

CX4240 Spring 2026

Convolution Neural Networks

Bo Dai
School of CSE, Georgia Tech
bodai@cc.gatech.edu

Multi-Layer Perceptron

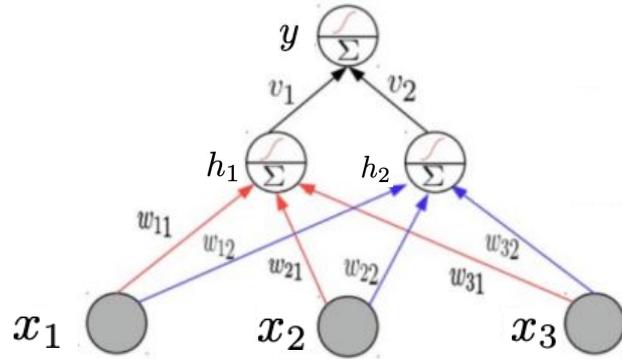
$$y = o(\mathbf{V}g(\mathbf{W}x))$$

$$\mathbf{V} = [v_1, v_2]$$

$$h = [h_1, h_2]^T = g(\mathbf{W}x)$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$x = [x_1, x_2, x_3]^T$$



Multi-Layer Perceptron

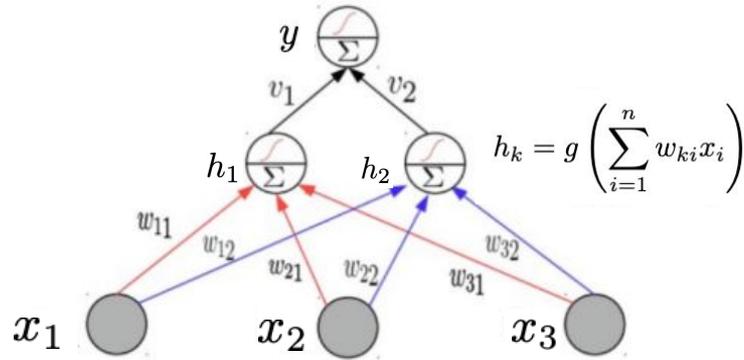
$$y = o(\mathbf{V}g(\mathbf{W}\mathbf{x}))$$

$$\mathbf{V} = [v_1, v_2]$$

$$\mathbf{h} = [h_1, h_2]^\top = g(\mathbf{W}\mathbf{x})$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$\mathbf{x} = [x_1, x_2, x_3]^\top$$



Multi-Layer Perceptron

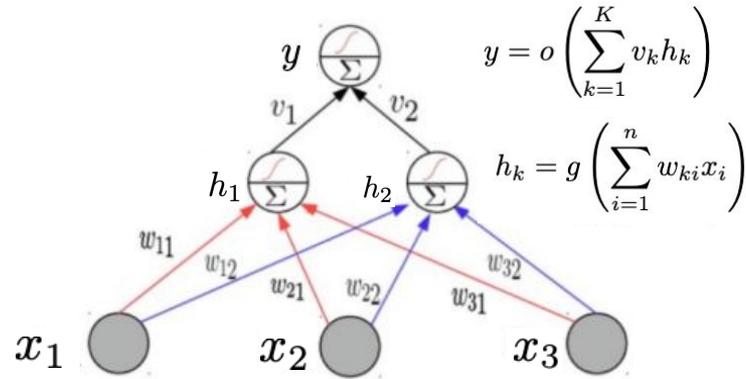
$$y = o(\mathbf{V}g(\mathbf{W}\mathbf{x}))$$

$$\mathbf{V} = [v_1, v_2]$$

$$\mathbf{h} = [h_1, h_2]^\top = g(\mathbf{W}\mathbf{x})$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$\mathbf{x} = [x_1, x_2, x_3]^\top$$



Multi-Layer Perceptron

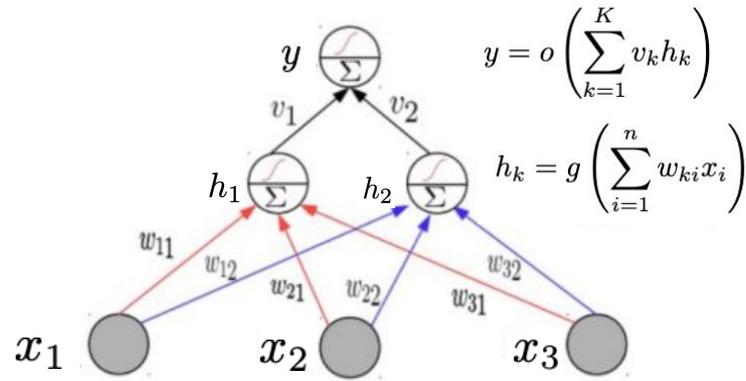
$$y = o(\mathbf{V}g(\mathbf{W}\mathbf{x}))$$

$$\mathbf{V} = [v_1, v_2]$$

$$\mathbf{h} = [h_1, h_2]^\top = g(\mathbf{W}\mathbf{x})$$

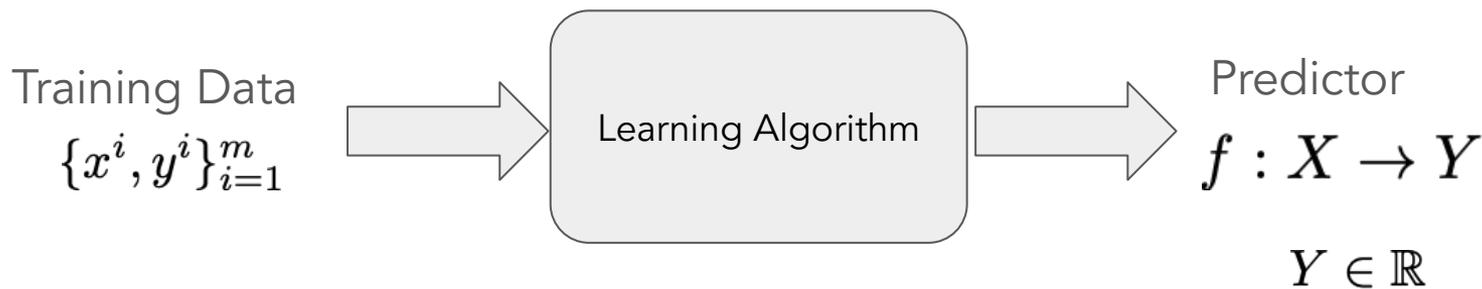
$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

$$\mathbf{x} = [x_1, x_2, x_3]^\top$$



Forward Pass

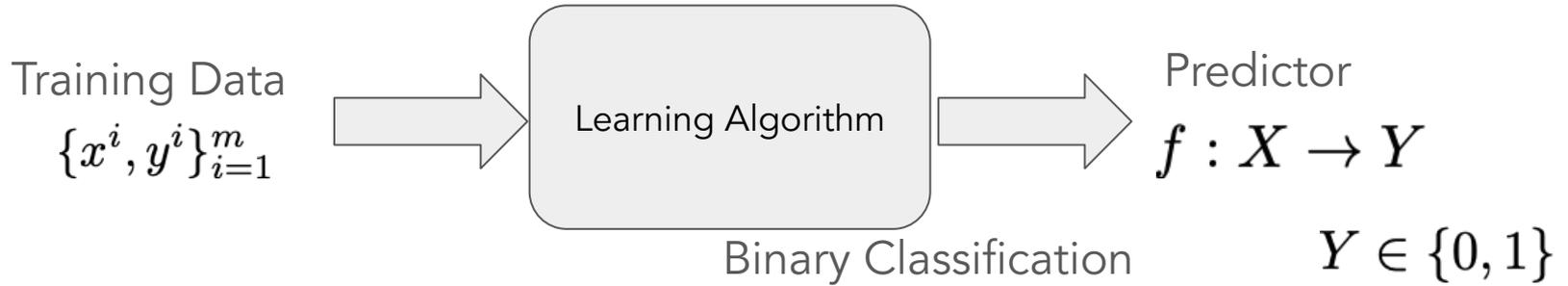
Regression Algorithms



Linear Regression Pipeline

1. Build probabilistic models:
Gaussian Distribution + Neural Network $y = o(\mathbf{V}g(\mathbf{W}x))$
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

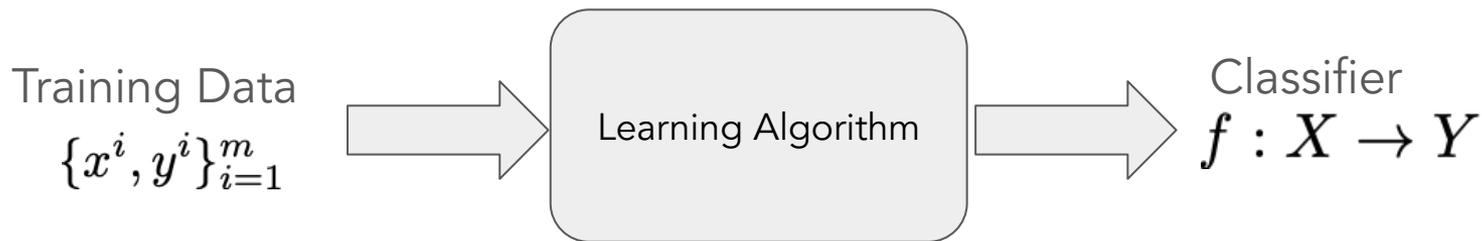
Binary Classification Algorithms



Binary Logistic Regression Pipeline

1. Build probabilistic models:
Bernoulli Distribution + Neural Network $y = o(\mathbf{V}g(\mathbf{W}x))$
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Multiclass Logistic Regression Algorithms



Multiclass Classification $Y \in \{0, 1, \dots, k\}$
Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
Categorical Distribution + Neural Network $y = o(\mathbf{V}g(\mathbf{W}x))$
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$
$$\theta = [V, W]$$

$$\ell(x^i, y^i, \theta) = (o(Vg(Wx^i)) - y^i)^2$$

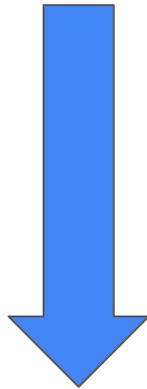
$$\ell(x^i, y^i, \theta) = -y^i \log \sigma(o(Vg(Wx^i)))$$
$$-(1 - y^i) \log(1 - \sigma(o(Vg(Wx^i))))$$

$$\ell(x^i, y^i, \theta) = - \sum_{j=1}^k y^i \log \frac{\exp(o(V_j g(Wx^i)))}{\sum_{c=1}^k \exp(o(V_c g(Wx^i)))}$$

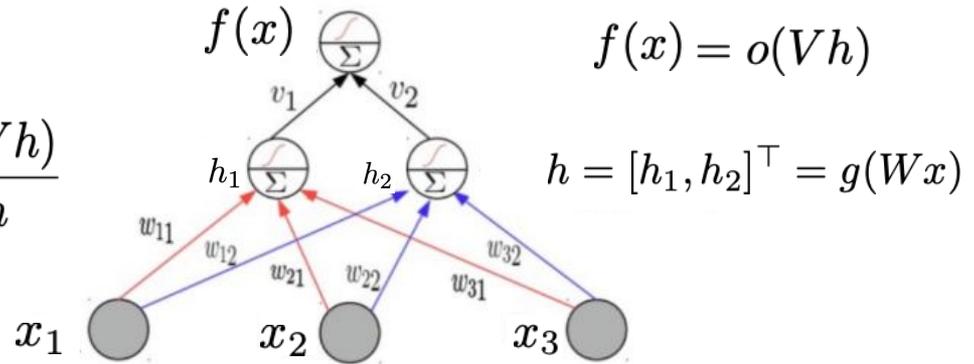
- (Stochastic) Gradient Descent

Backpropagation: Chain Rule on Neural Network

$$\ell(x^i, y^i, \theta) = (f(x^i, \mathbf{V}, \mathbf{W}) - y^i)^2$$



$$\frac{\partial \ell(x^i, y^i, \theta)}{\partial f}$$
$$\frac{\partial o(\mathbf{V}h)}{\partial \mathbf{V}} \frac{\partial o(\mathbf{V}h)}{\partial h}$$
$$\frac{\partial h}{\partial \mathbf{W}}$$



Backward Pass

(Stochastic) Gradient Descent

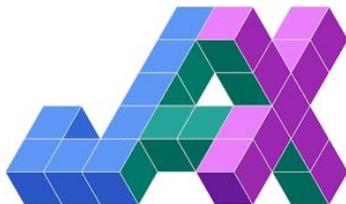
- Initialize parameter θ^0
- Sample $\{x^i, y^i\}_{i=1}^B$
- Do
$$\theta^{t+1} \leftarrow \theta^t - \eta \sum_{i=1}^B \nabla_{\theta} \ell(x^i, y^i, \theta^t) - (\lambda \nabla \Omega(\theta^t))$$

Auto-differentiation Packages

PyTorch



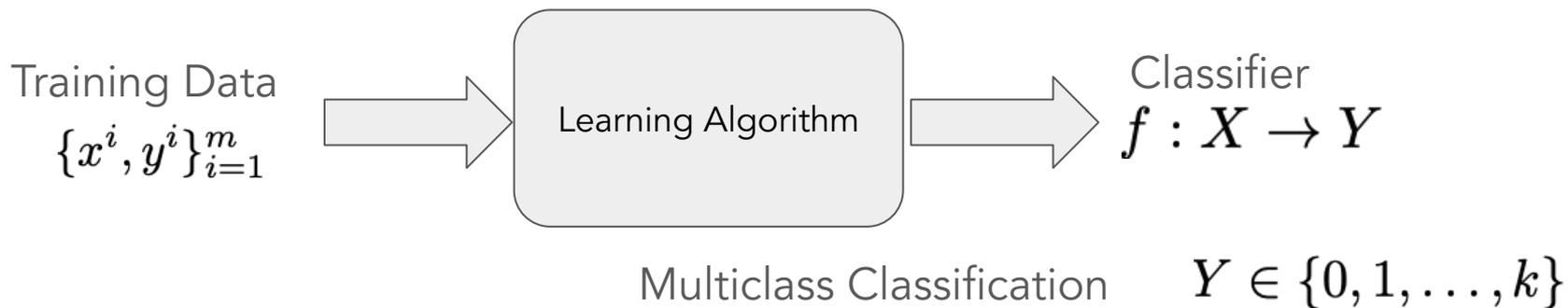
JAX



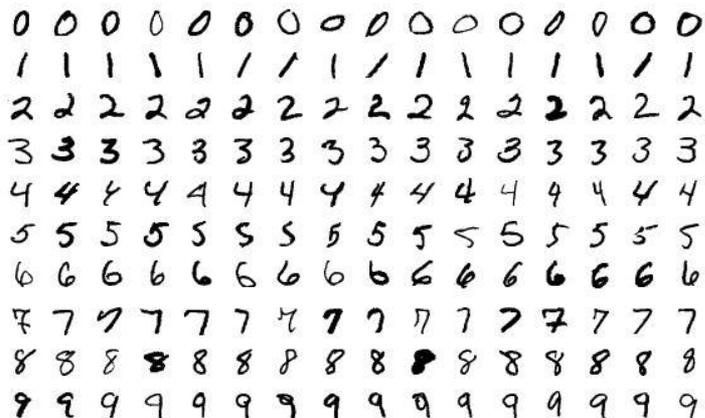
Tensorflow



Multiclass Logistic Regression Algorithms with CNN

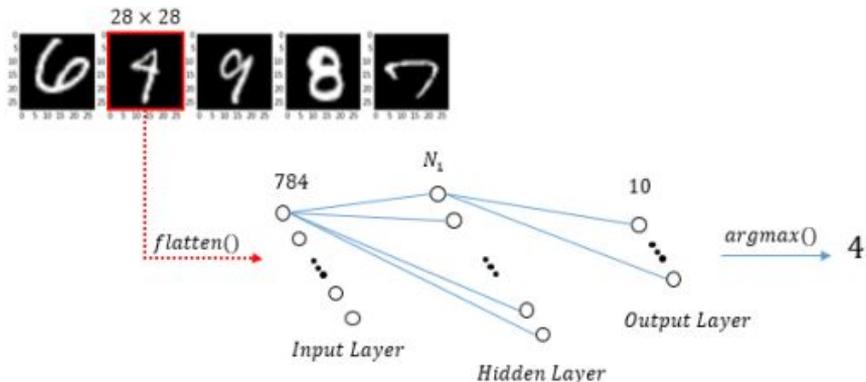


Revisit MLP for MNIST

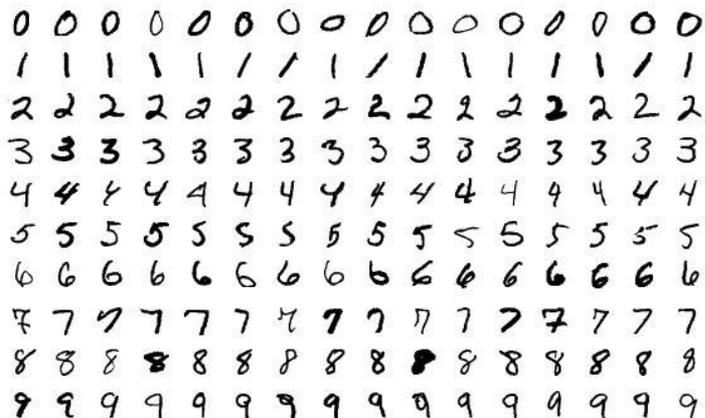


- 60,000 images
- 28x28 pixels = 784
- Grayscale, from 0 to 255 → Converted to [0,1]

MNIST hand-written character recognition

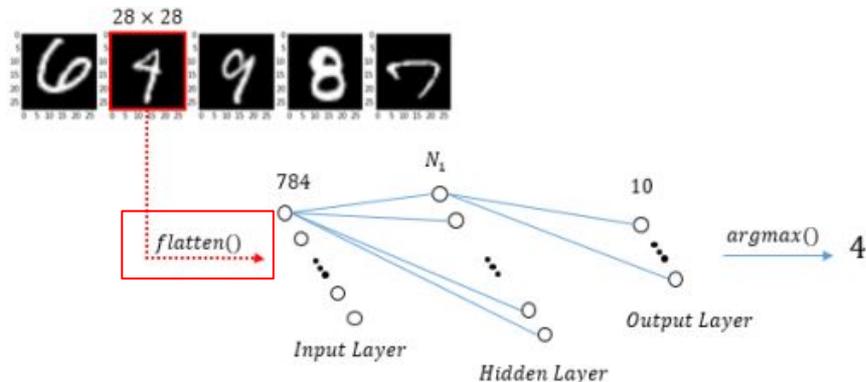


Revisit MLP for MNIST

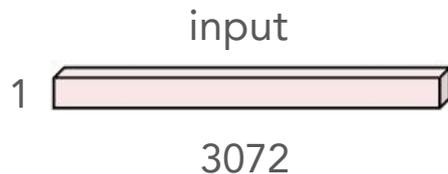
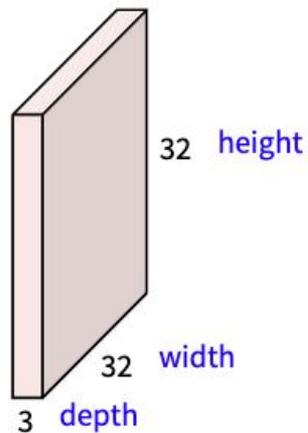
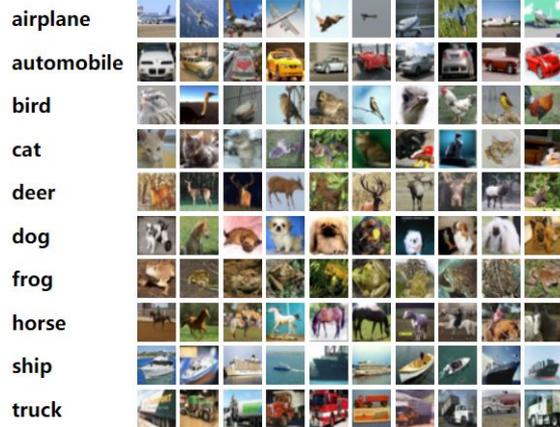


- 60,000 images
- 28x28 pixels = 784
- Grayscale, from 0 to 255 → Converted to [0,1]

MNIST hand-written character recognition



MLP for CIFAR10



32x32x3 image \rightarrow stretch to 3072x1

MLP: CIFAR10

airplane



automobile



bird



cat



deer



dog



frog



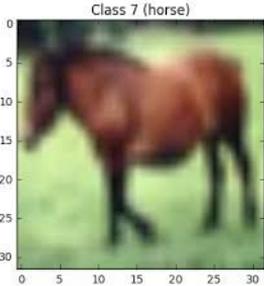
horse



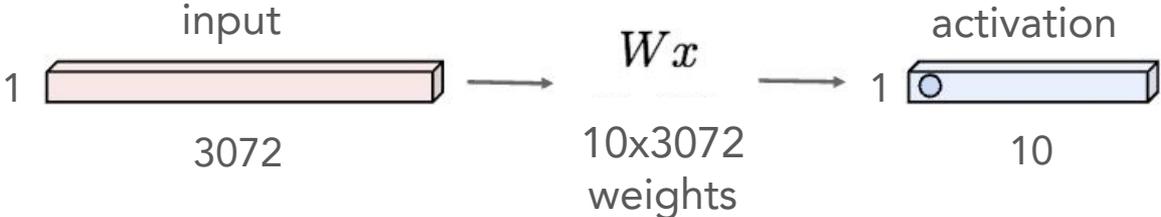
ship



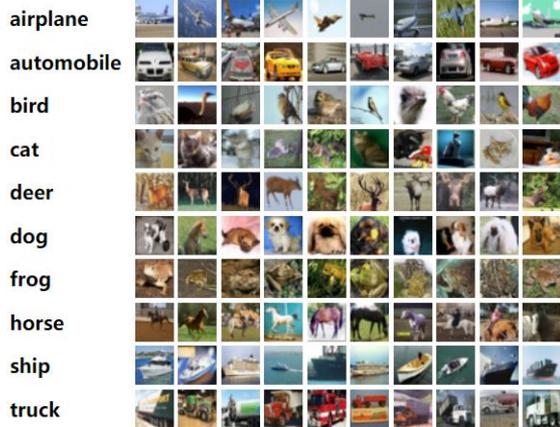
truck



32x32x3 image → stretch to 3072x1

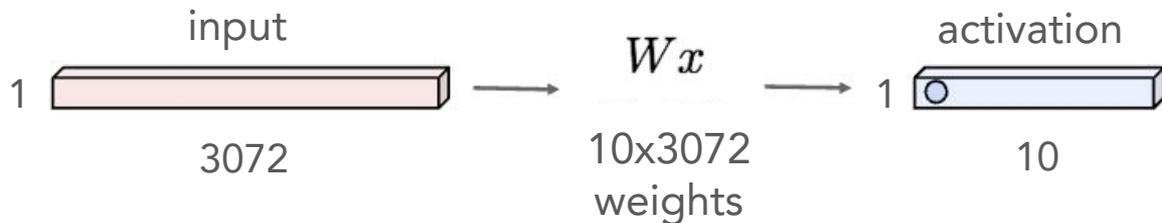


MLP: CIFAR10

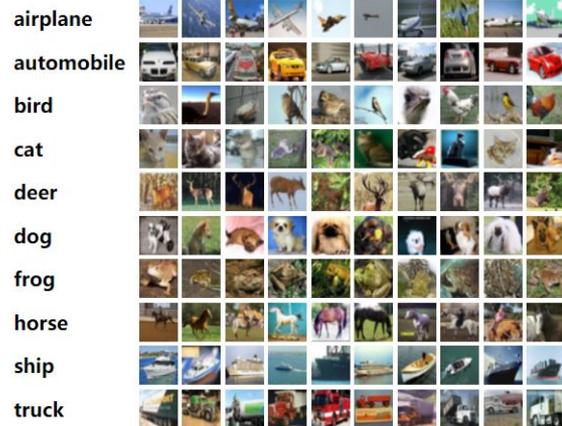


32x32x3 image \rightarrow stretch to 3072x1

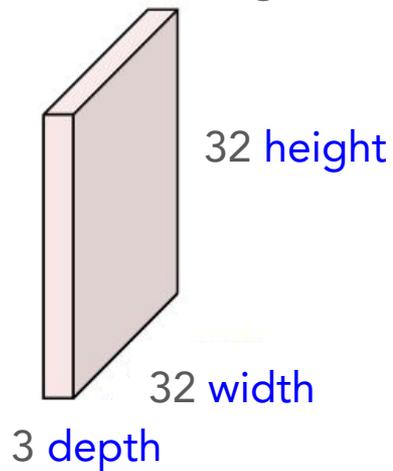
ruin the spatial structure



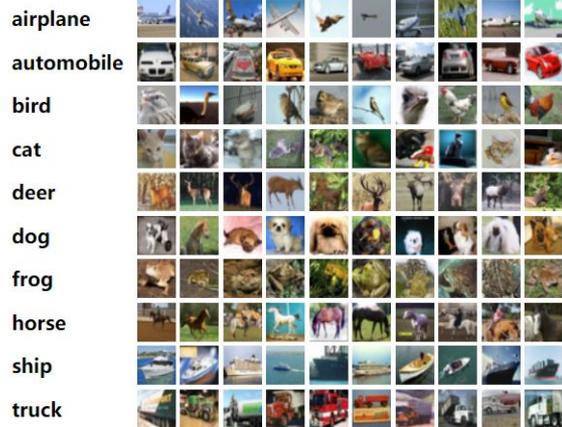
Convolution Layer



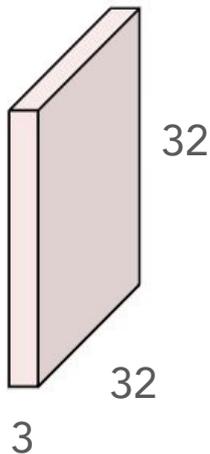
32x32x3 image → preserve spatial structure



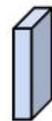
Convolution Layer



32x32x3 image

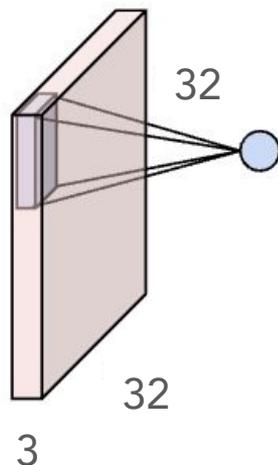
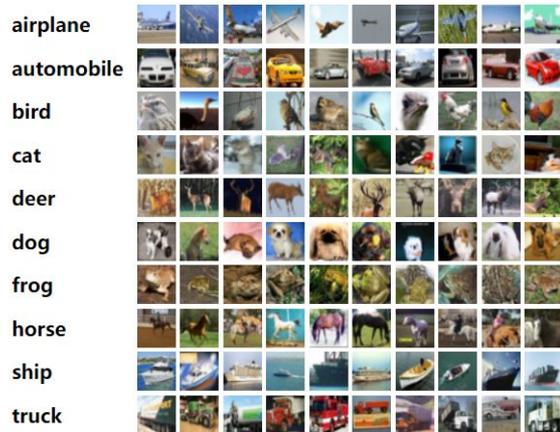


5x5x3 filter



Convolve the filter with the image, i.e., "slide over the image spatially, computing dot product"

Convolution Layer



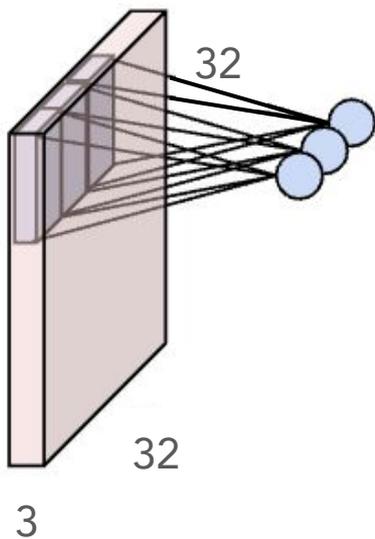
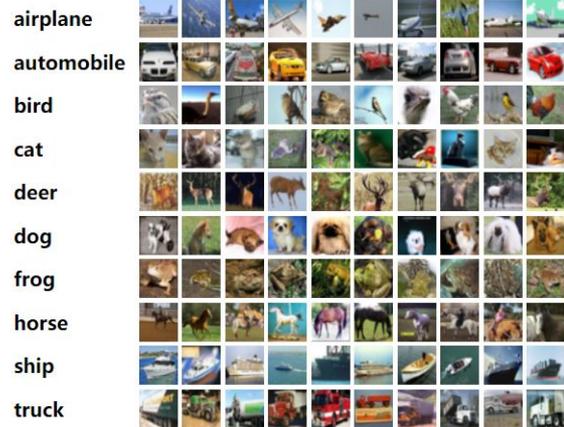
32x32x3 image
5x5x3 filter w

1 number:

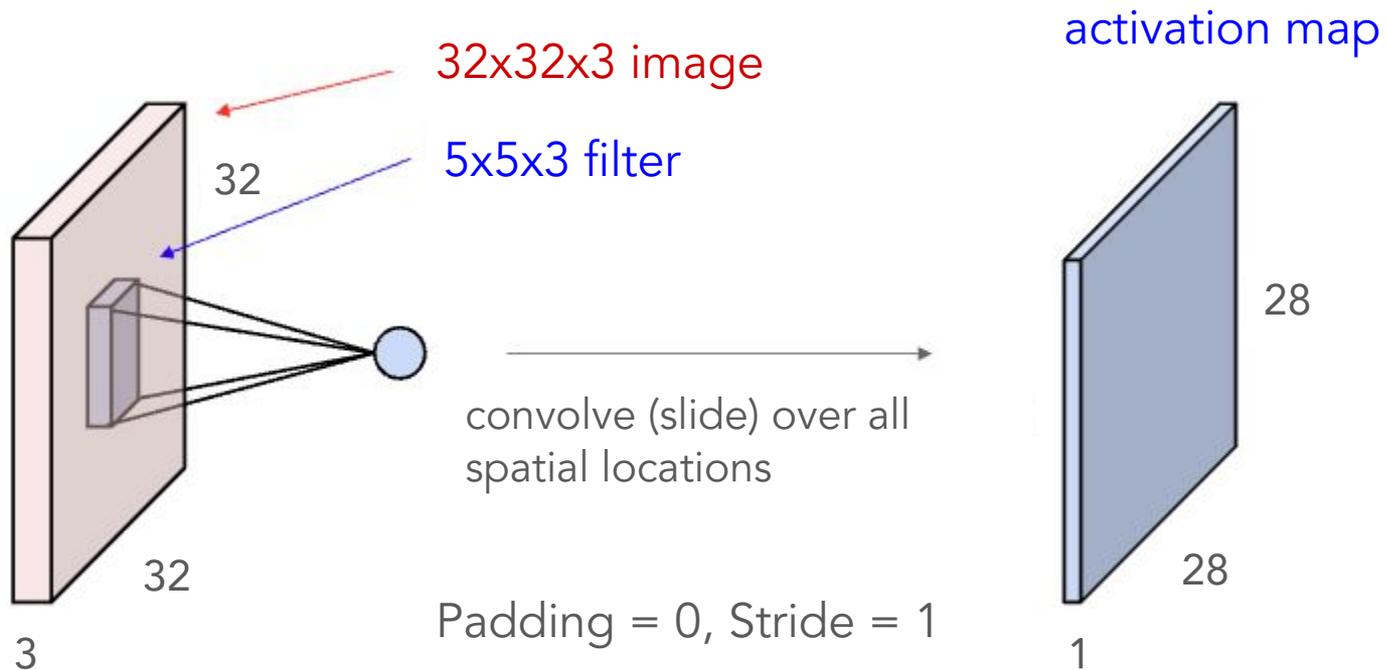
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. $5*5*3 = 75$ -dimensional dot product + bias)

$$w^T x + b$$

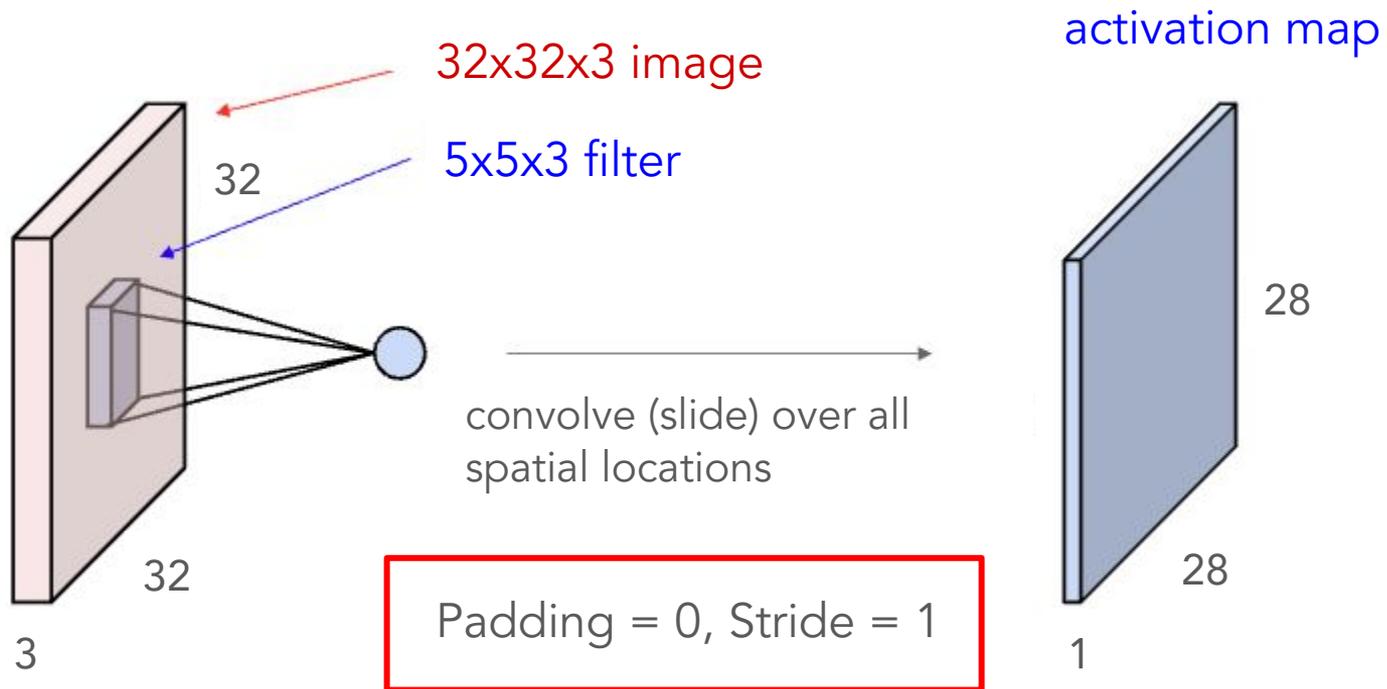
Convolution Layer



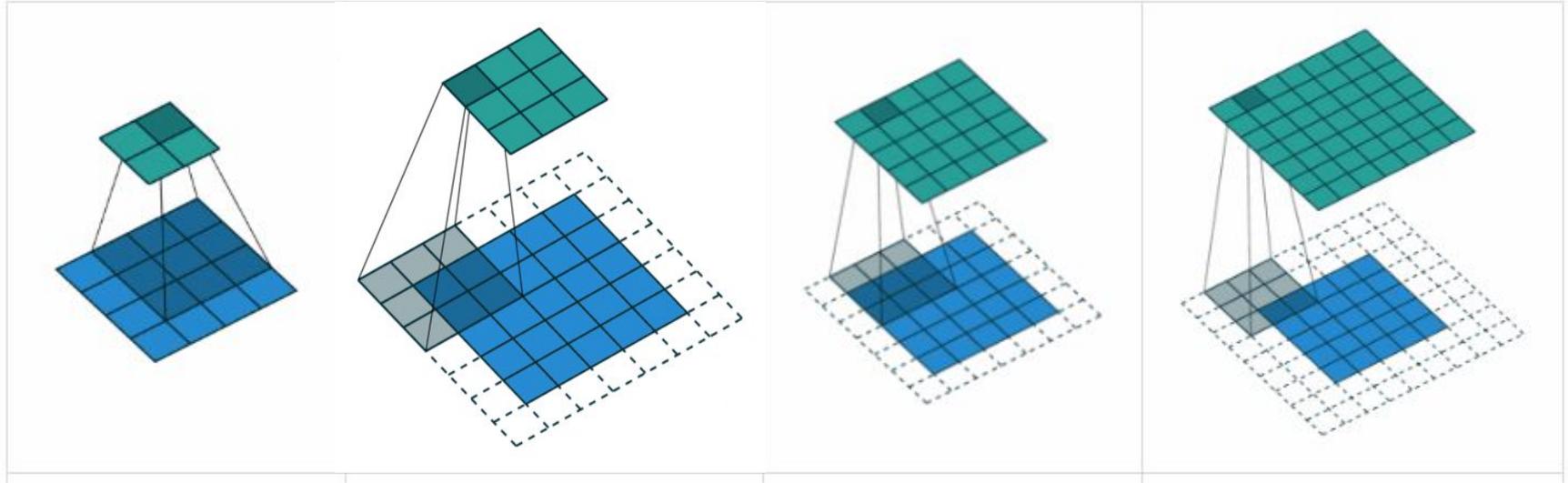
Convolution Layer



Convolution Layer



Convolution Layer



padding = 0, stride = 1

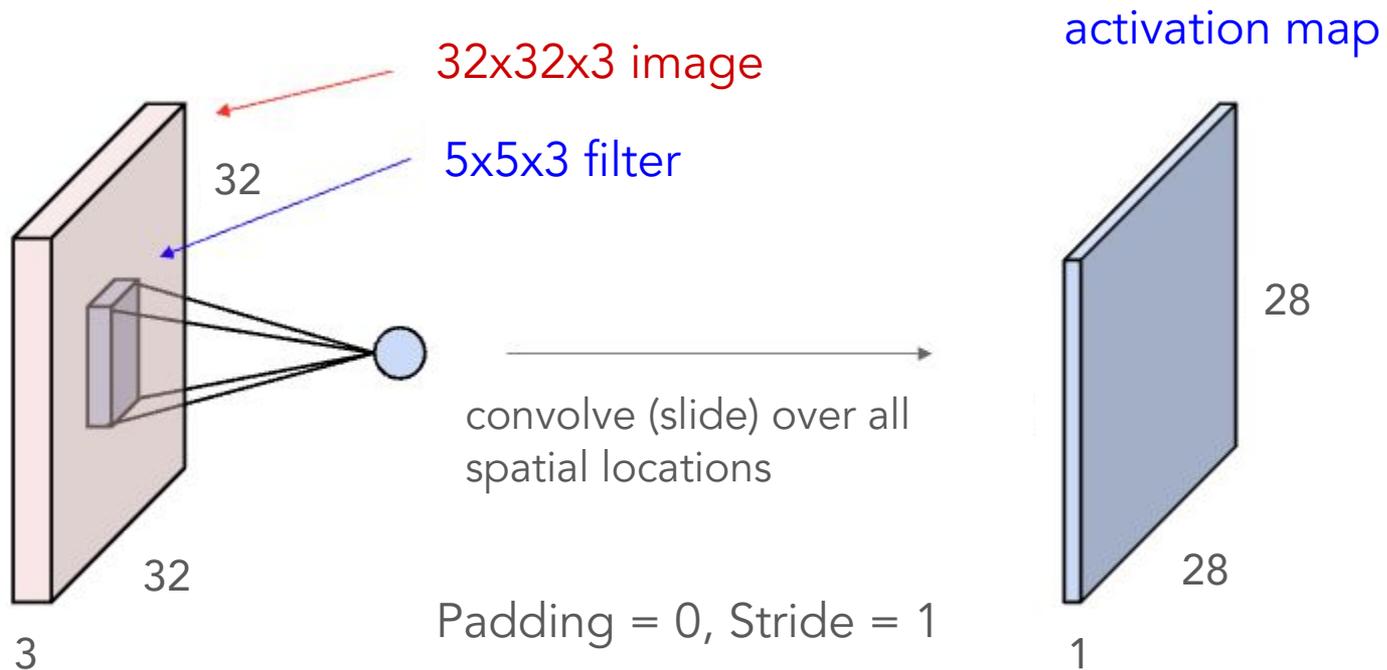
padding = 1, stride = 2

padding = 1, stride = 1

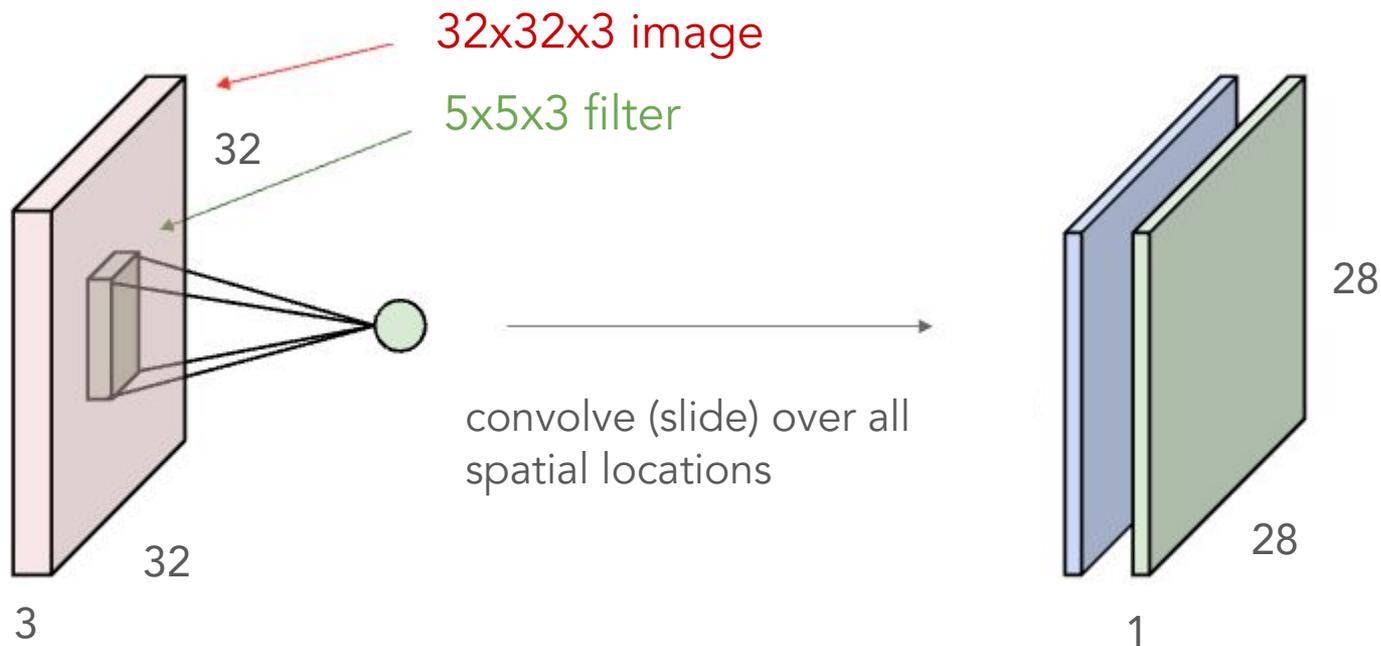
padding = 2, stride = 1

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

Convolution Layer

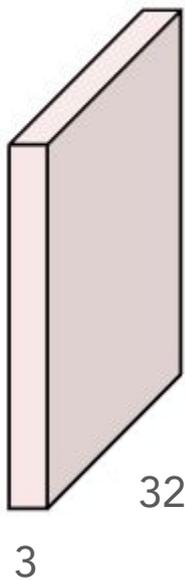


Convolution Layer



Convolution Layer

32x32x3 image



32

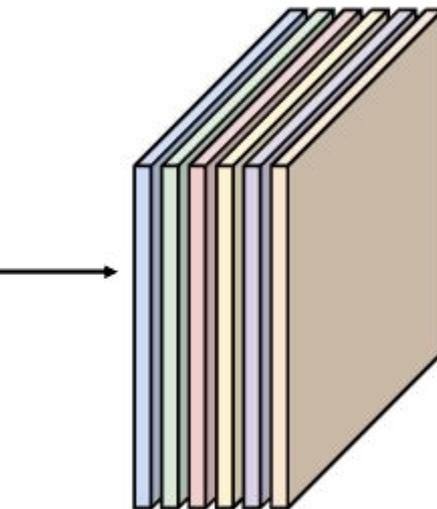
32

3

Consider 6 filters,
each 5x5x3



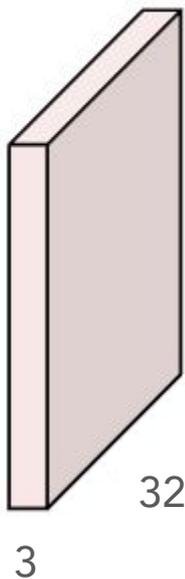
6x5x5x3
filters



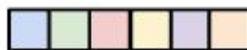
Stack activations to get a
6x28x28 output image!

Convolution Layer

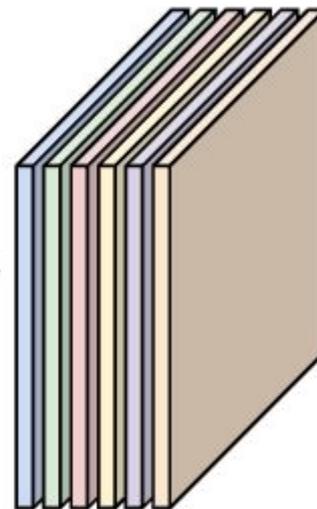
32x32x3 image



Don't forget bias terms!



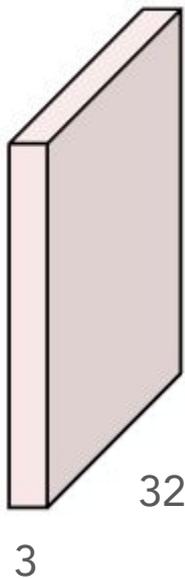
6x3x5x5 filters



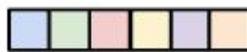
Stack activations to get a 6x28x28 output image!

Convolution Layer

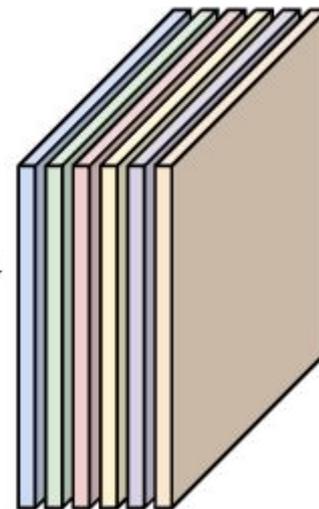
32x32x3 image



Don't forget bias terms!



6x3x5x5 filters

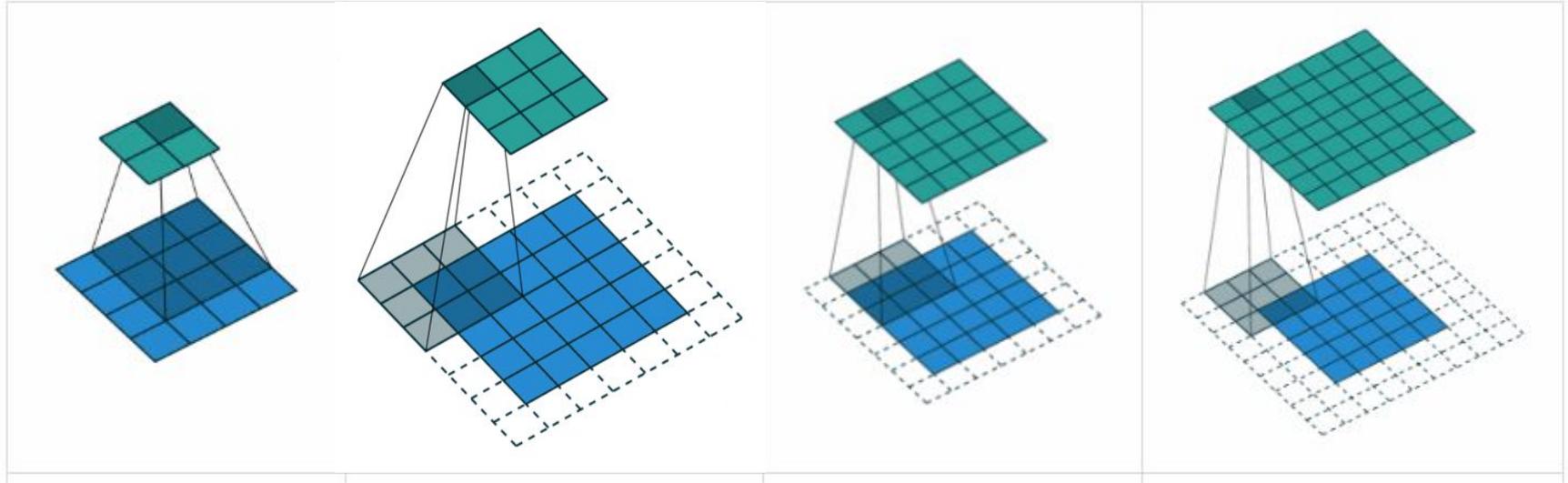


Activation Function!

ReLU

Stack activations to get a 6x28x28 output image!

Convolution Layer



padding = 0, stride = 1

padding = 1, stride = 2

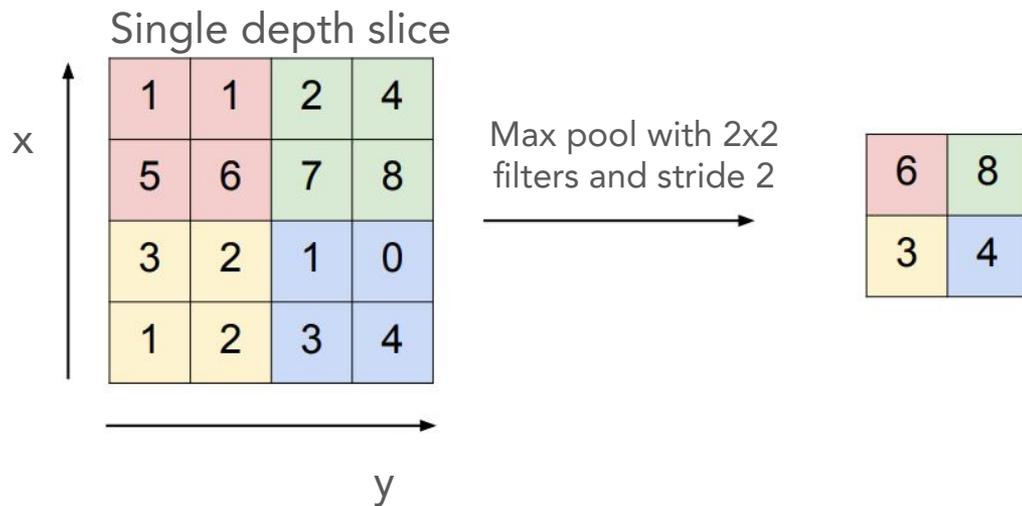
padding = 1, stride = 1

padding = 2, stride = 1

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

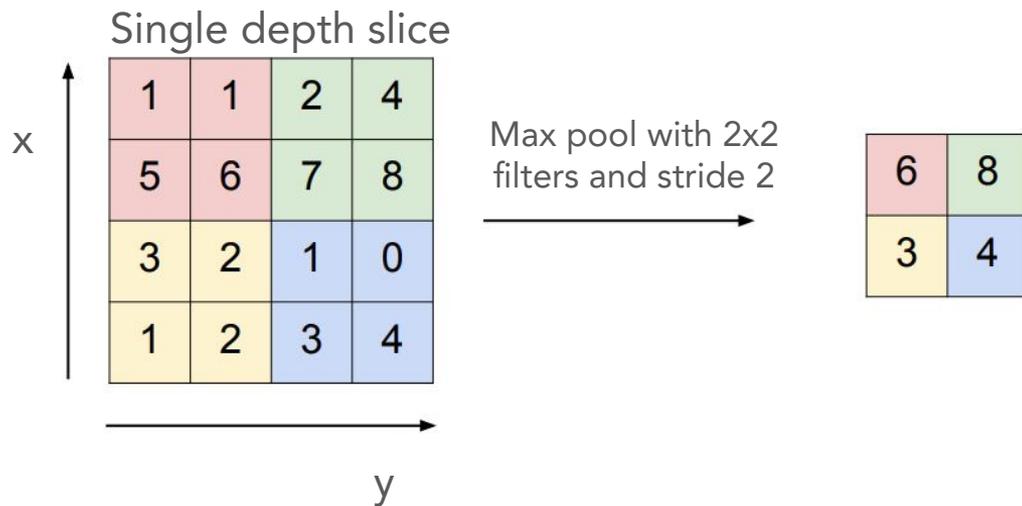
Pooling (Subsampling)

MAX POOLING



Pooling (Subsampling)

MAX POOLING

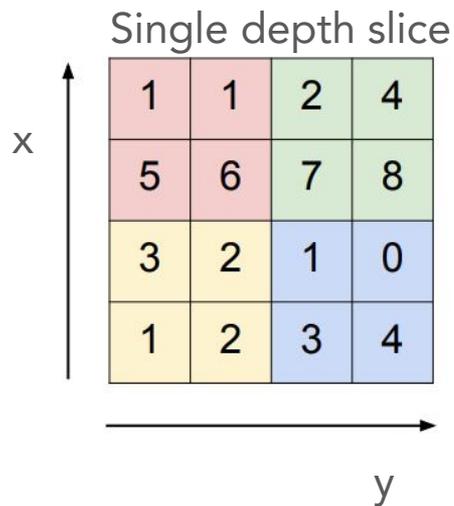


$$(W_x - F) / \text{Stride} + 1$$

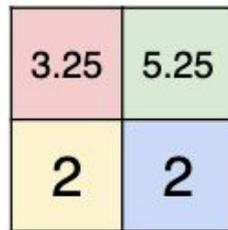
$$(W_y - F) / \text{Stride} + 1$$

Pooling (Subsampling)

MEAN POOLING



Mean pool with 2x2
filters and stride 2

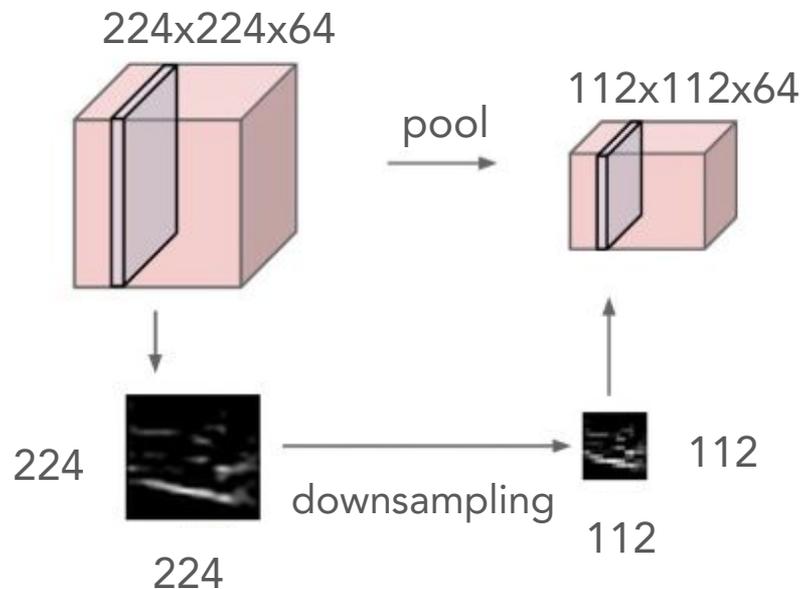


$$(W_x - F) / \text{Stride} + 1$$

$$(W_y - F) / \text{Stride} + 1$$

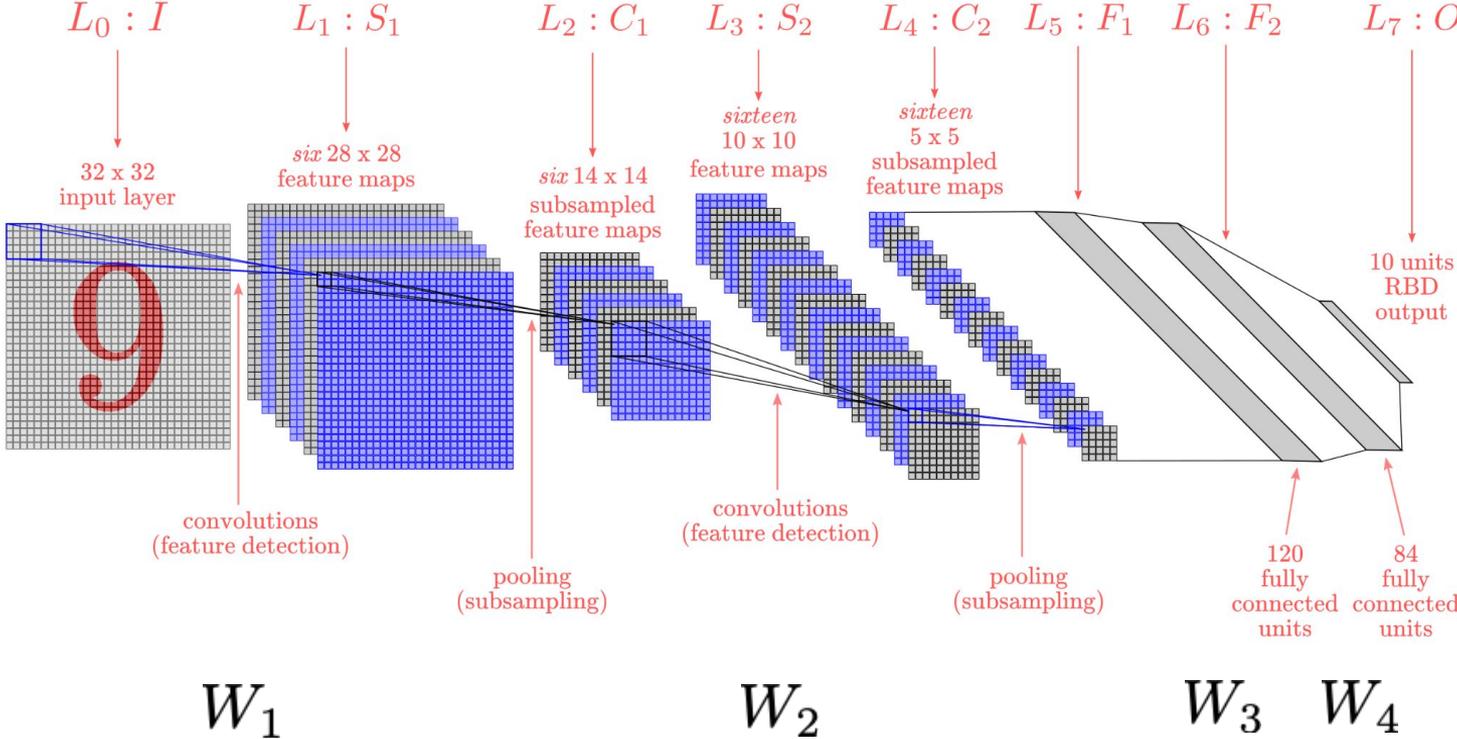
Average
Pooling

Pooling (Subsampling)



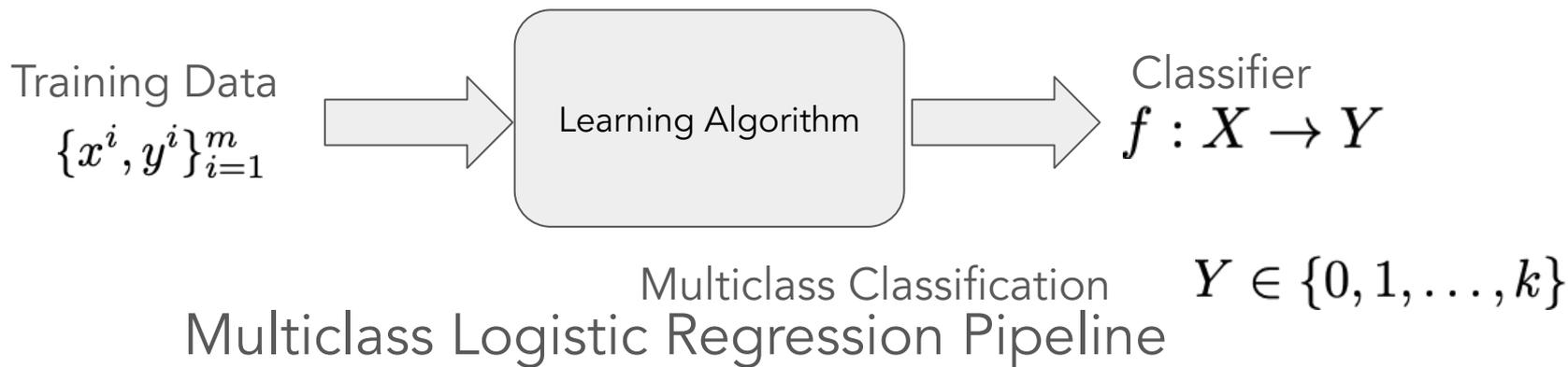
- Pooling layers simplify / subsample / compress the information in the output from the convolutional layer
- Reduce parameters

Put Everything Together



[LeNet-5, LeCun 1980]

Multiclass Logistic Regression Algorithms



1. Build probabilistic models:
Categorical Distribution + Conv NN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Select Optimizer

$$L(\theta) = \sum_{i=1}^m \ell(x^i, y^i, \theta) + \lambda \Omega(\theta)$$

$$\theta = [W_1, W_2, W_3, W_4]$$

$$\ell(x^i, y^i, \theta) = - \sum_{j=1}^k y^j \log \frac{\exp(o(V_j g(Wx^i)))}{\sum_{c=1}^k \exp(o(V_c g(Wx^i)))}$$

$f(x^i; \theta)$

$f(x^i; \theta)$

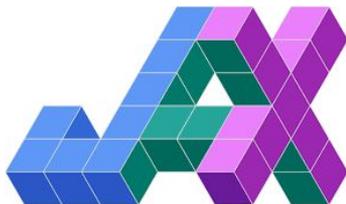
- (Stochastic) Gradient Descent

Auto-differentiation Packages

PyTorch



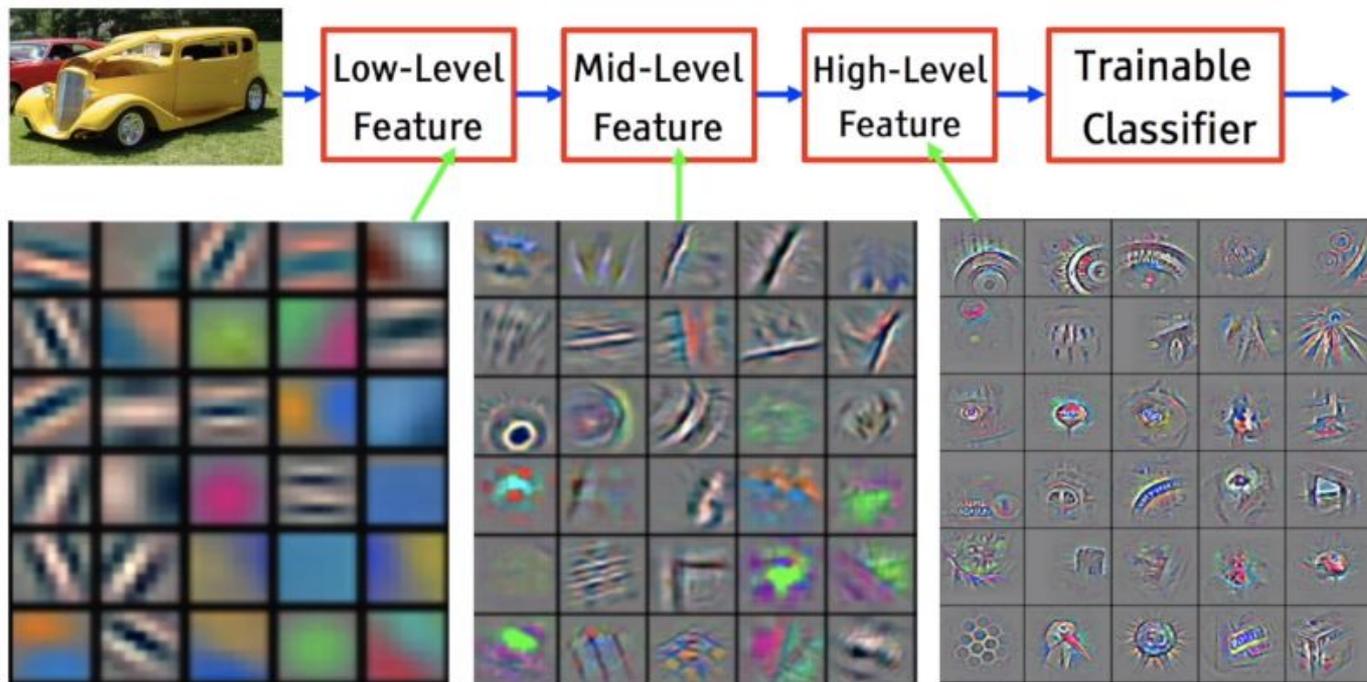
JAX



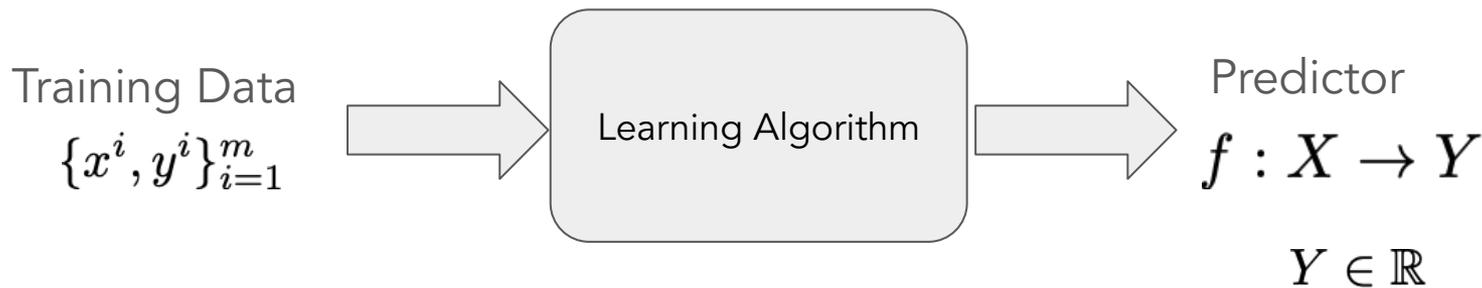
Tensorflow



Convolution Features



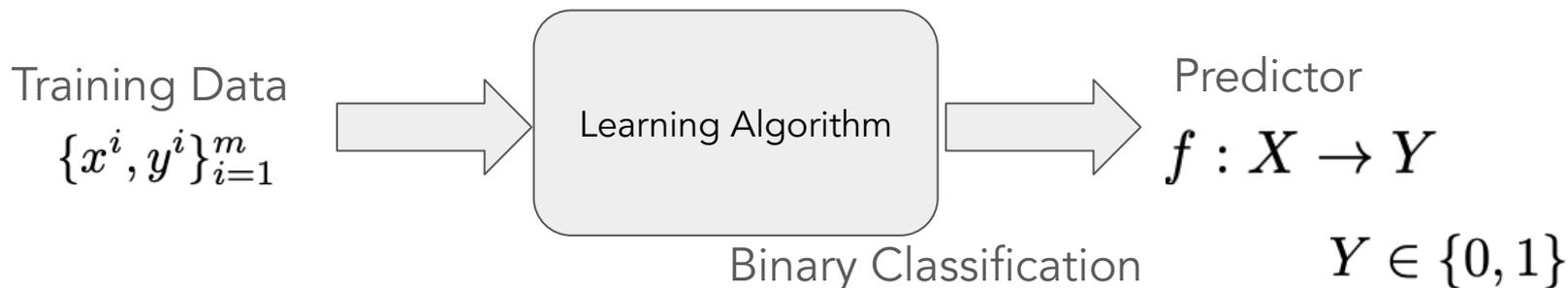
Regression Algorithms



Linear Regression Pipeline

1. Build probabilistic models:
Gaussian Distribution + Conv NN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

Binary Classification Algorithms

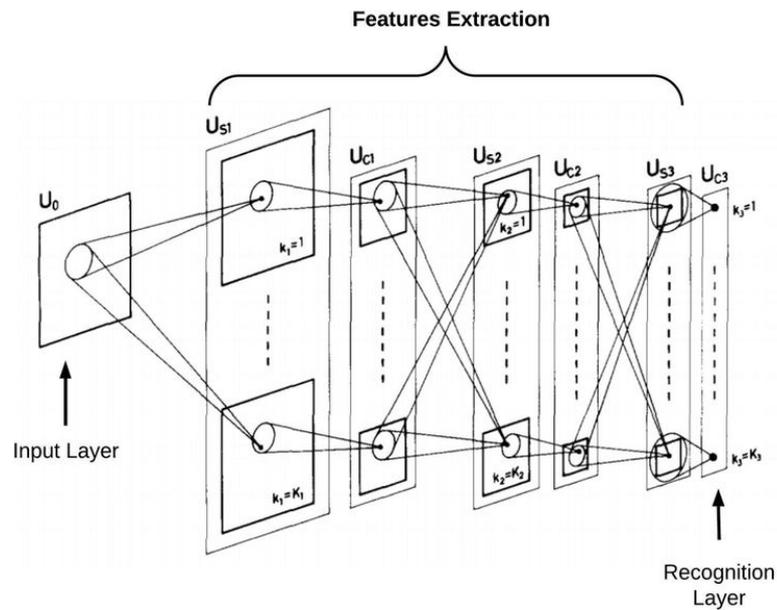


Binary Logistic Regression Pipeline

1. Build probabilistic models:
Bernoulli Distribution + Conv NN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

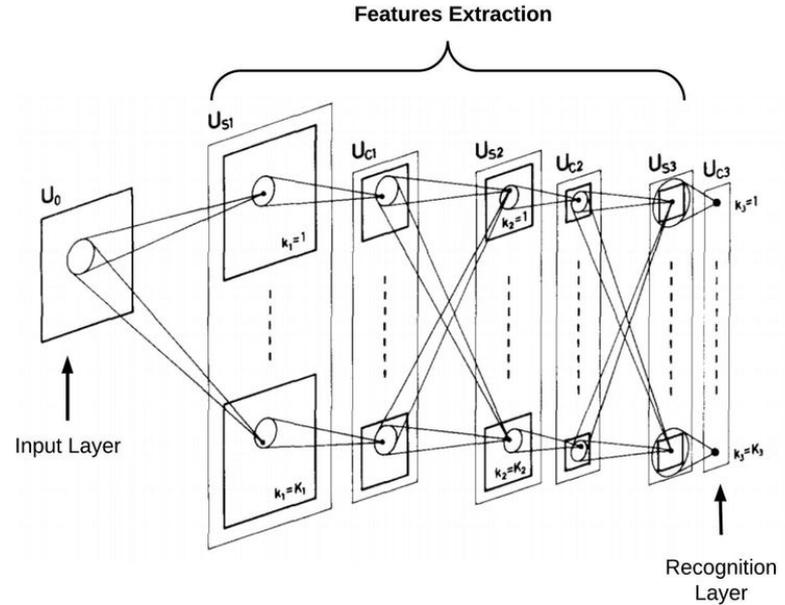
A Bit of History

- [Neocognitron](#) [Fukushima 1980]



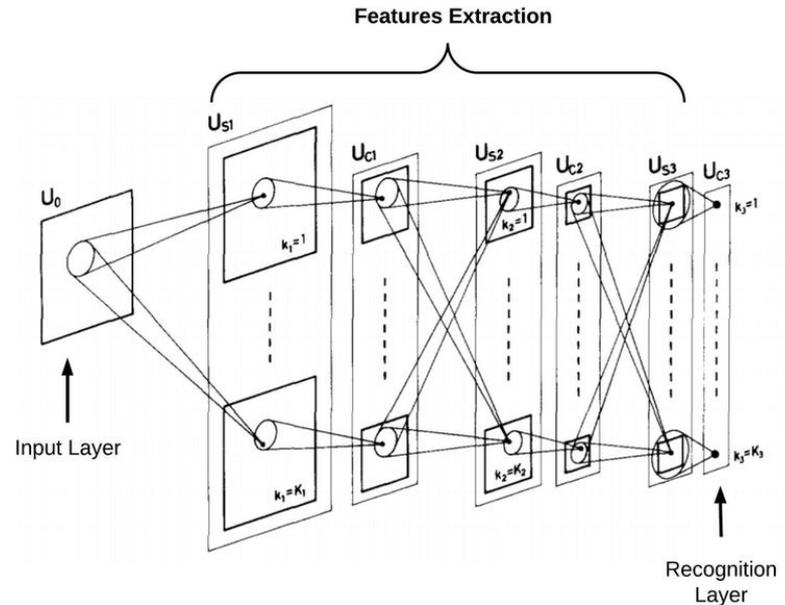
A Bit of History

- Neocognitron [[Fukushima 1980](#)]
- Backpropagation [[Rumelhart et al., 1986](#), [Werbos 1974](#)]



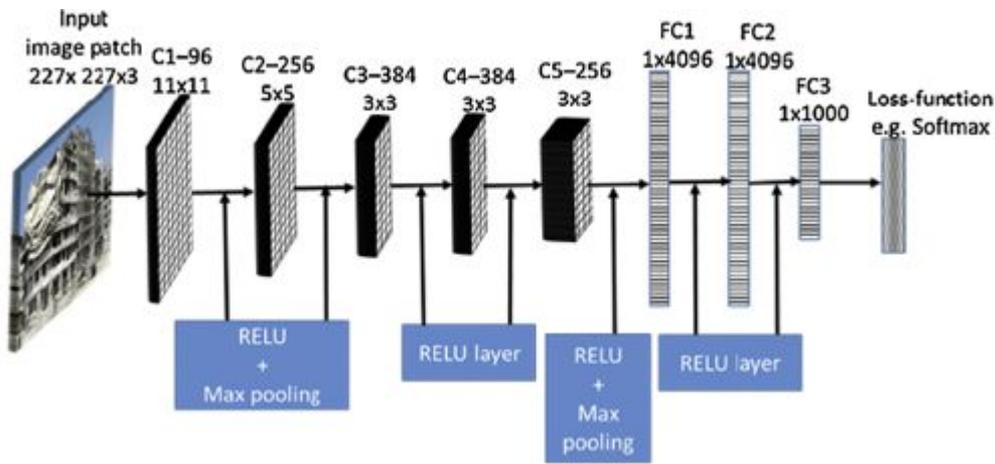
A Bit of History

- Neocognitron [[Fukushima 1980](#)]
- Backpropagation [[Rumelhart et al., 1986](#), [Werbos 1974](#)]
- BP on CNN [[LeCun et al., 1989](#)]

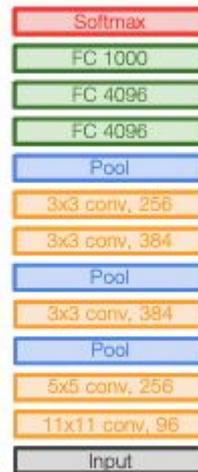
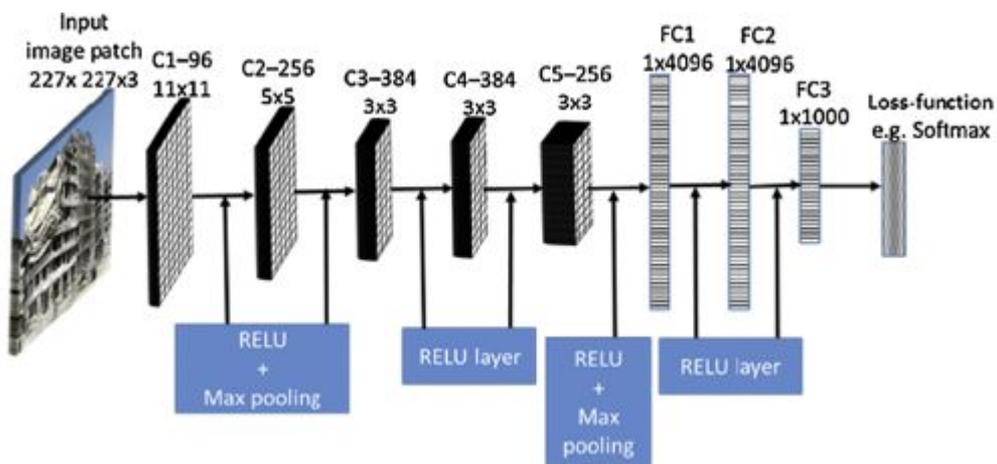


AlexNet (2012)

- AlexNet was one of the first deep convolutional on ImageNet competition with an accuracy of **84.7%** as compared to the second-best with an accuracy of **73.8%**.
- The activation used is the Rectified Linear Unit (ReLU).
- The training set had 1.2 million images. It was trained for 90 epochs, which took **five to six days** on two NVIDIA GTX 580 3GB GPUs.



AlexNet (2012)



AlexNet

Image Classification on ImageNet

Leaderboard

Community Models

Dataset

View

Top 1 Accuracy



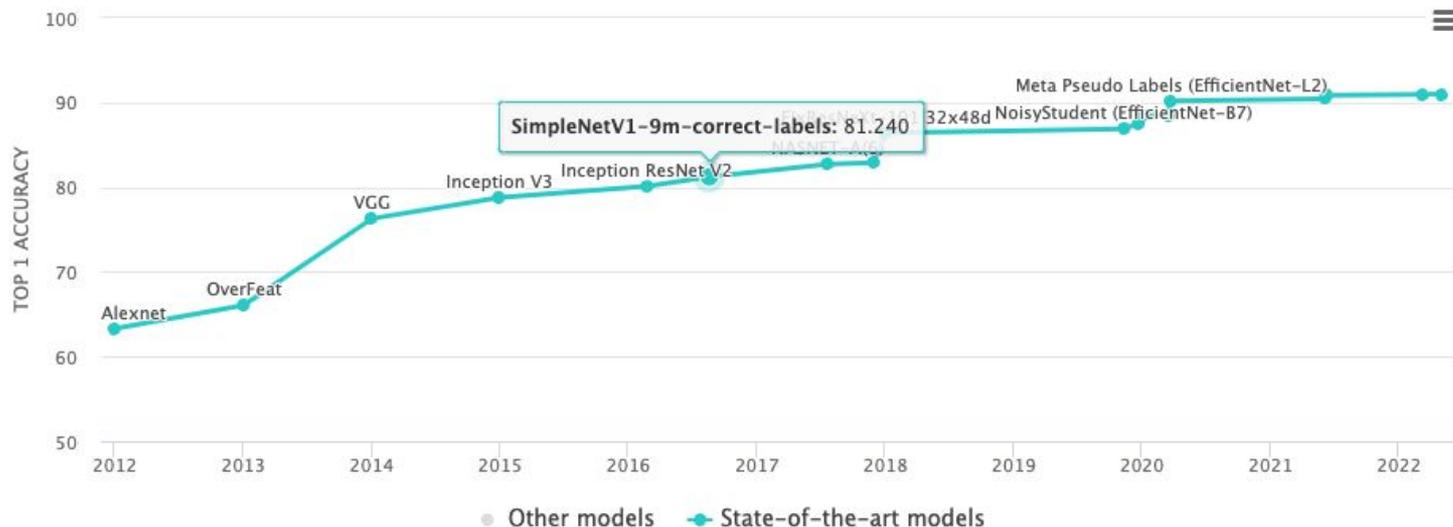
by

Date



for

All models



Q&A

HW2 Due Today!

No HW until Mar. 11th
Use the time wisely for your projects