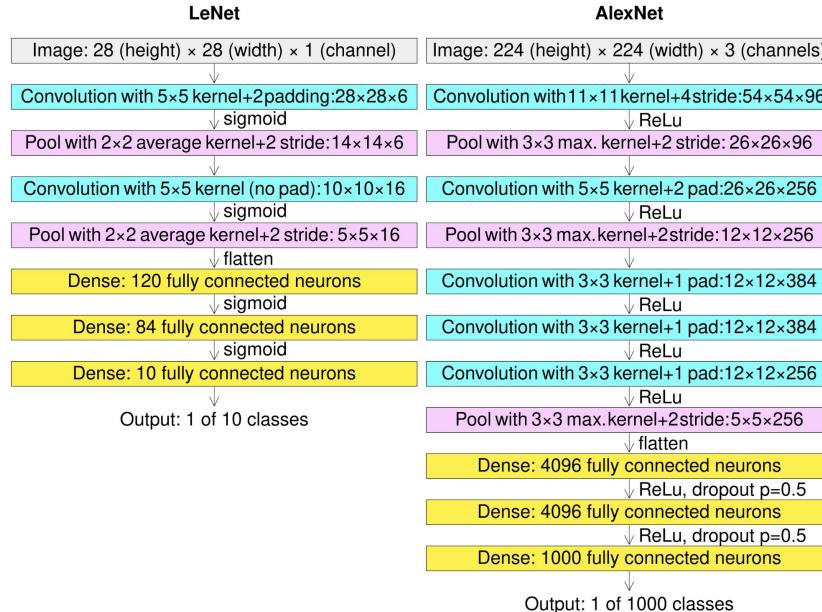
Example: AlexNet



AlexNet consists of 5 Convolutional Layers and 3 Fully Connected Layers

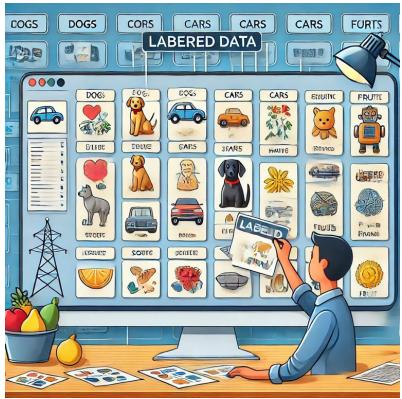
Convolutional Neural Network (CNN)

Example: Comparison of the LeNet and AlexNet convolution, pooling, and dense layers



Training ConvNets

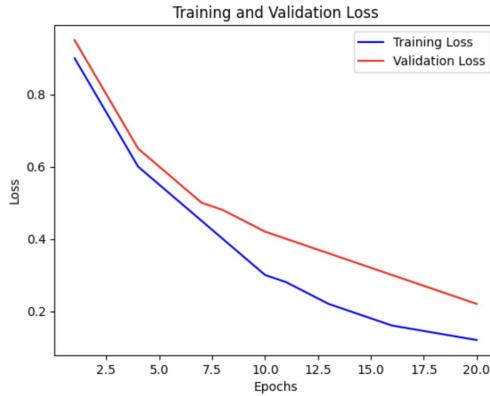
Q: How do you actually train these things?



Gather labeled data



Build a ConvNet architecture



Minimize the loss

Training ConvNets

- **Training a convolutional neural network**
 - **Loading the Data Set**
 - **Normalization, Reshape and Label Encoding**
 - **Train Test Split**
 - **Convolutional Neural Network**
 - Convolution Operation
 - Same Padding
 - Max Pooling
 - Flattening
 - Full Connection

Training ConvNets

- Training a convolutional neural network
 - Loading the Data Set

```
# read train
train = pd.read_csv("/Users/xiaotingzhou/Documents/GitHub/AI_Civil/Dataset/Digit Recognizer/train.csv")
print(train.shape)
train.head()
```

Python

(42000, 785)

	label	pixel0	pixel1	pixel2	pixel3	pixel4	pixel5	pixel6	pixel7	pixel8	...	pixel774	pixel775	pixel776	pixel777	pixel778	pixel779	pixel780	pixel781	pixel782	pixel783
0	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
2	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
3	4	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	

5 rows x 785 columns

Training ConvNets

- Training a convolutional neural network
 - Normalization, Reshape and Label Encoding

```
# Normalize the data
X_train = X_train / 255.0
test = test / 255.0
print("x_train shape: ",X_train.shape)
print("test shape: ",test.shape)
```

```
x_train shape: (42000, 28, 28, 1)
test shape: (28000, 28, 28, 1)
```

```
# Reshape
X_train = X_train.reshape(-1,28,28,1)
test = test.reshape(-1,28,28,1)
print("x_train shape: ",X_train.shape)
print("test shape: ",test.shape)
```

```
x_train shape: (42000, 28, 28, 1)
test shape: (28000, 28, 28, 1)
```

```
# Label Encoding
# from keras.utils import np_utils # convert to one-hot-encoding
from tensorflow.keras.utils import to_categorical

Y_train = to_categorical(Y_train, num_classes = 10)
```

Training ConvNets

- Training a convolutional neural network

- Train Test Split

We split the data into train and test sets.

- test size is 10%.
- train size is 90%.

```
# Split the train and the validation set for the fitting
from sklearn.model_selection import train_test_split
X_train, X_val, Y_train, Y_val = train_test_split(X_train, Y_train, test_size = 0.1, random_state=2)
print("x_train shape",X_train.shape)
print("x_test shape",X_val.shape)
print("y_train shape",Y_train.shape)
print("y_test shape",Y_val.shape)
```

```
x_train shape (37800, 28, 28, 1)
x_test shape (4200, 28, 28, 1)
y_train shape (37800, 10, 10)
y_test shape (4200, 10, 10)
```

Training ConvNets

- Training a convolutional neural network

- Convolutional Neural Network

- Convolution Operation
 - Same Padding
 - Max Pooling
 - Flattening
 - Full Connection

Create Model

- conv => max pool => dropout =>
- conv => max pool => dropout =>
- fully connected (2 layer)**
- **Dropout:** Dropout is a technique
where randomly selected neurons are
ignored during training

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout, Flatten, Conv2D, MaxPool2D, Reshape
from tensorflow.keras.optimizers import RMSprop, Adam
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.callbacks import ReduceLROnPlateau
from tensorflow.keras import backend as K

# Define custom loss function if necessary (e.g., categorical crossentropy for multi-class segmentation)
def custom_loss_function(y_true, y_pred):
    return K.categorical_crossentropy(y_true, y_pred)

# Model definition
model = Sequential()

model.add(Conv2D(filters=8, kernel_size=(5,5), padding='Same', activation='relu', input_shape=(28, 28, 1)))
model.add(MaxPool2D(pool_size=(2,2)))
model.add(Dropout(0.25))

model.add(Conv2D(filters=16, kernel_size=(3,3), padding='Same', activation='relu'))
model.add(MaxPool2D(pool_size=(2,2), strides=(2,2)))
model.add(Dropout(0.25))

# Fully connected
model.add(Flatten())
model.add(Dense(256, activation='relu'))
model.add(Dropout(0.5))

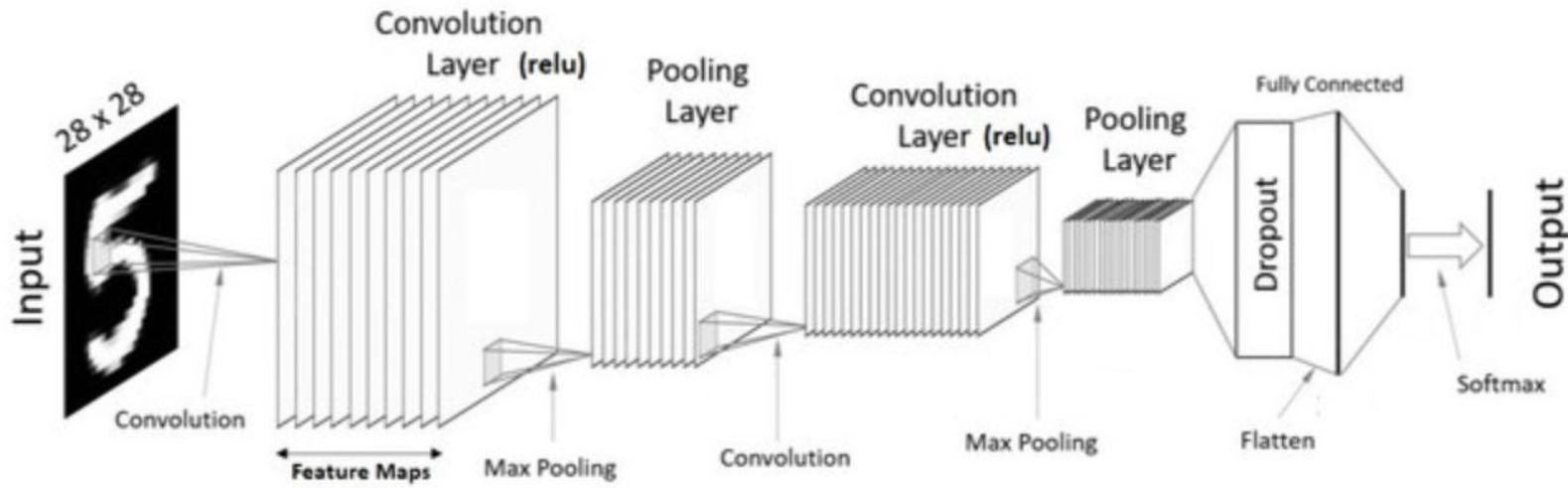
# Change to Dense to output a (10, 10) matrix, then reshape to (10, 10, 1)
model.add(Dense(10*10, activation='softmax')) # Flattening into 100 units
model.add(Reshape((10, 10, 1))) # Reshape to (10, 10, 1)

# Compile the model
model.compile(optimizer=Adam(), loss='categorical_crossentropy', metrics=['accuracy'])

# Print model summary to verify architecture
model.summary()
```

Python: CNN in Python

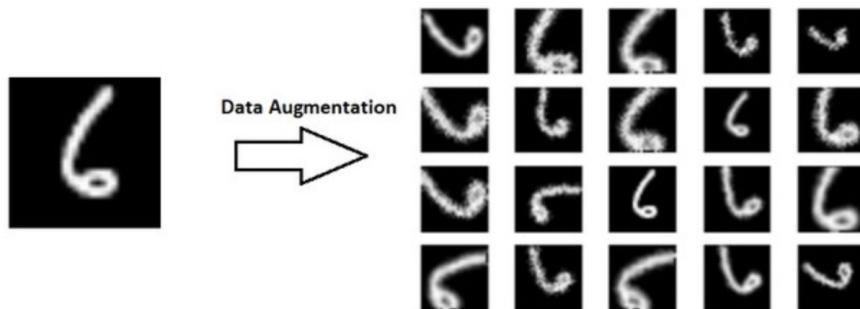
Example: Digit Recognizer



Python: CNN in Python

Example: Digit Recognizer

- **Implementing with Keras**
 - Create Model
 - Define Optimizer
 - Compile Model
 - Epochs and Batch Size
 - Data Augmentation
 - Fit the Model
 - Evaluate the Model



Python: CNN in Python

Example: Digit Recognizer

Create Model

- conv => max pool =>
dropout => conv =>
max pool => dropout
=> **fully connected**
(2 layer)
- **Dropout:** Dropout is
a technique where
randomly selected
neurons are ignored
during training

Layer (type)	Output Shape	Param #
conv2d_6 (Conv2D)	(None, 28, 28, 8)	208
max_pooling2d_6 (MaxPooling2D)	(None, 14, 14, 8)	0
dropout_9 (Dropout)	(None, 14, 14, 8)	0
conv2d_7 (Conv2D)	(None, 14, 14, 16)	1,168
max_pooling2d_7 (MaxPooling2D)	(None, 7, 7, 16)	0
dropout_10 (Dropout)	(None, 7, 7, 16)	0
flatten_3 (Flatten)	(None, 784)	0
dense_6 (Dense)	(None, 256)	200,960
dropout_11 (Dropout)	(None, 256)	0
dense_7 (Dense)	(None, 100)	25,700
reshape (Reshape)	(None, 10, 10, 1)	0

Python: CNN in Python

Example: Digit Recognizer

