

EN.601.482/682 Machine Learning: Deep Learning

Recap

Silvio Amir

samir@jhu.edu

Department of Computer Science
Johns Hopkins University

Image Classification



5.1 cat
2.3 car
-0.7 frog

1
0
0

One-hot encoding

Image Classification

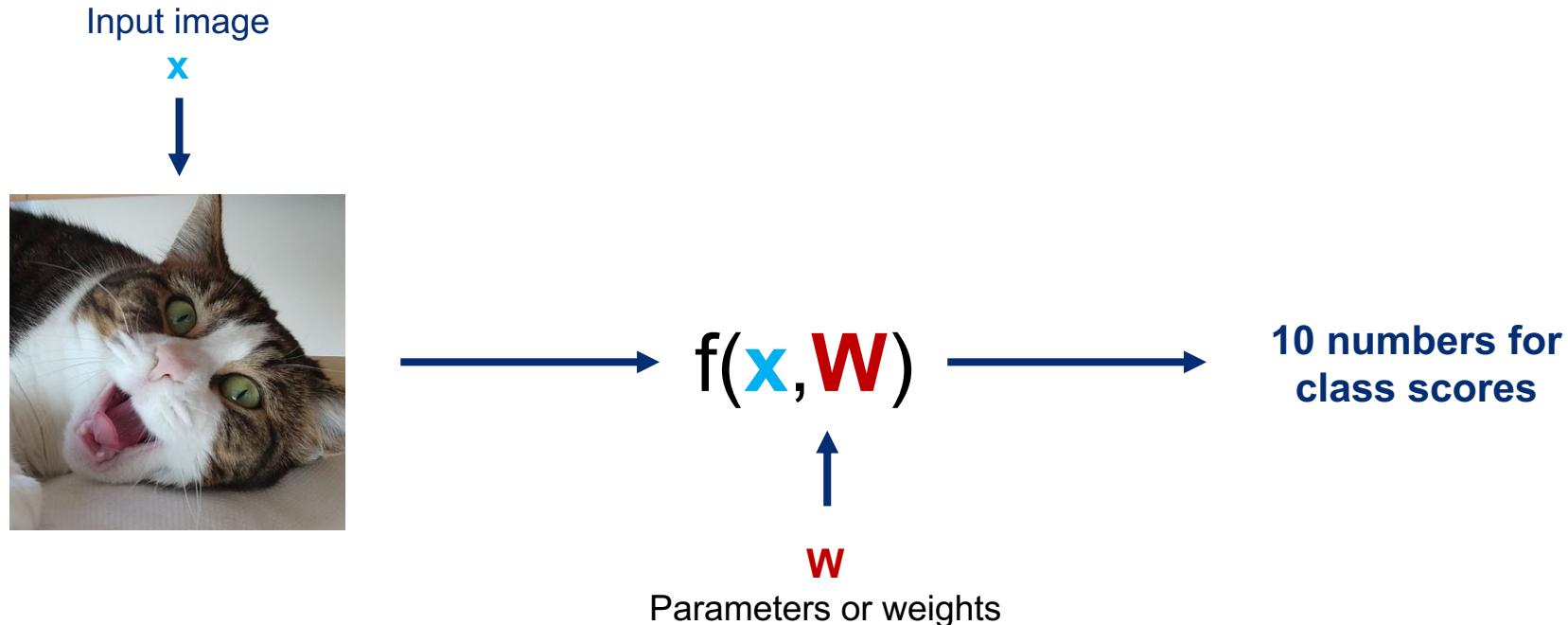


Image Classification

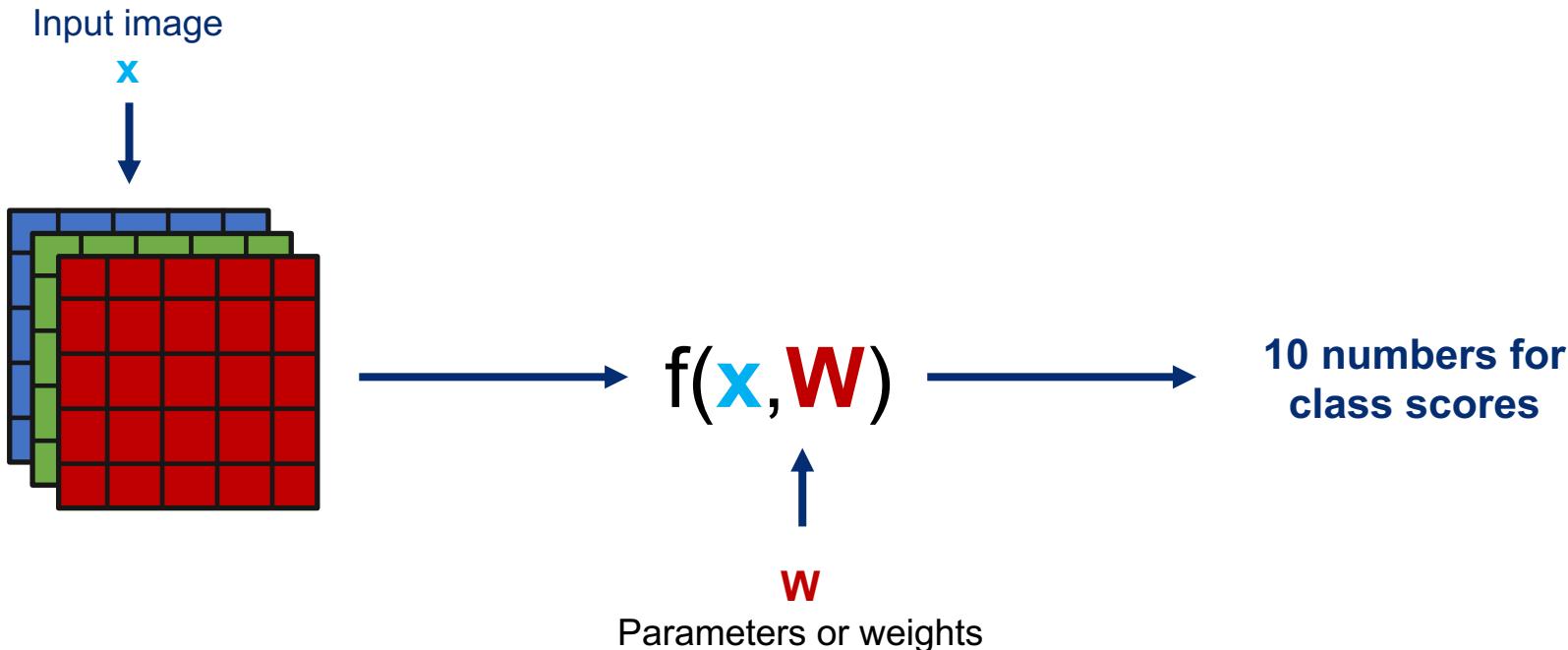
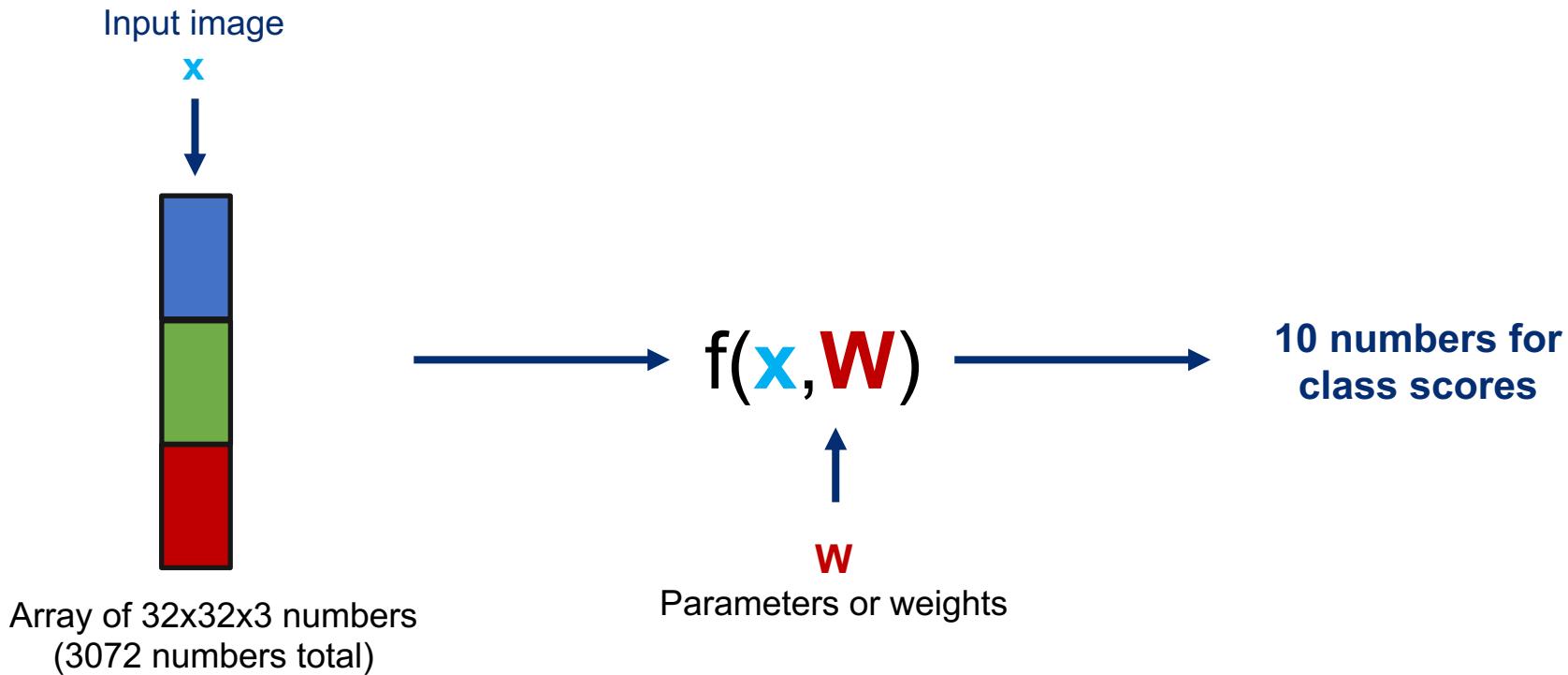
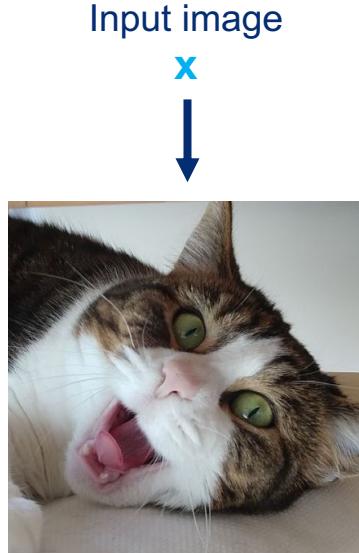


Image Classification



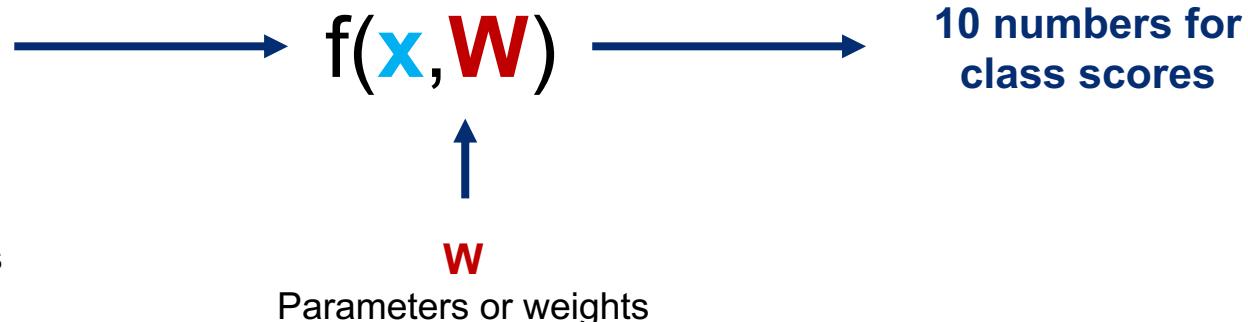
Linear Classifier



Array of $32 \times 32 \times 3$ numbers
(3072 numbers total)

$$f(x, W) = Wx + b$$

10 x 1 3072 x 1
10 x 3072 Bias term: constant!
 → Input independent preference



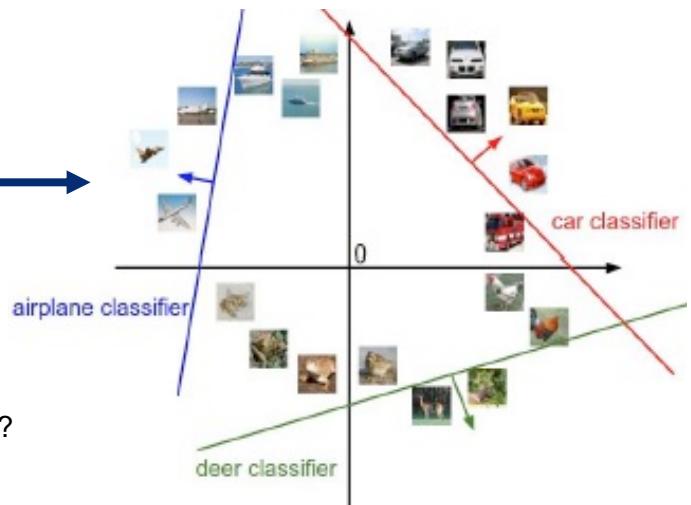
Training

Scores: $s = f(x_i, W) = Wx_i + b$



$$f(\mathbf{x}, \mathbf{W})$$

\mathbf{W}
How to estimate?



Loss Function

Loss function: quantifies classifier performance

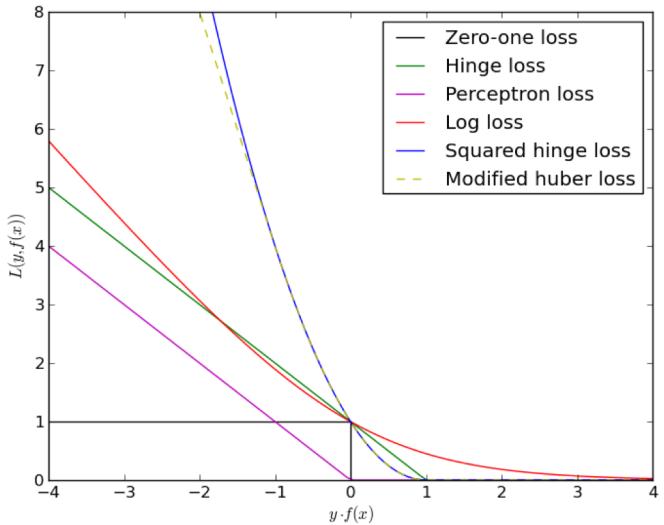
Given a dataset $\{(x_i, y_i)\}_{i=1}^N$ where

x_i is the image (or feature representation)

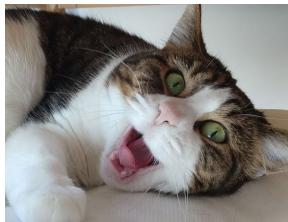
y_i is the corresponding label (integer)

Loss over the dataset: sum over all examples

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$



The Mean Squared Error



$$\text{Loss: } L_i = (y_i - s_i)^2$$

True label

compare

3.2

5

$$L_i = (5 - 3.2)^2 = 3.24$$

Negative Log Likelihood



True label for cat
(one-hot encoding)

Cat	3.2	1.00
Car	5.1	0.00
Bird	-1.7	0.00

These images are CC0 1.0 public domain.

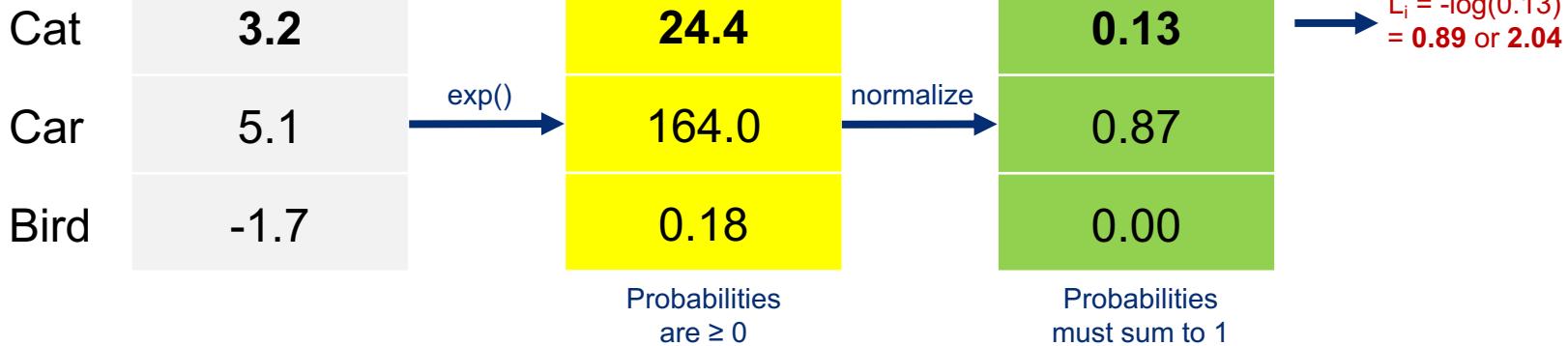


Negative Log Likelihood

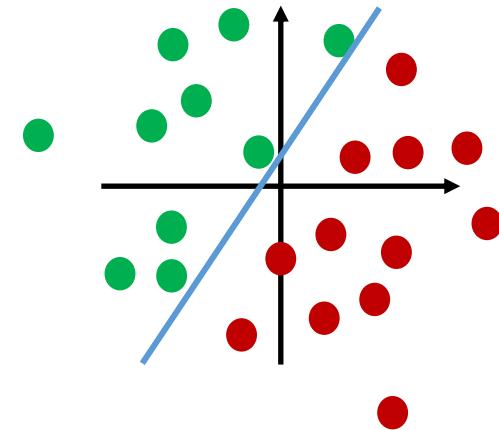
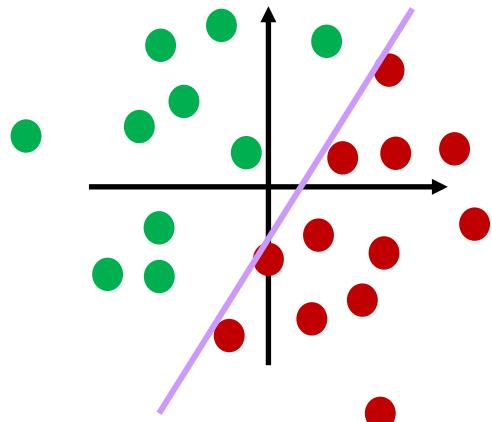
$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$



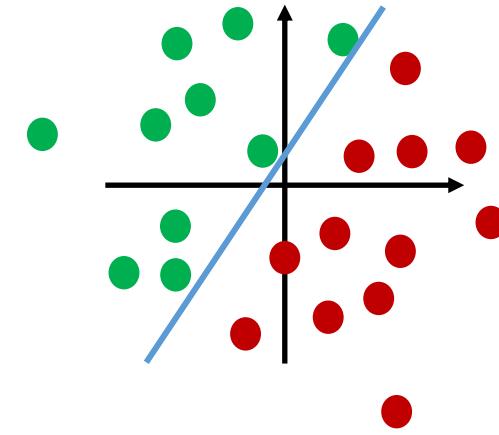
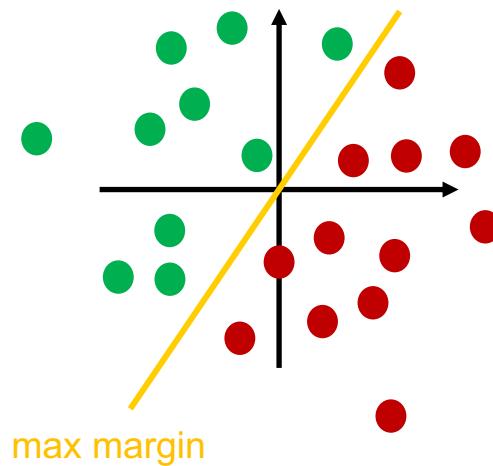
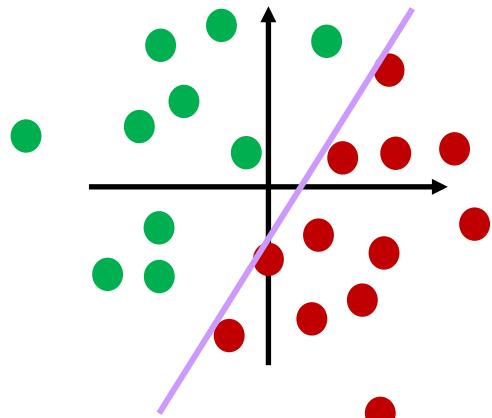
Loss: $L_i = -\log P(Y = y_i|X = x_i)$ → Log-likelihood of true class



Which W to choose?



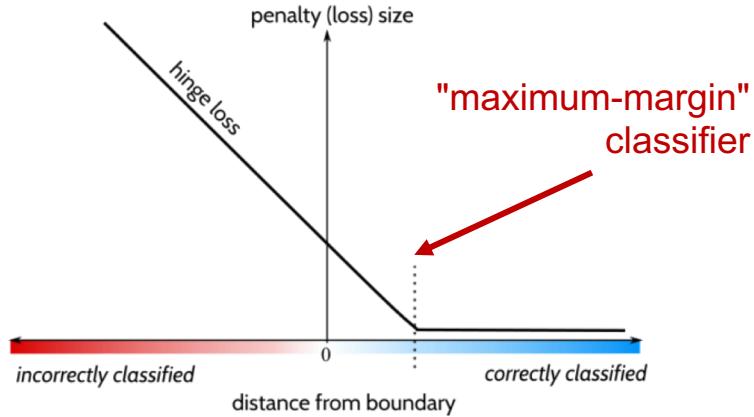
Which W to choose?



The Hinge (SVM) Loss



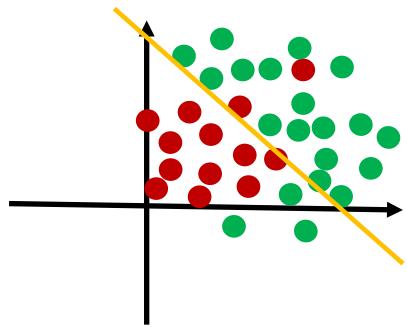
Cat	3.2	1.3	2.2
Car	5.1	4.9	2.5
Bird	-1.7	2.0	-3.1



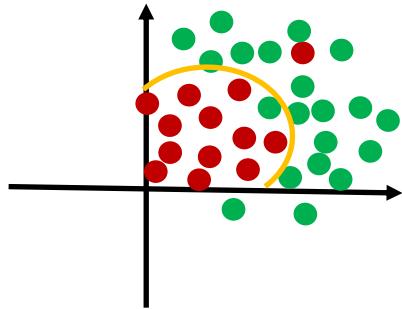
Loss:

$$\begin{aligned} L_i &= \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases} \\ &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \end{aligned}$$

Bias and Variance



High bias



High Variance

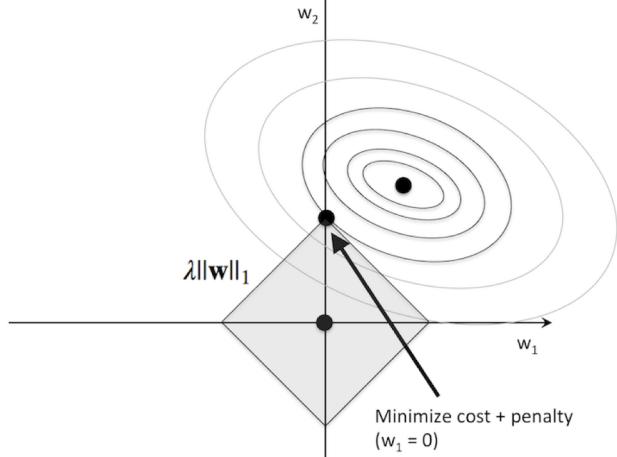
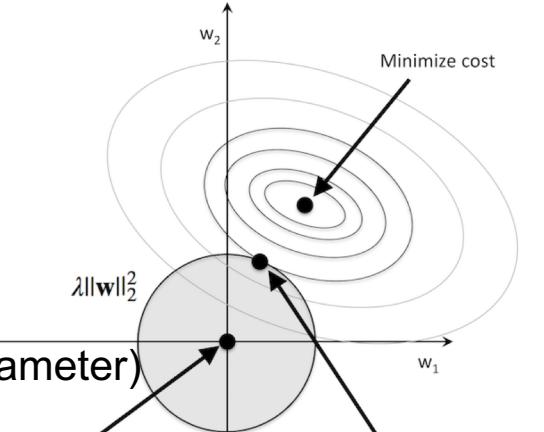
Norm Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_i L_i(f(x_i, W), y_i)}_{\text{Data term: Predictions must match annotations}} + \lambda R(W)$$

Regularization
 λ strength (hyperparameter)

Simple regularizers

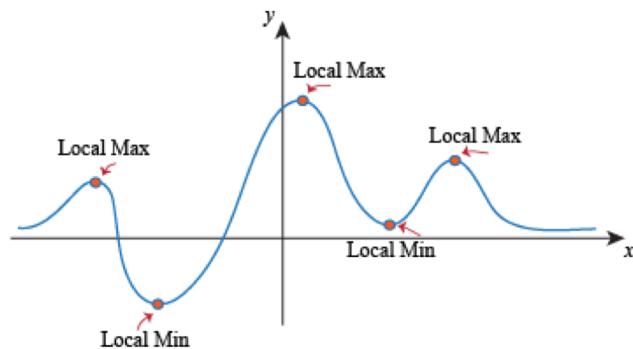
- L2 (magnitude): $R(W) = \sum_{k,l} W_{k,l}^2$
- L1 (sparsity): $R(W) = \sum_{k,l} |W_{k,l}|$
- Versions, e.g. Elastic net, Huber,...



Images from rasbt.github.io.

Gradient Descent

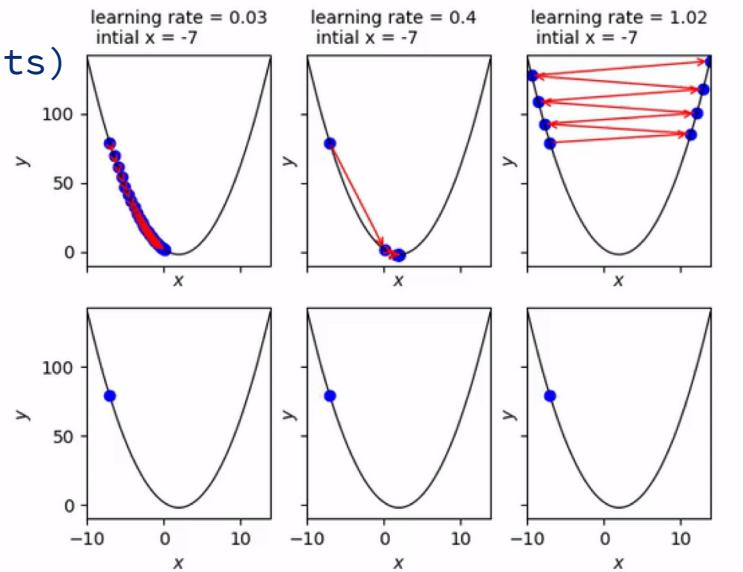
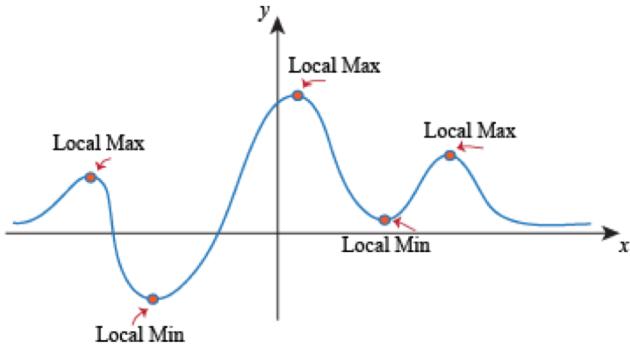
Minimize the loss function with gradient descent



Gradient Descent

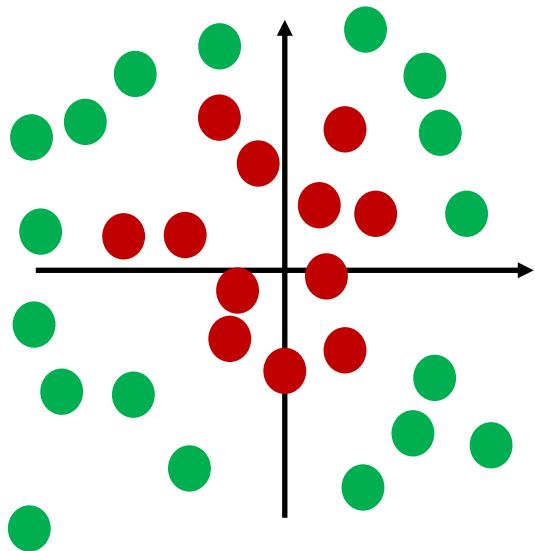
Minimize the loss function with gradient descent

```
while not_converged:  
    gradient = eval_gradient(loss,data,weights)  
    weights += - l_rate * gradient
```



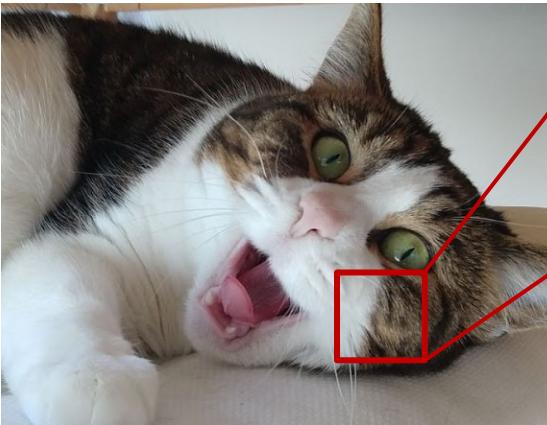
[Image from jed-ai.github.io](https://jed-ai.github.io). Too large step sizes can cause divergence!

Real World Data



Challenges of Image Data

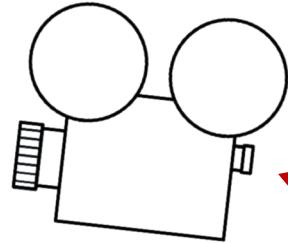
Why is this challenging?



What the human sees

1185	132	189	332	594	99	186	89	96	163	132	119	104	97	93	87		
81	98	182	388	594	79	96	183	99	165	123	138	120	105	94	85		
79	89	98	385	528	105	87	98	95	99	133	112	109	105	95	85		
89	91	81	83	520	581	127	388	99	98	102	98	96	98	101	93		
186	91	61	64	69	91	88	85	101	107	104	98	76	80	66	95		
114	124	95	25	95	66	64	65	64	80	122	128	74	84	91	91		
132	137	144	342	342	81	95	65	63	20	64	34	63	162	93	45	82	
129	137	144	348	599	95	86	79	62	65	60	63	58	73	64	101	91	
125	133	148	337	319	121	117	84	85	79	89	85	54	64	72	88	91	
127	125	131	347	533	127	126	211	111	98	89	75	61	84	72	84	91	
113	114	189	323	559	148	131	218	113	169	188	92	74	65	72	79	91	
83	93	98	97	593	147	133	218	113	114	132	189	196	95	77	99	91	
83	77	88	81	77	79	182	323	117	115	132	129	125	138	313	97	91	
132	65	82	88	78	71	88	381	128	125	118	181	187	131	131	109	91	
133	65	75	88	99	71	62	88	128	126	128	185	91	98	128	119	91	
137	65	71	87	106	66	66	45	76	126	126	187	92	94	105	112	91	
124	97	91	88	117	120	120	94	41	51	65	93	99	95	182	147	91	
134	94	97	91	112	120	120	120	124	384	76	88	94	181	93	91	184	91
157	179	157	128	93	98	214	232	112	97	69	55	79	93	88	91	91	
139	128	134	242	539	169	185	218	121	134	134	97	65	53	85	91	91	
129	132	96	96	337	559	144	128	219	104	167	182	93	87	81	72	79	91
123	187	96	88	83	112	153	249	522	169	186	75	28	187	312	99	91	
122	121	182	88	82	86	94	317	505	148	188	182	98	78	82	187	91	
122	164	148	182	71	56	38	84	93	168	119	139	182	61	68	85	91	

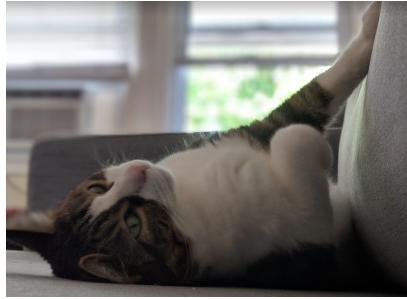
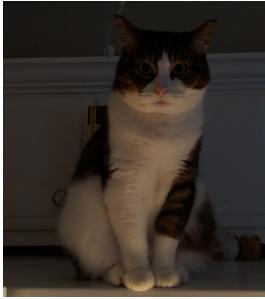
What the computer sees



When camera moves,
all pixel values change!

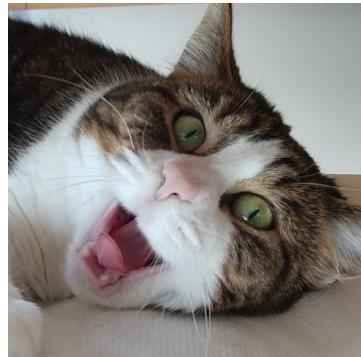
Challenges of Image Data

Why is this challenging?

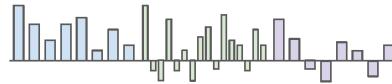


Features

Higher-level representation via feature extraction



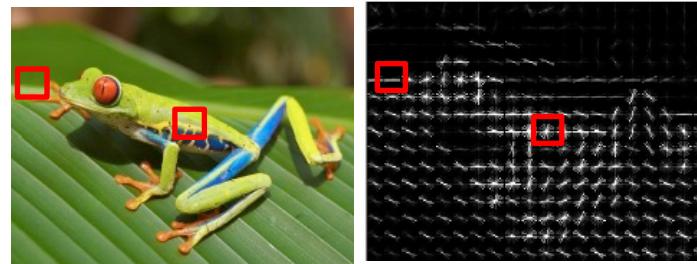
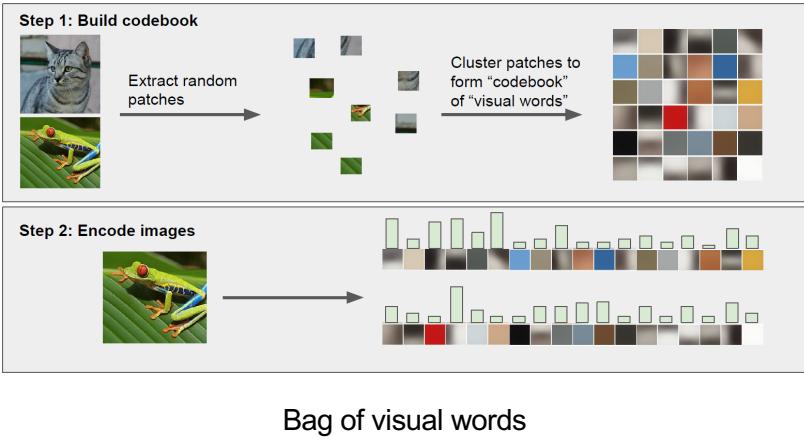
Feature
extractor



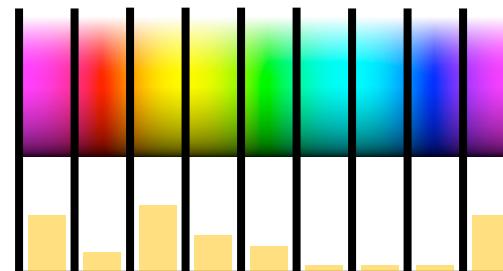
class scores

$$f(\mathbf{x}, \mathbf{W})$$

Features for Image Classification



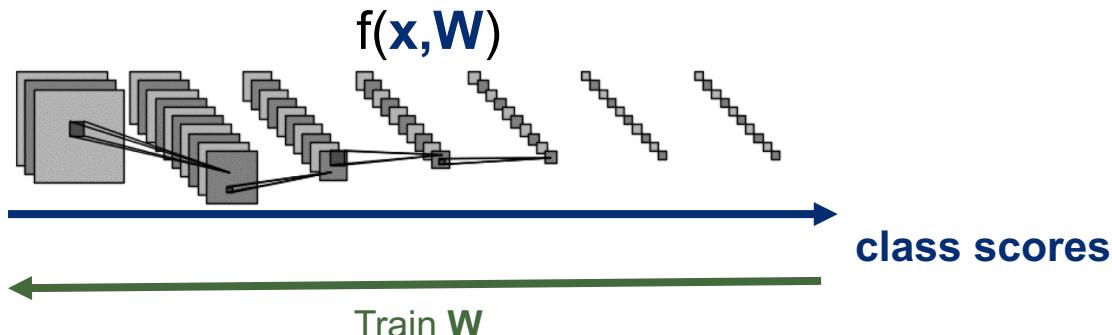
Oriented gradients



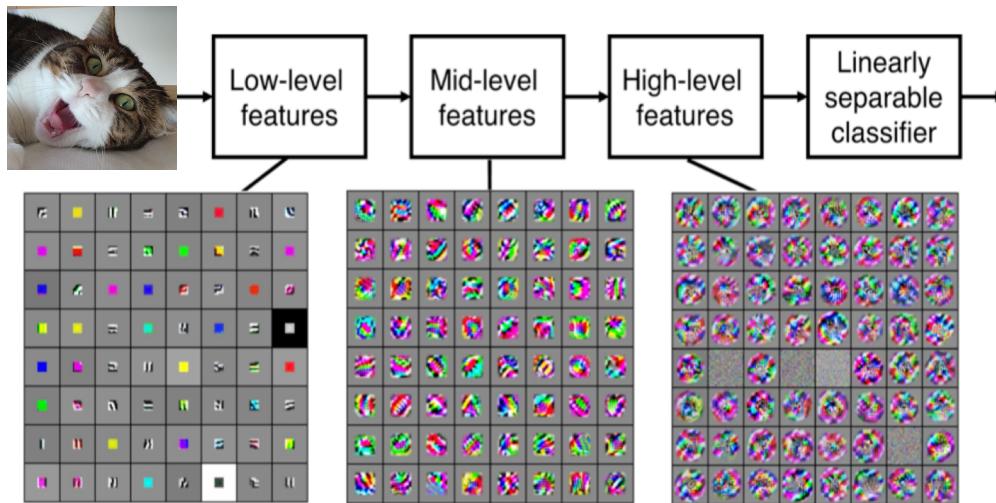
Color histogram



Learning Feature Extractors



Learning Feature Extractors



Recap

Neural Networks

Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

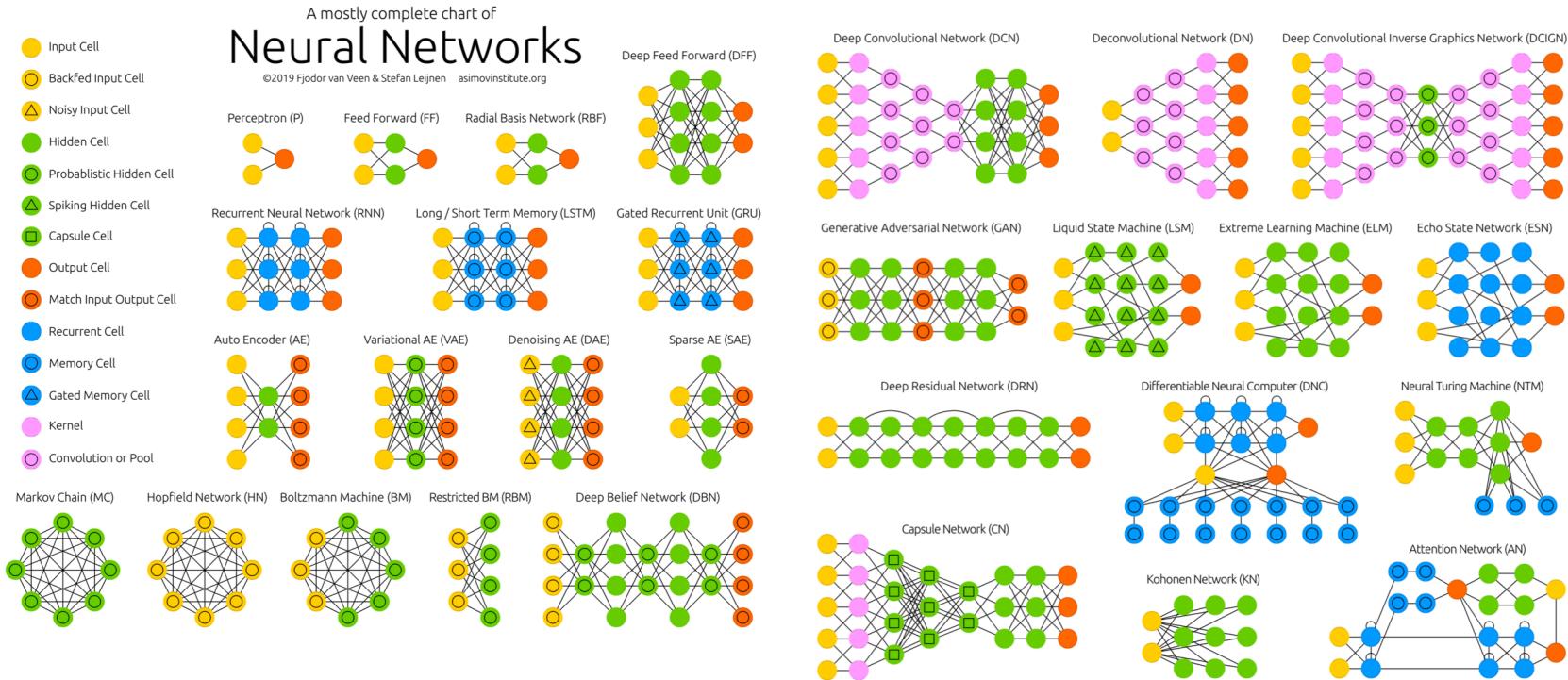
Softmax
Sigmoid

Optimization

Gradient Descent



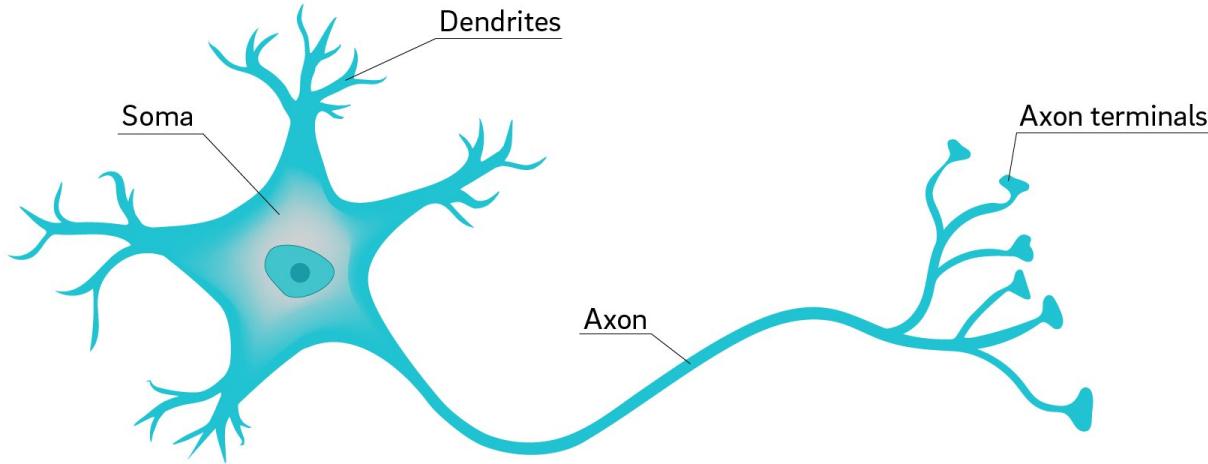
Neural Network Zoo (2016)



<https://www.asimovinstitute.org/neural-network-zoo/>

Neuron

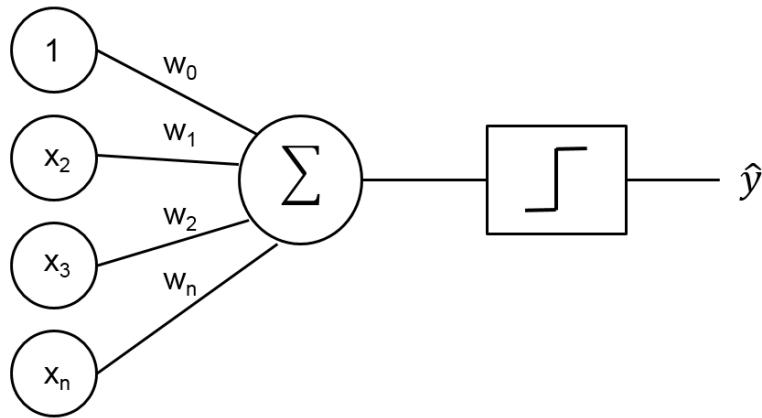
- Neural networks imply connections to neurons



- Signal comes in through dendrites into soma
- Then connects via the axon to other neurons
- Fires if input exceeds a certain threshold

Neuron

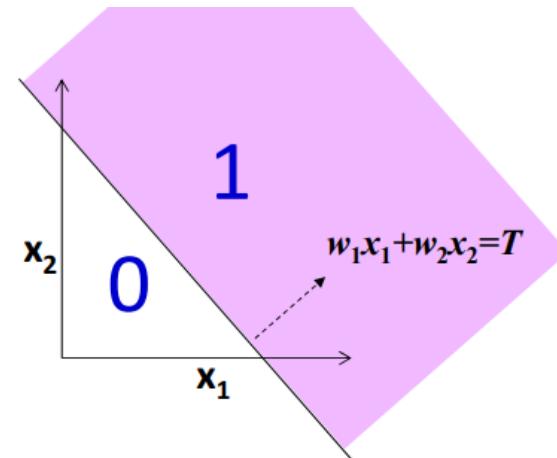
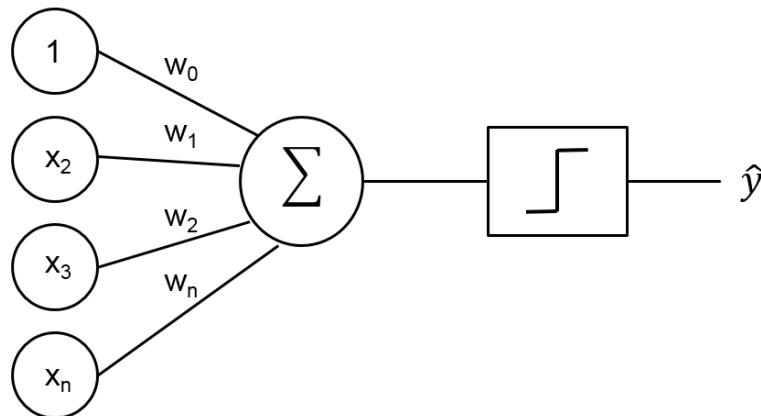
- Neural networks imply connections to neurons



- Signal comes in through dendrites into soma
- Then connects via the axon to other neurons
- Fires if input exceeds a certain threshold

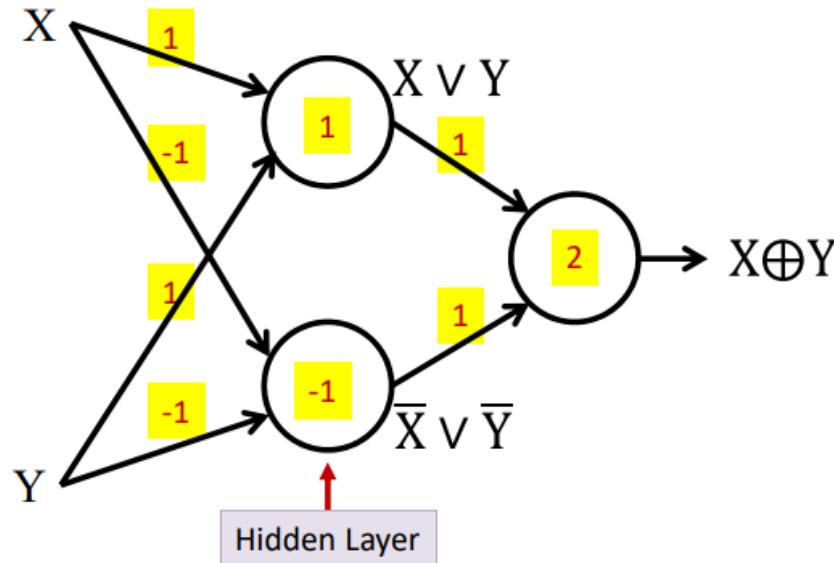
Perceptron

The perceptron can model almost all Boolean gates but in the end it is still a linear classifier



A Single Perceptron is Not Enough

Individual perceptrons are weak → Networked elements required!

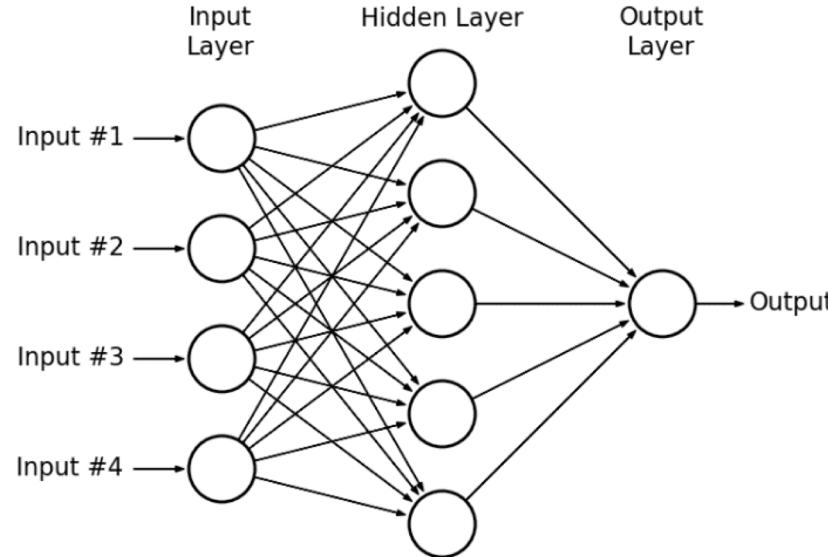


$(X \text{ or } Y) \text{ and } (\text{Not } X \text{ or } \text{Not } Y) \rightarrow \text{XOR}$

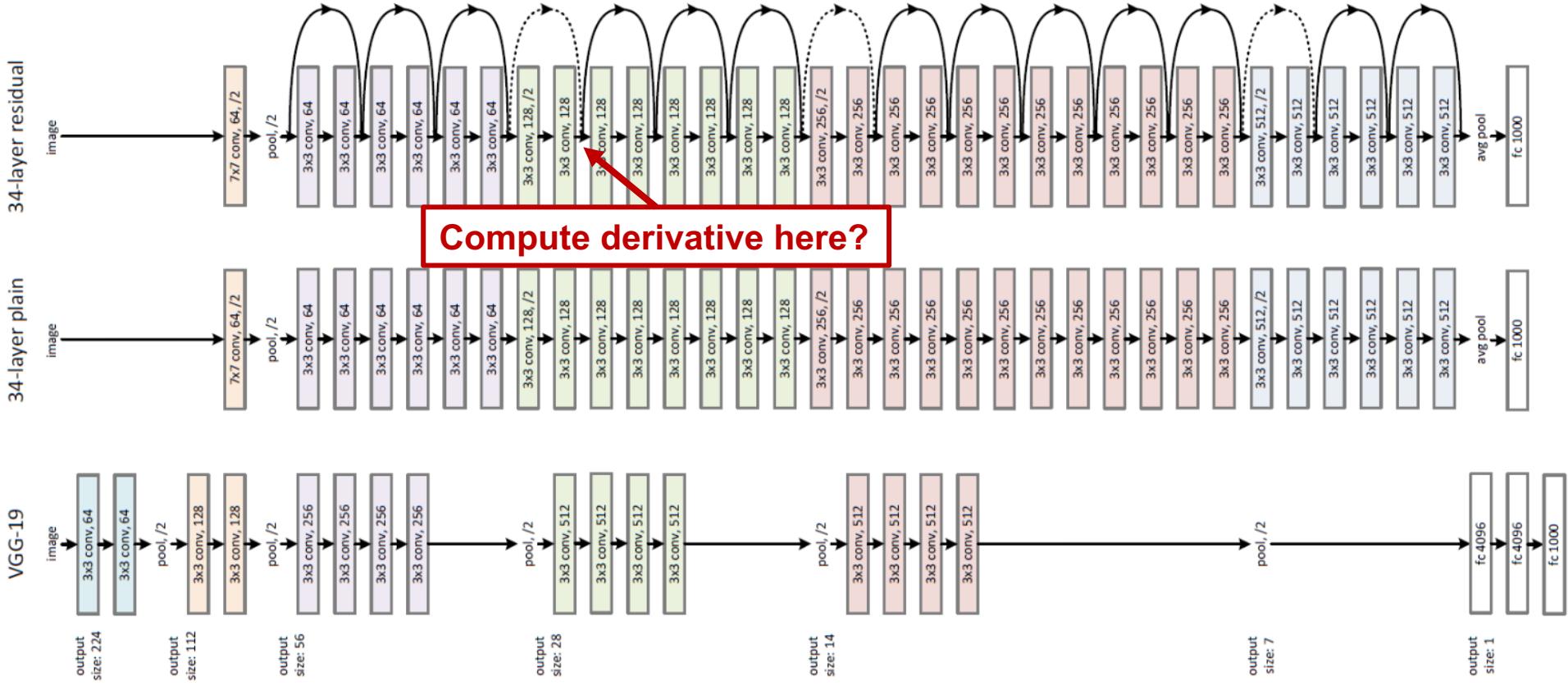
[Minsky, M., & Papert, S. \(1969\). Perceptron Expanded Edition.](#)

Multi-layer Perceptron

Multi-layer perceptrons can compose arbitrarily complex functions.



Computational Graphs

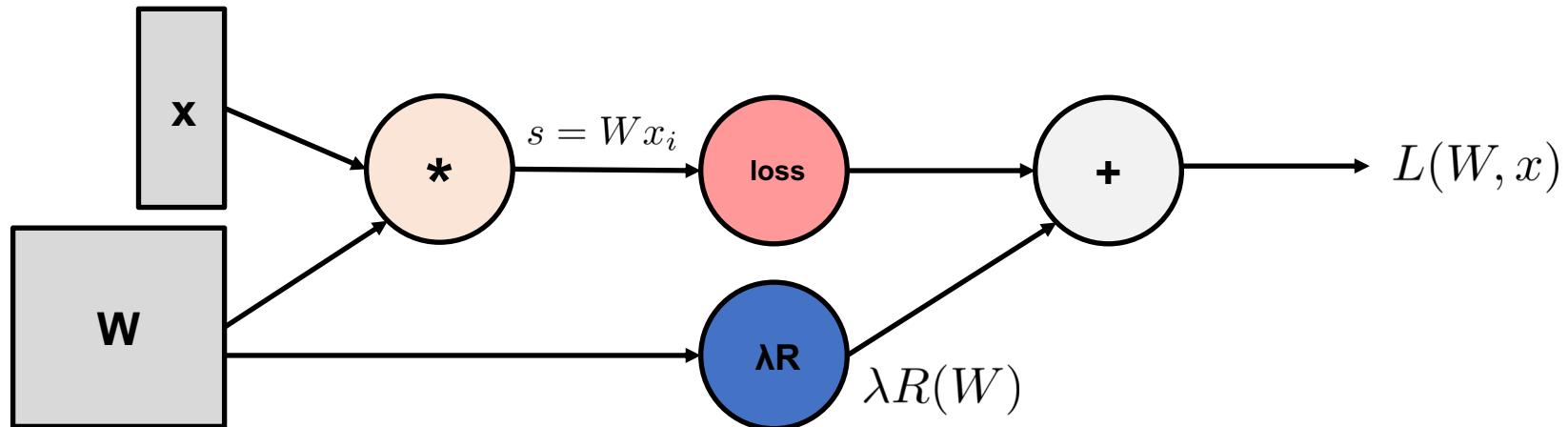


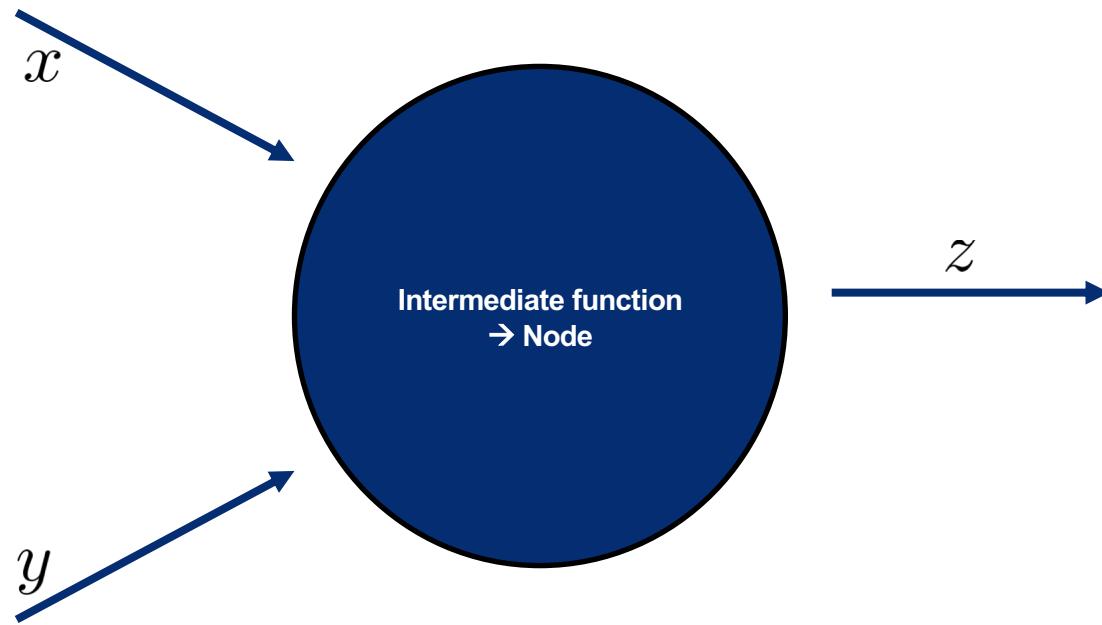
Computational Graphs

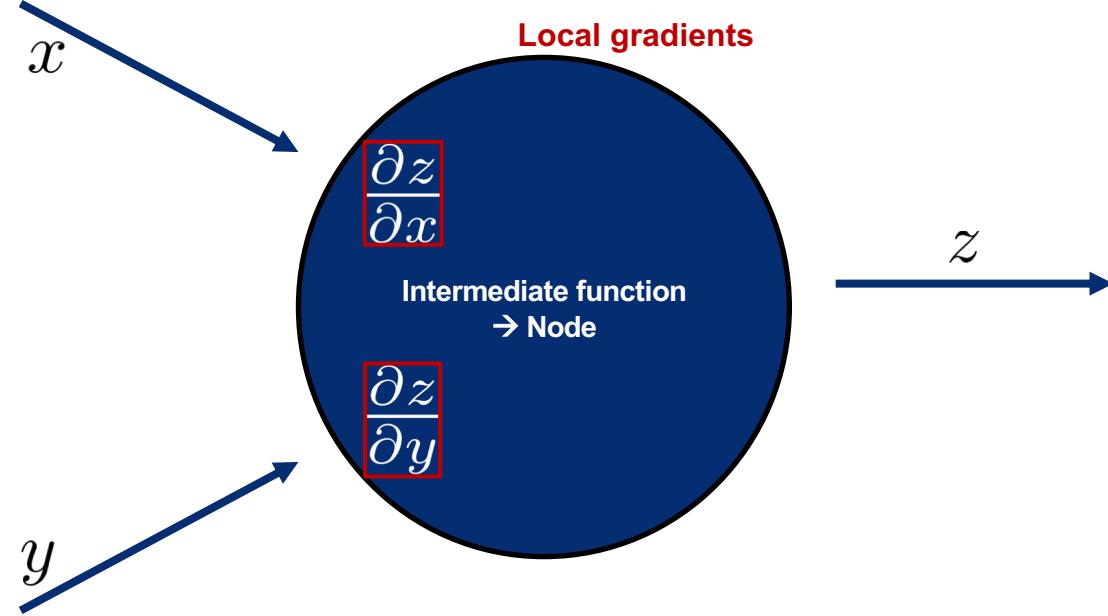
Score function: $s = f(x_i, W) = Wx_i$

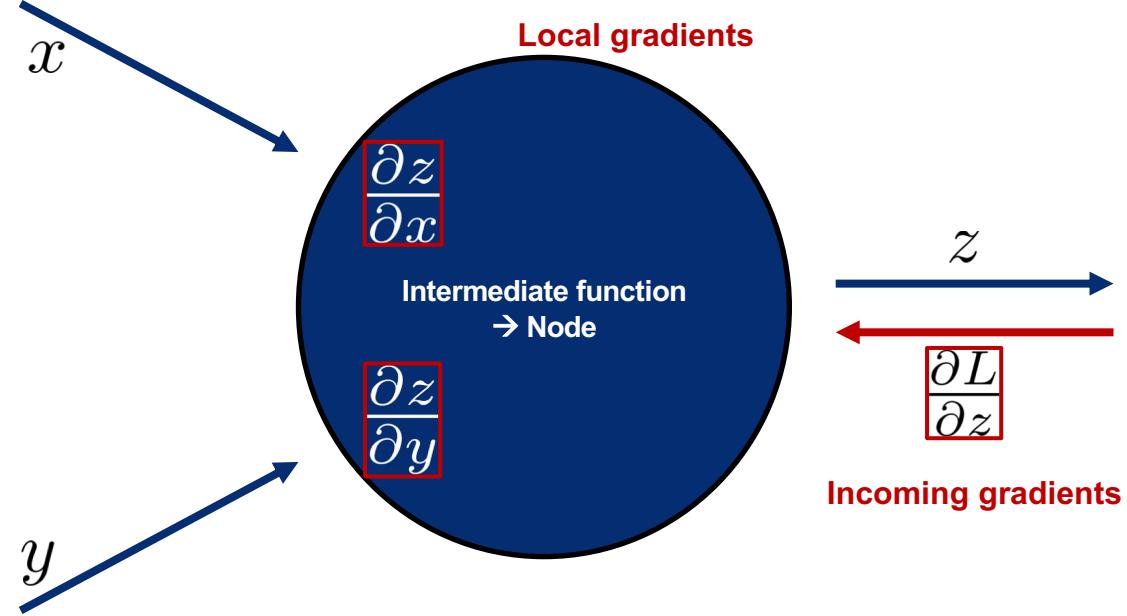
Loss function: $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$

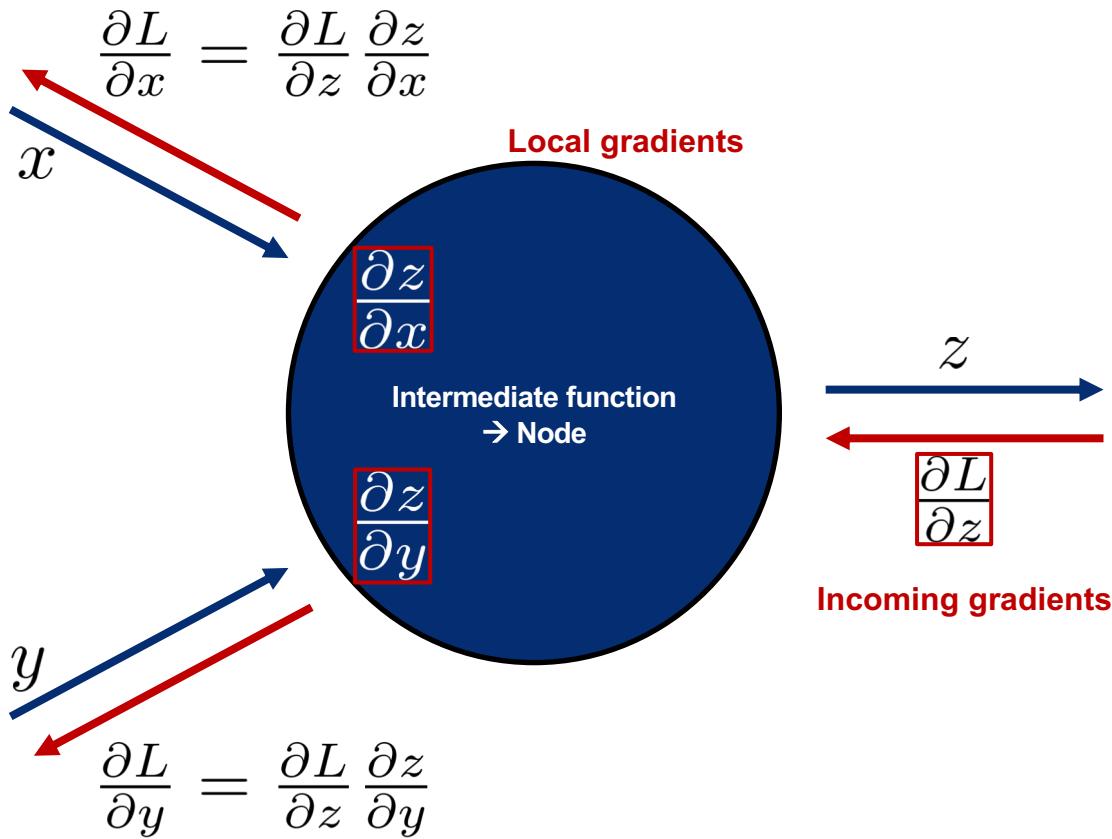
$$L(W, x) = \frac{1}{N} \sum_i L_i (f(x_i, W), y_i) + \lambda R(W)$$

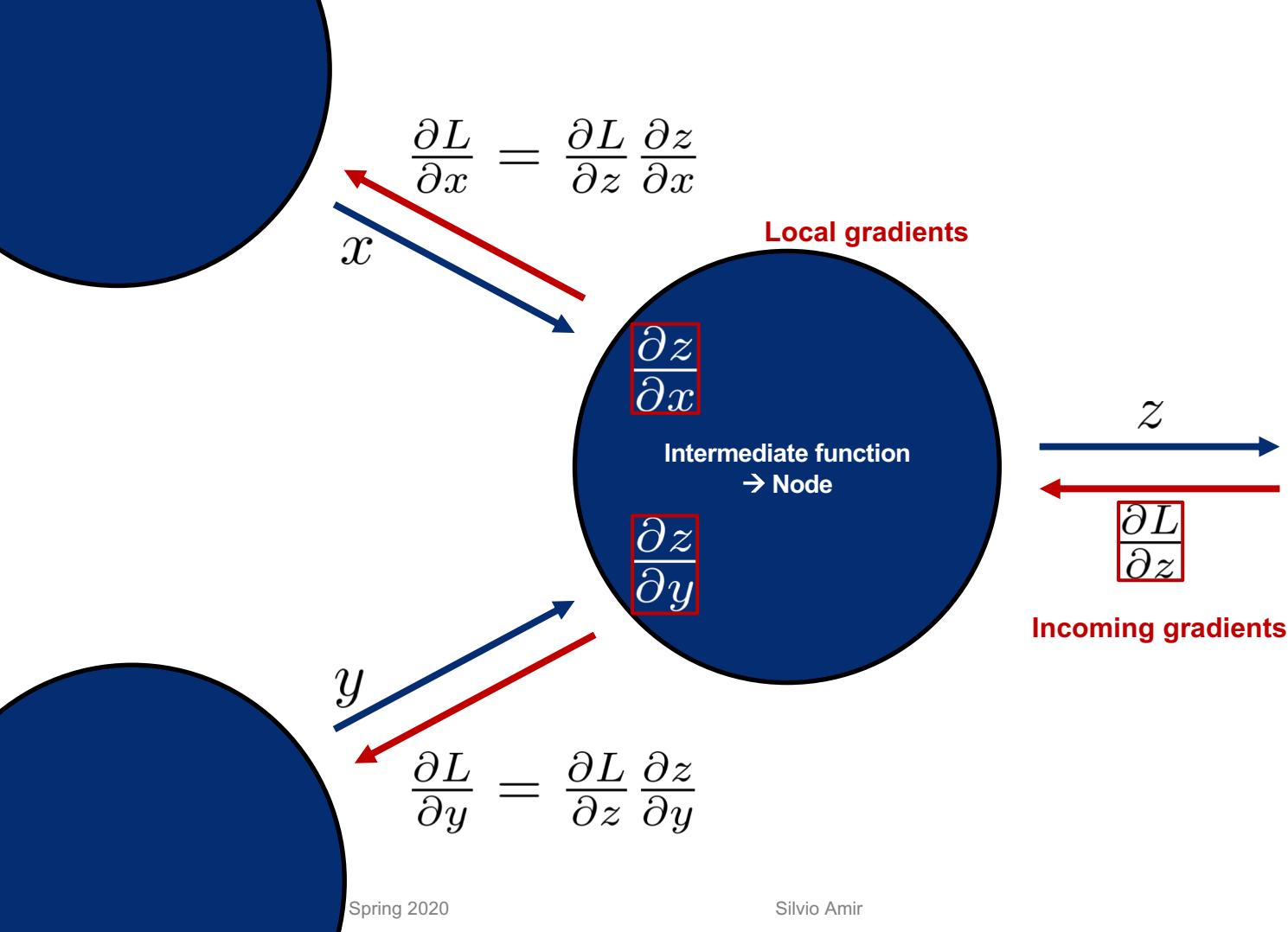












Patterns in Flow

Add node

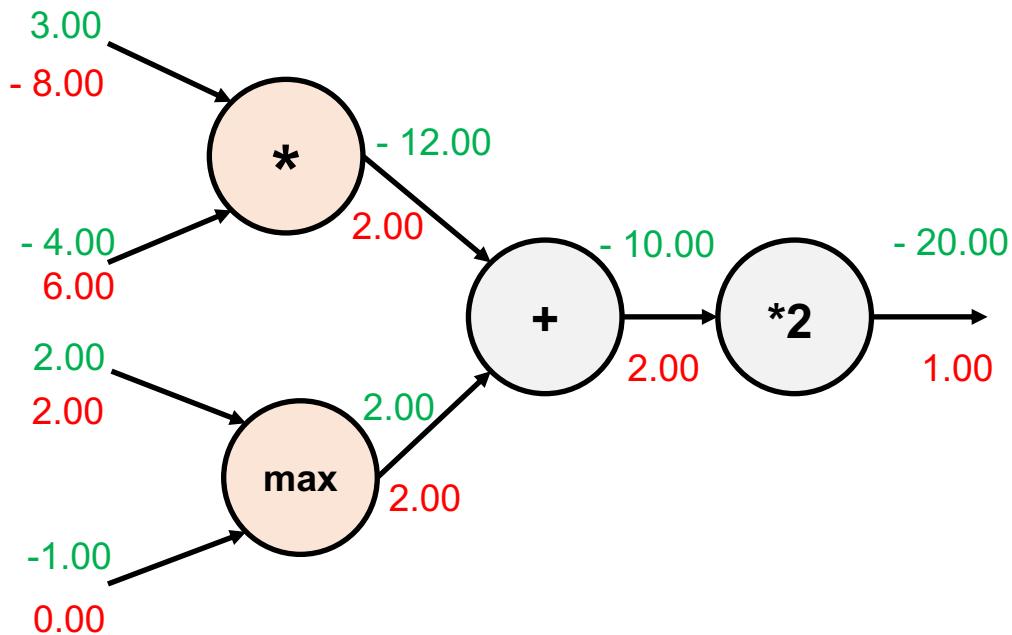
→ Distributes gradient

Max node

→ Routes gradient

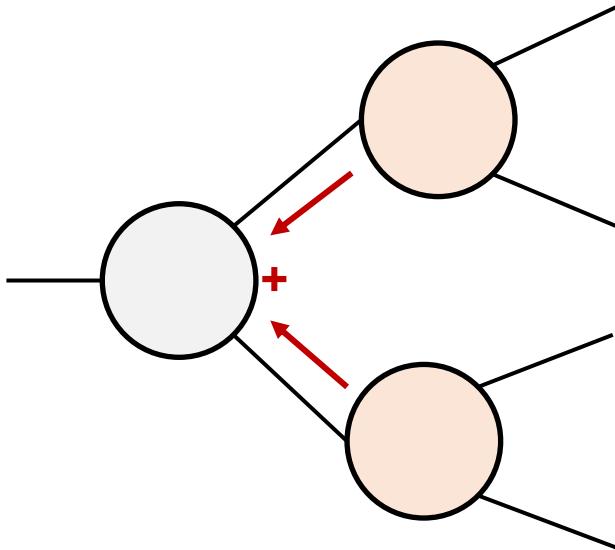
Mul node

→ Switches gradient



Patterns in Flow

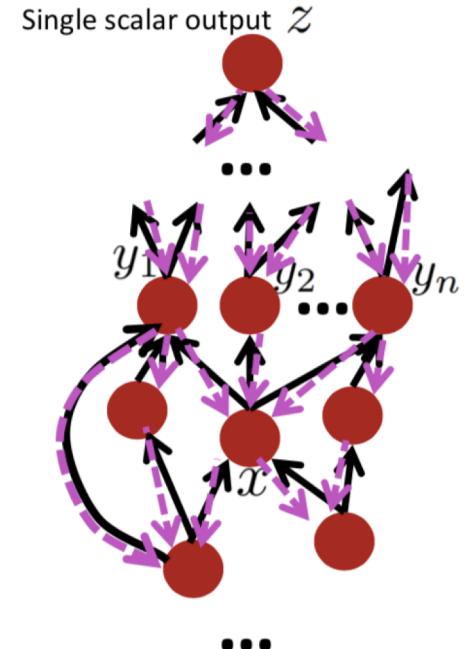
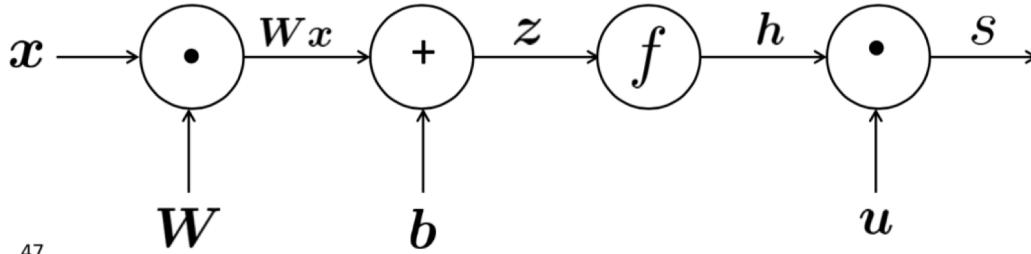
Gradients add at branches



Backpropagation Algorithm

Compute derivatives of the loss wrt to all parameters of the network

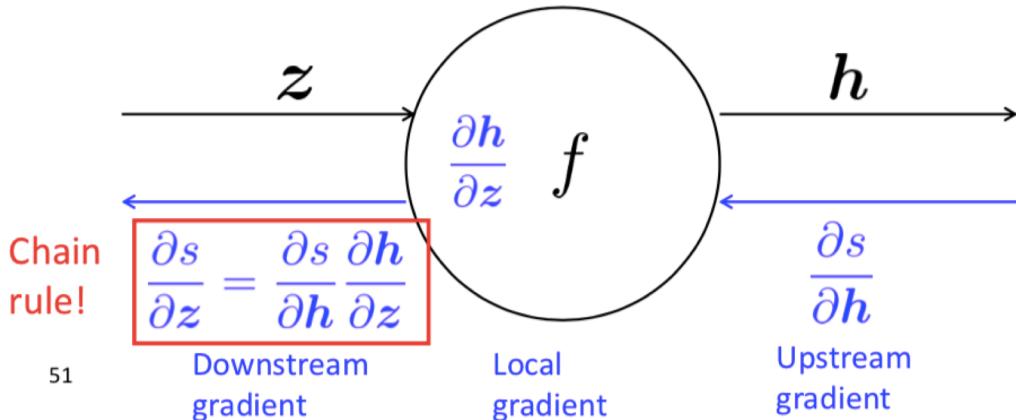
- **breakdown the networks computation into small steps**
- define local derivatives of outputs wrt to inputs
- recursively apply the chain rule to compose derivatives
- re-use intermediate values computed in previous steps



Backpropagation Algorithm

Compute derivatives of the loss wrt to all parameters of the network

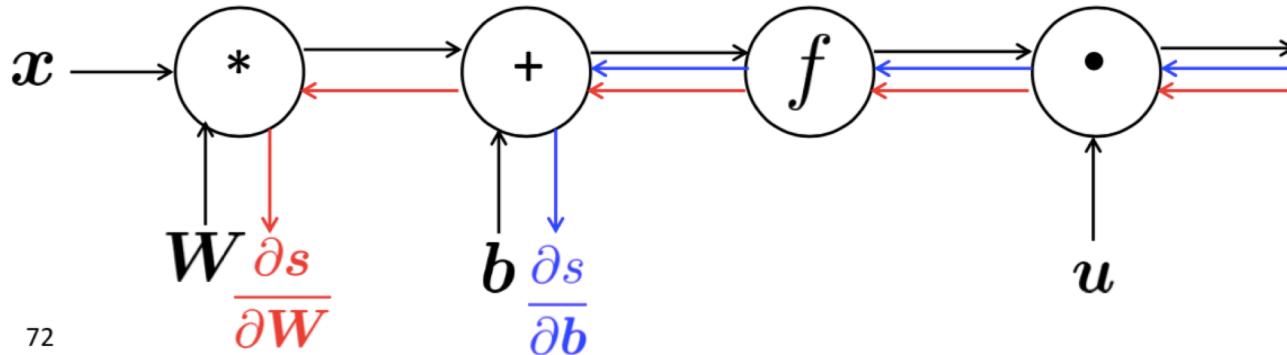
- breakdown the networks computation into small steps
- **define local derivatives of outputs wrt to inputs**
- recursively apply the chain rule to compose derivatives
- re-use intermediate values computed in previous steps



Backpropagation Algorithm

Compute derivatives of the loss wrt to all parameters of the network

- breakdown the networks computation into small steps
- define local derivatives of outputs wrt to inputs
- **recursively apply the chain rule to compose derivatives**
- re-use intermediate values computed in previous steps

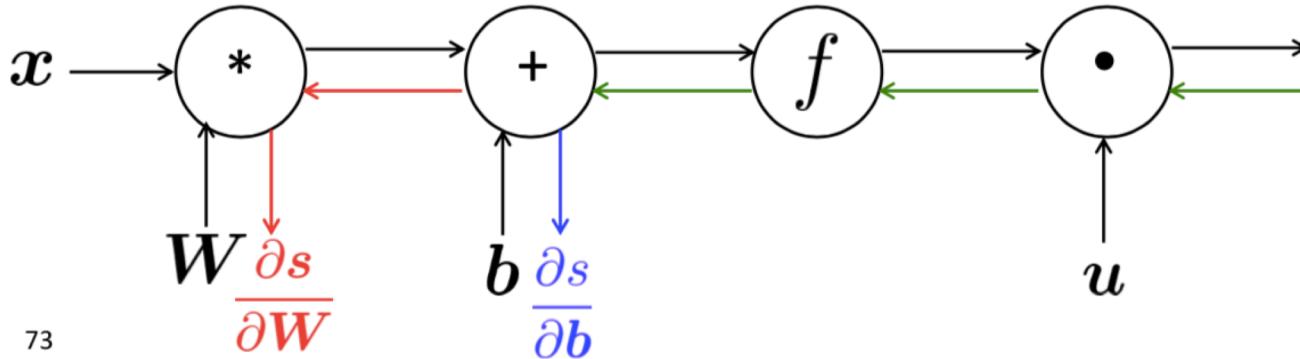


72

Backpropagation Algorithm

Compute derivatives of the loss wrt to all parameters of the network

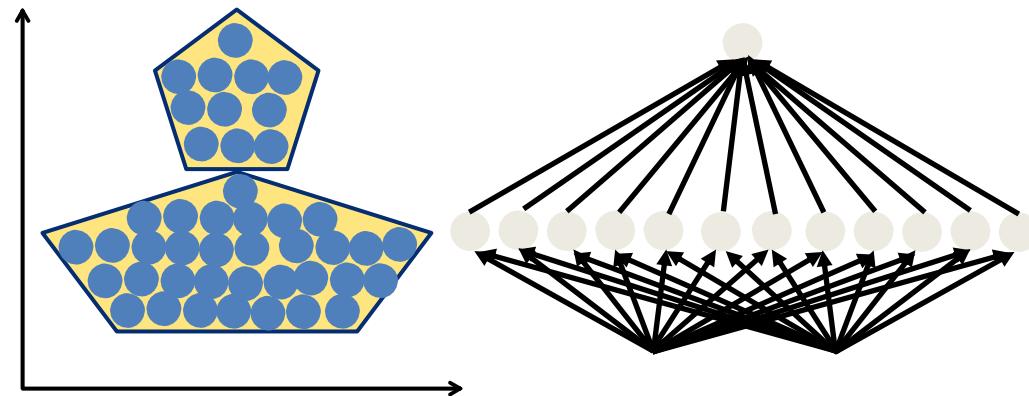
- breakdown the networks computation into small steps
- define local derivatives of outputs wrt to inputs
- recursively apply the chain rule to compose derivatives
- **re-use intermediate values computed in previous steps**



73

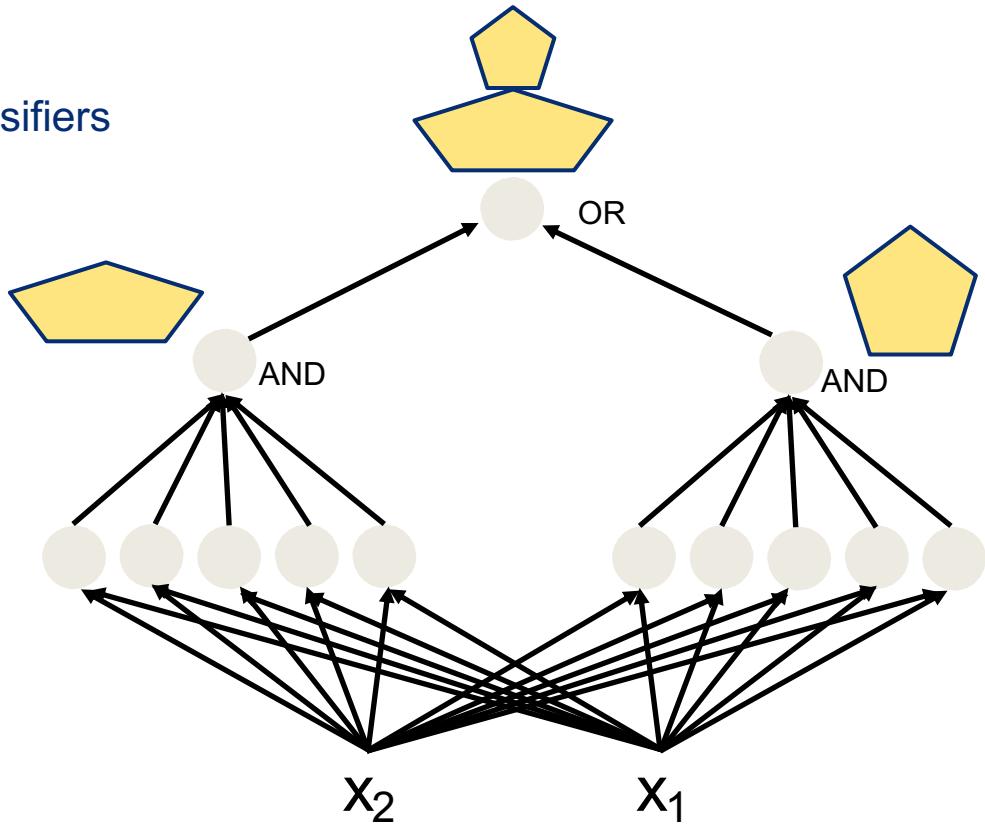
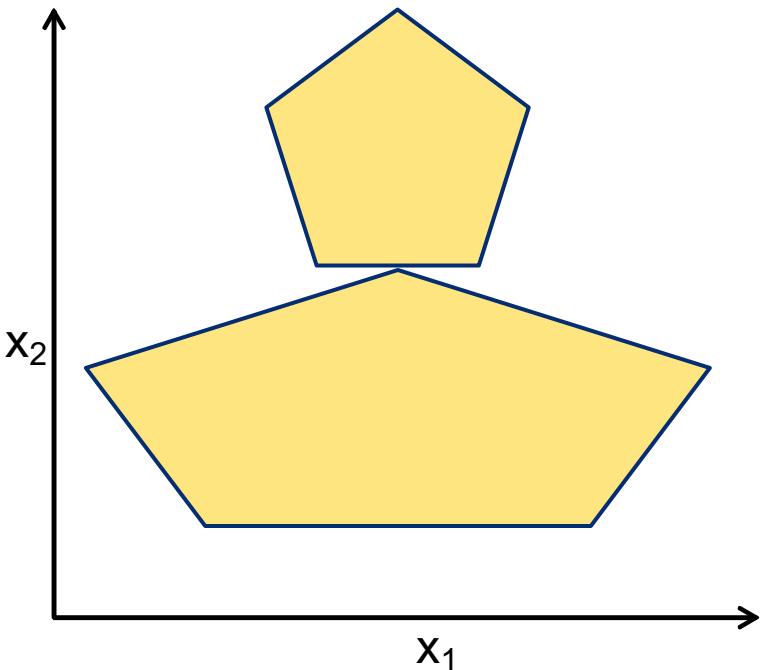
Multilayer Perceptron

Multilayer Perceptrons are Universal Function Approximators



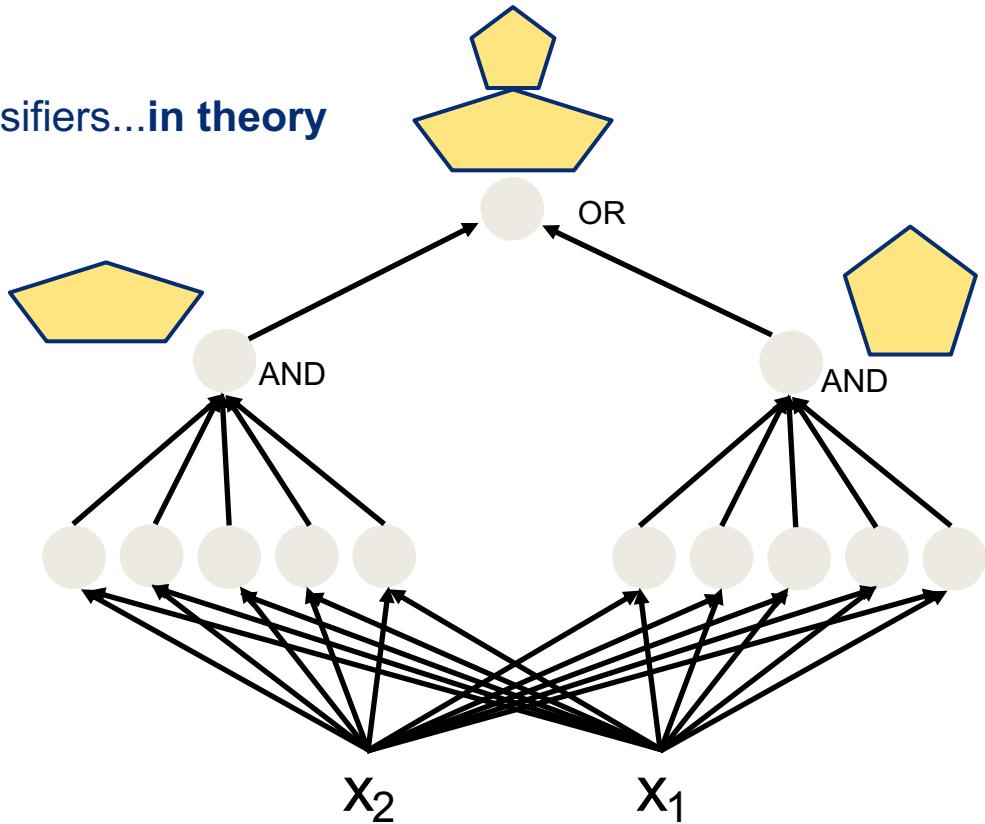
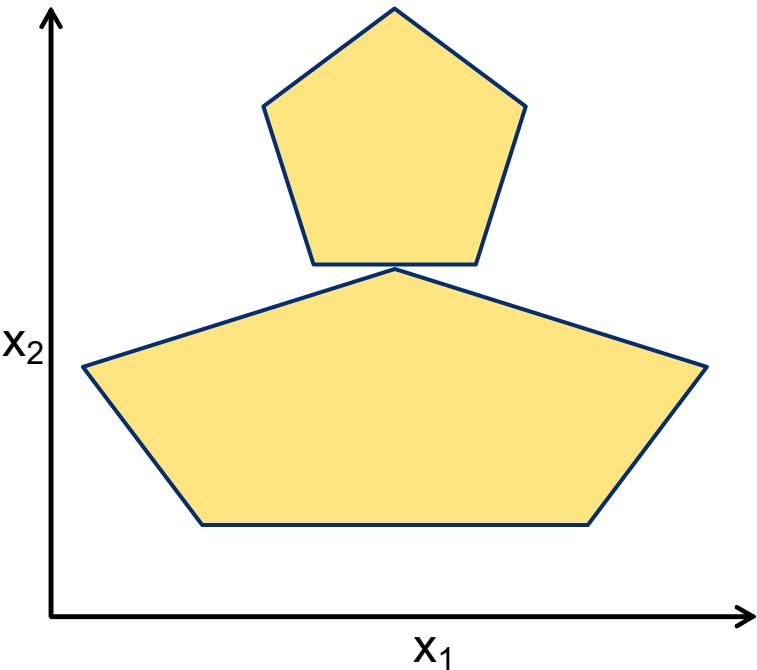
Multilayer Perceptron

Multilayer Perceptrons are Universal Classifiers



Multilayer Perceptron

Multilayer Perceptrons are Universal Classifiers...in theory



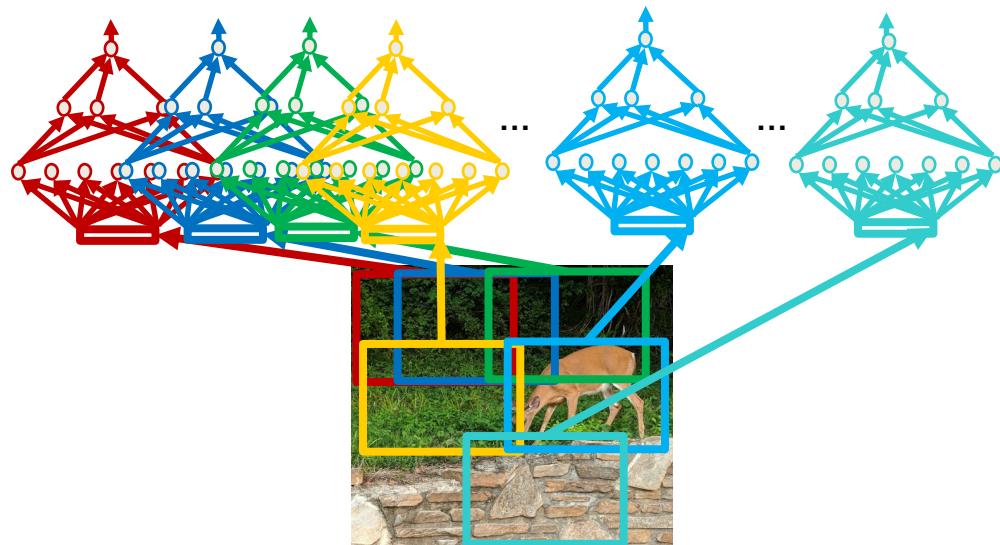
Pattern Recognition

Conventional MLPs are sensitive to location, but often
the **location** of a pattern **is not important!**



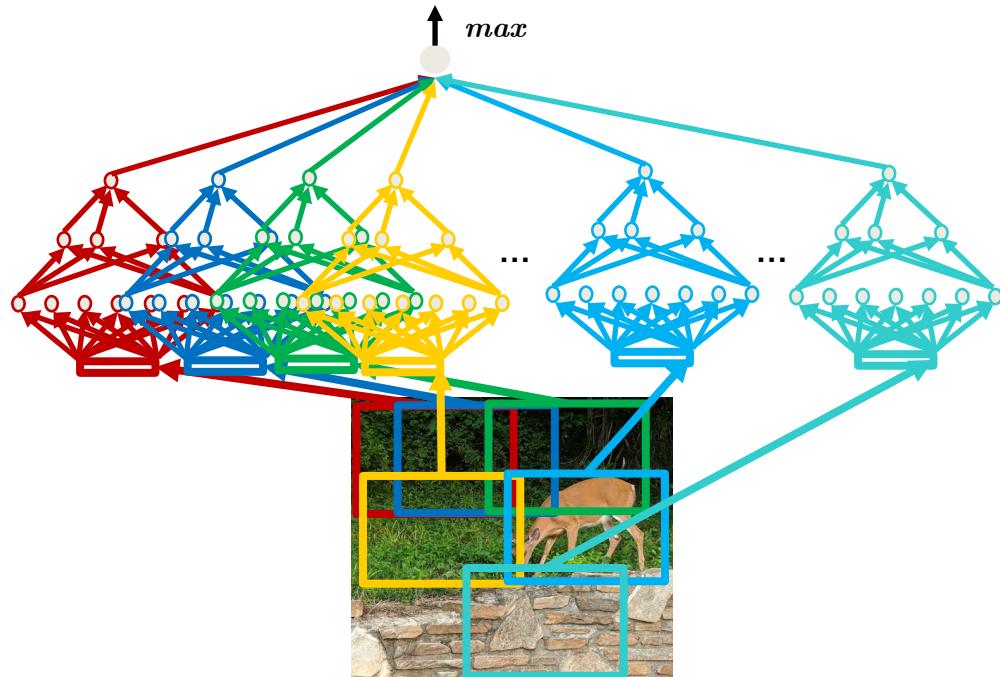
Scanning with a MLP

Shift invariance can be achieved by scanning with a MLP



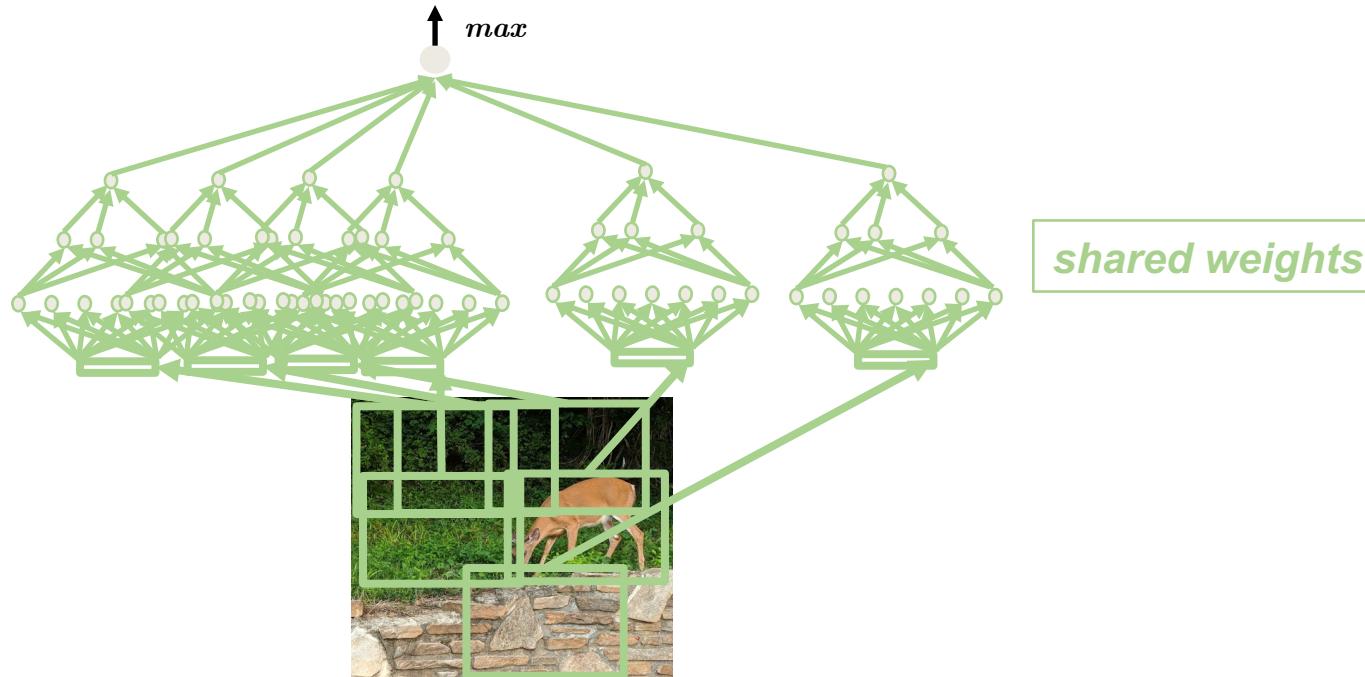
Scanning with a MLP

Shift invariance can be achieved by scanning with a MLP
Combine multiple local predictions



Scanning with a MLP

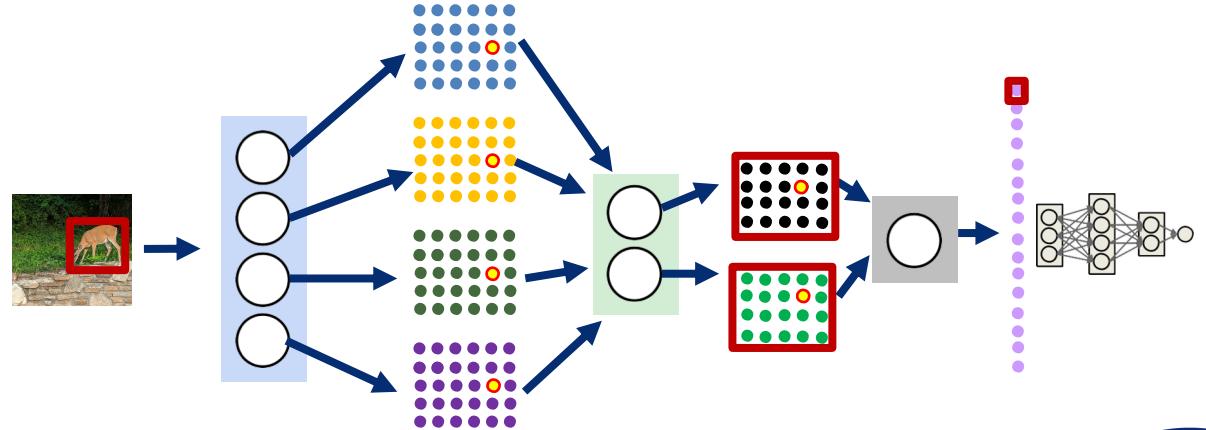
Shift invariance can be achieved by scanning with a MLP
Combine multiple local predictions



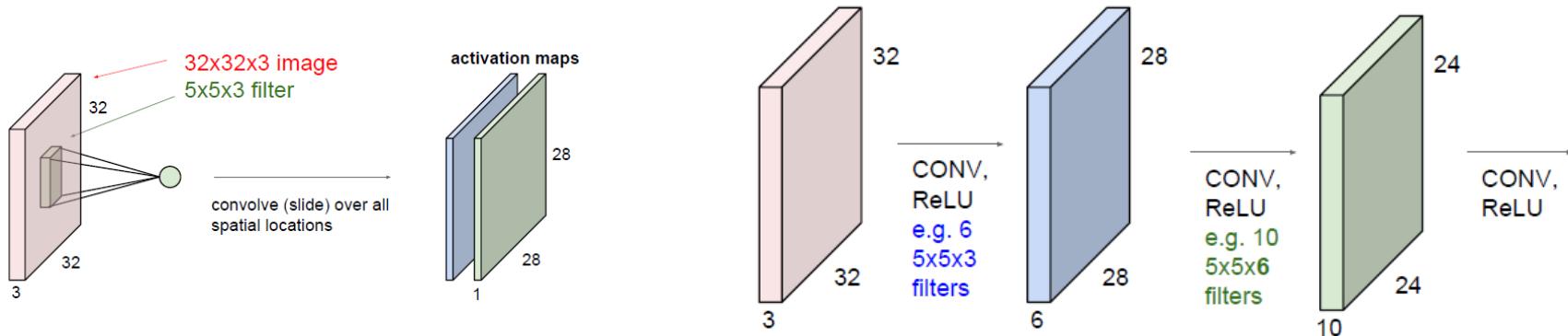
Scanning with a MLP

Better approach: scan the input one layer at a time

- Hierarchical features
- Compositionality



Convolutional Layers



Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

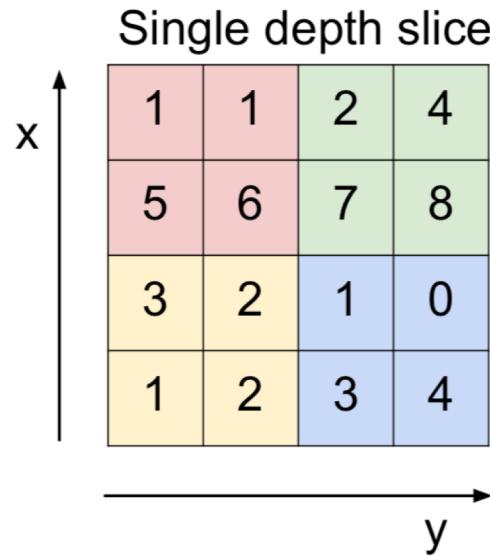
Output volume size:
 $(32+2*2-5)/1+1 = 32$ spatially, so
32x32x10

Number of parameters in this layer?
each filter has $5*5*3 + 1 = 76$ params (+1 for bias)
 $=> 76*10 = 760$



Pooling Layers

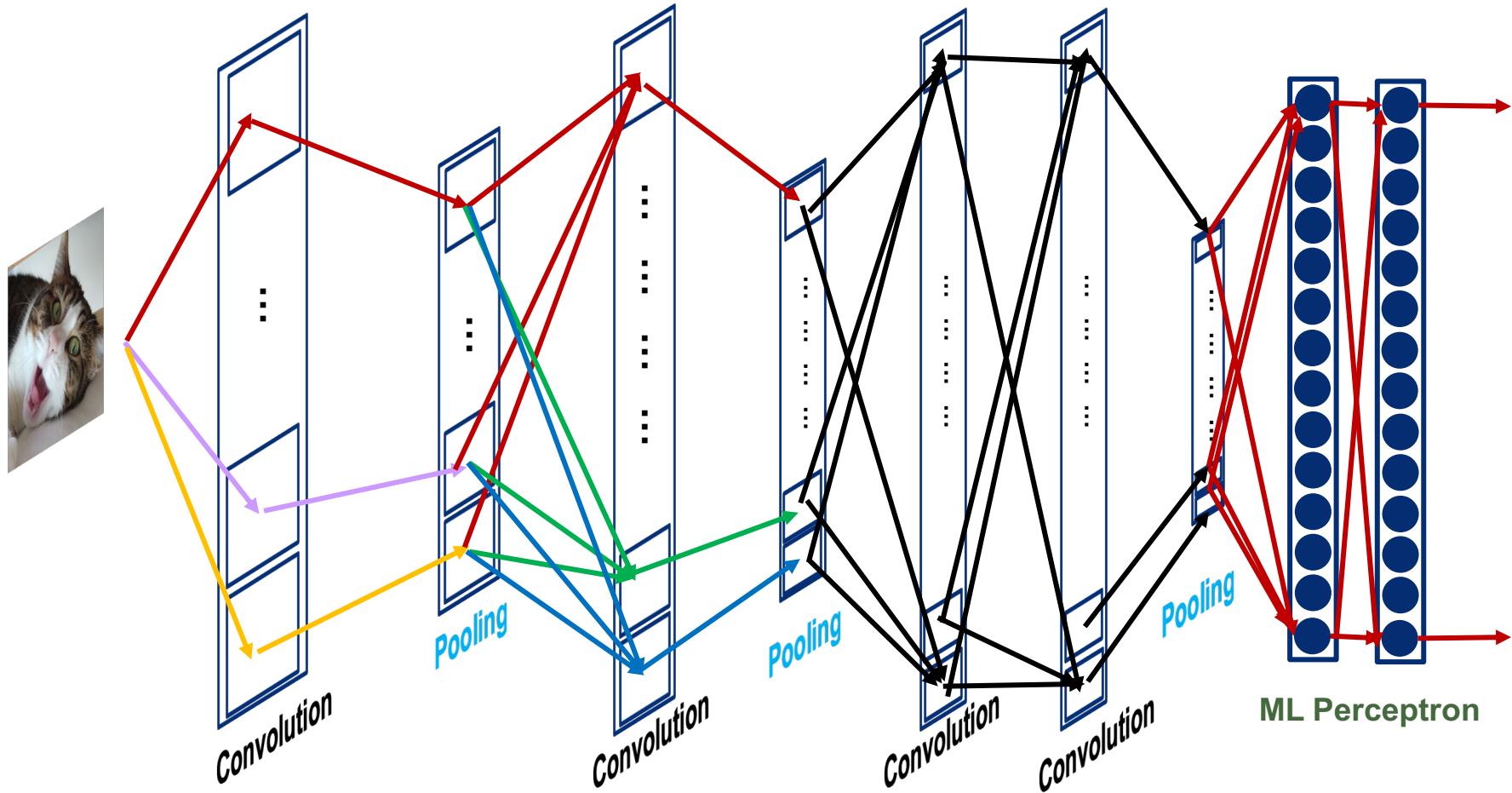
Max Pooling



max pool with 2x2 filters
and stride 2

6	8
3	4

Convolutional Neural Networks



Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid

Optimization

Gradient Descent

Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
Multilayer Perceptron
Fully Connected (Dense) Layers
Shared Weight Layers,
Convolutional Layers
Pooling Layers

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid

Optimization

Gradient Descent
Backpropagation



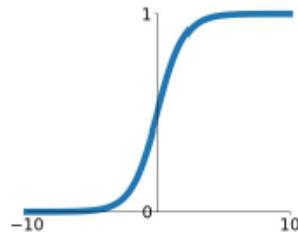
Recap

Deep Neural Networks

The Zoo of Common Activation Functions

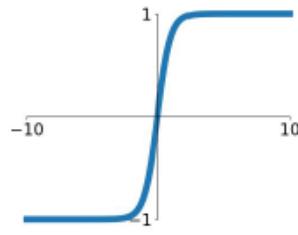
Sigmoid

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



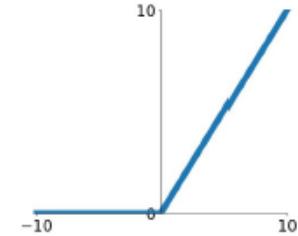
tanh

$$\tanh(x)$$



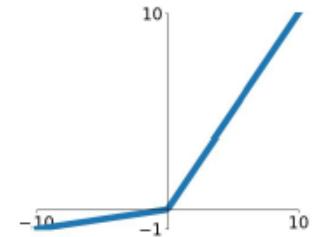
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

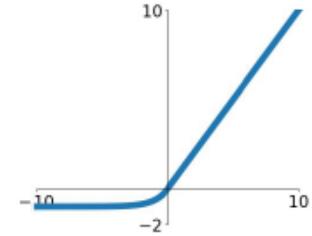


Maxout

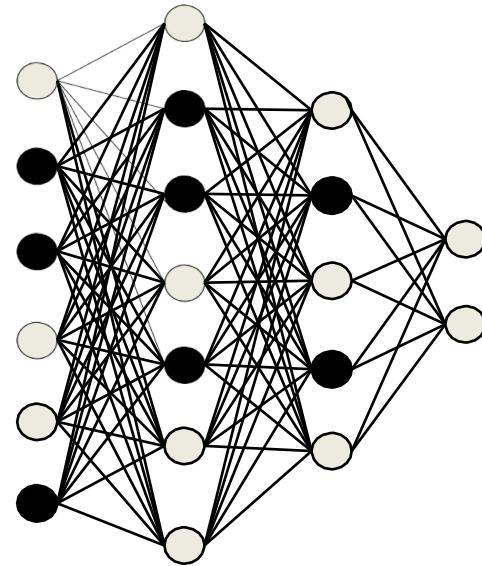
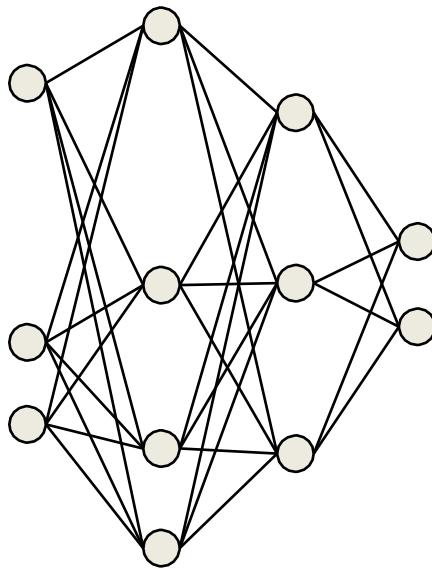
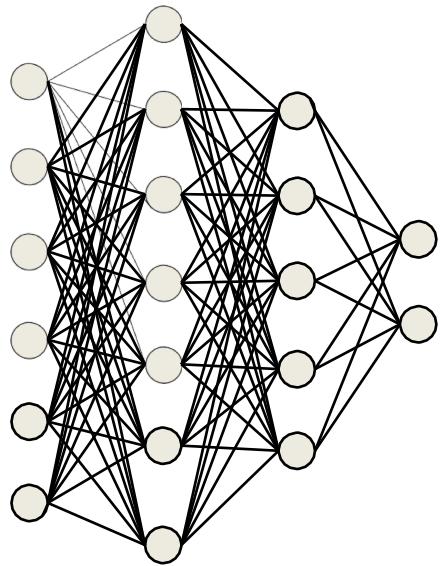
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Dropout



During training

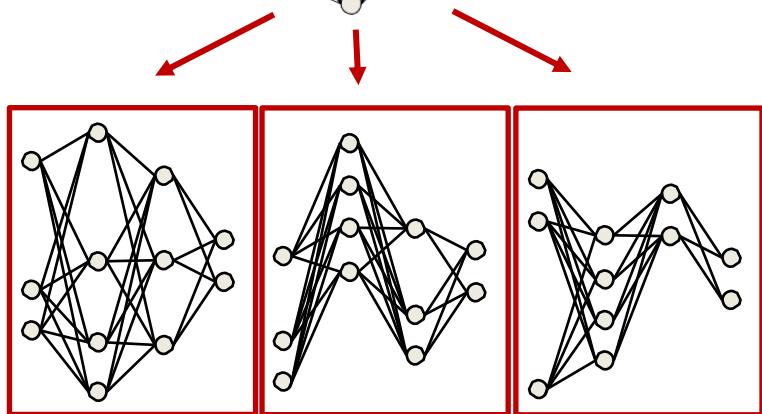
- At each iteration, in each layer, “knock out” each neuron with probability $1-\alpha$
- In practice, we do not drop connections but set inputs/outputs to zero

Dropout

Inference

- Scale weights or activations

→ Learns a N -neuron network that averages over a sample of 2^N possible sub-networks



Batch Normalization

1. Compute empirical mean and variance for each channel

$$E[x^{(k)}], \text{Var}[x^{(k)}]$$

2. Normalize to unit Gaussian

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

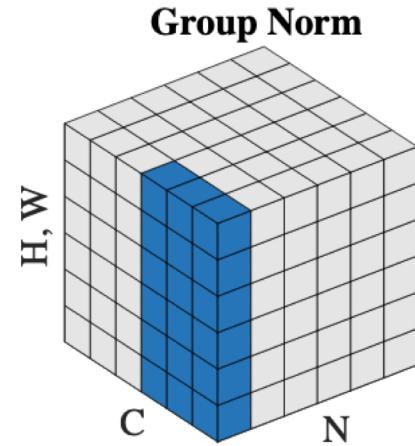
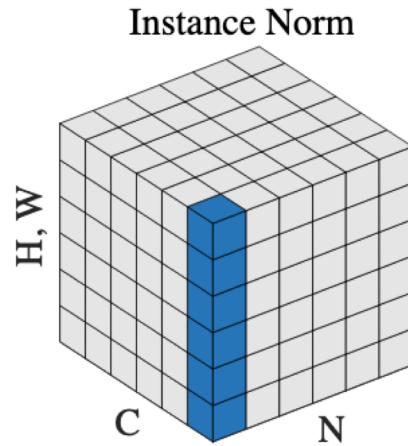
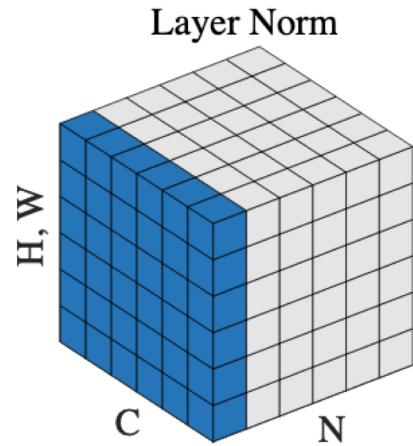
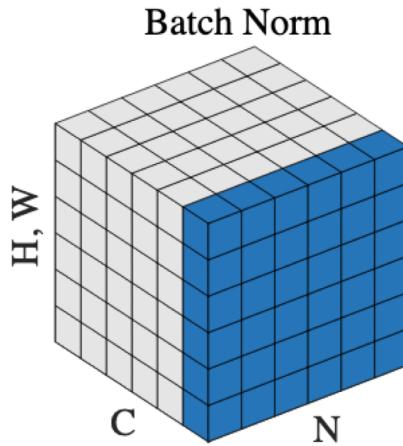
3. Squash output to beneficial range

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

These are parameters and are learned during training.

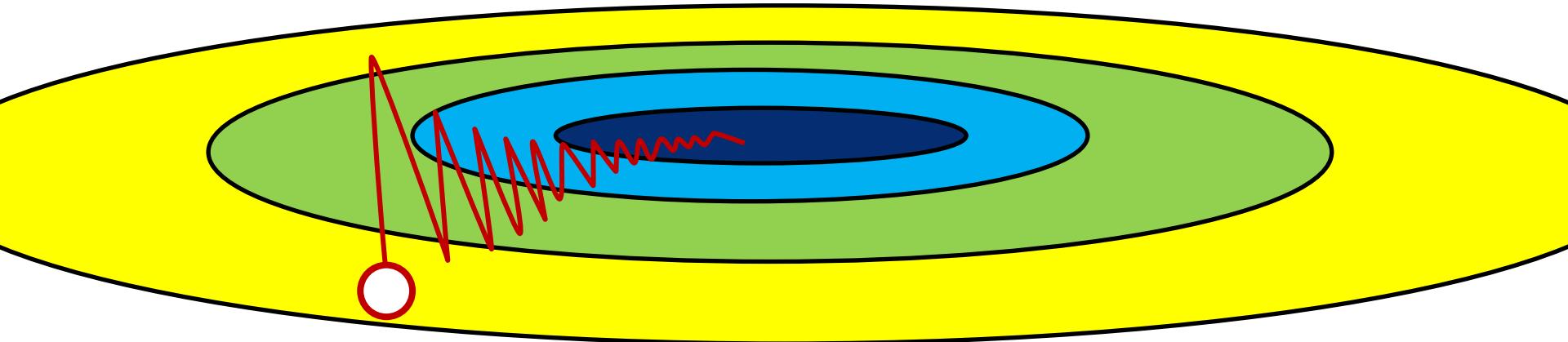


Normalization Methods



Problems with SGD Optimization

- Slow progress along shallow dimension, jitter along the other
- Problematic: Neural networks have millions of parameters!



Why is this?

Condition number: Ratio of largest to smallest singular value of the Hessian
→ If large, then loss function at this point badly conditioned

Another Problem

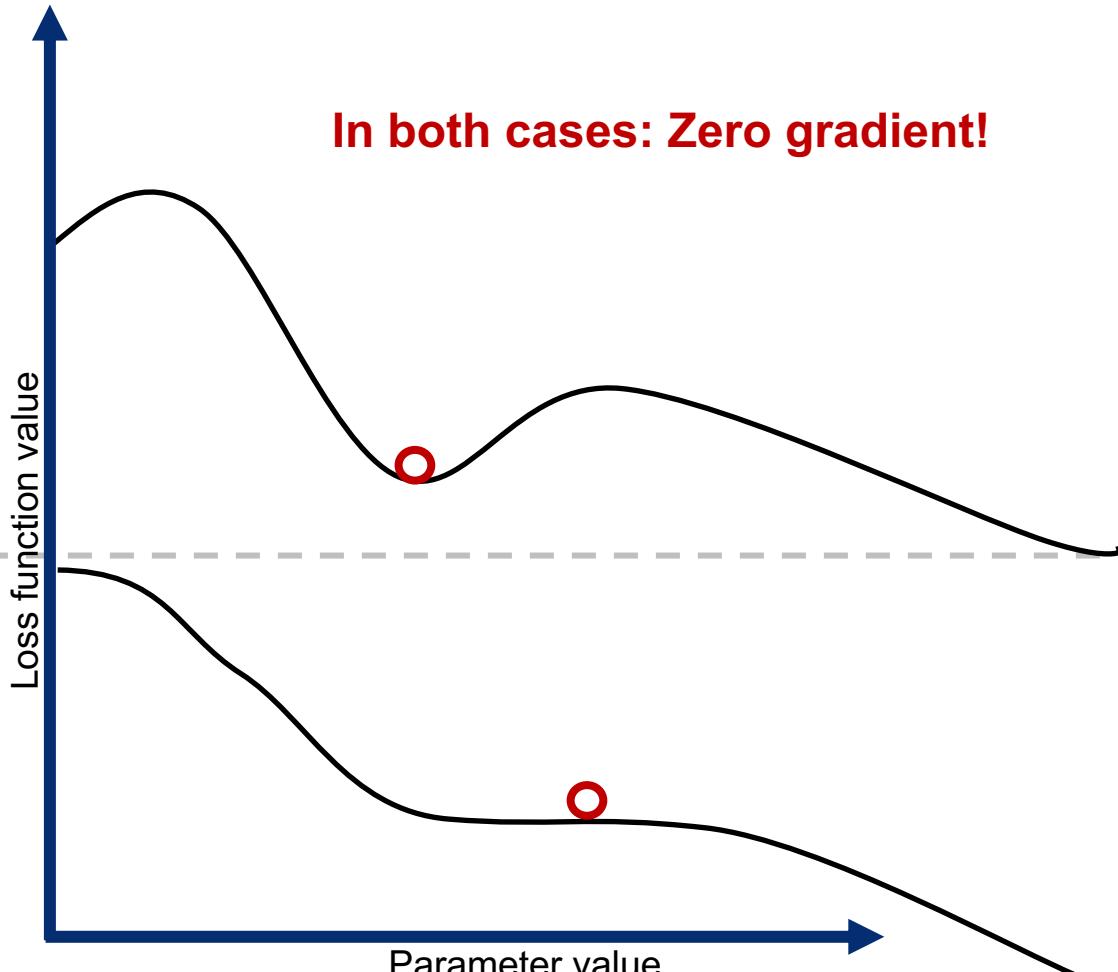
Local minimum

In every direction, loss will go up.

Saddle point

In some direction loss will go up,
in other direction loss will go down.

In high dimensional space, this
scenario is much more common.



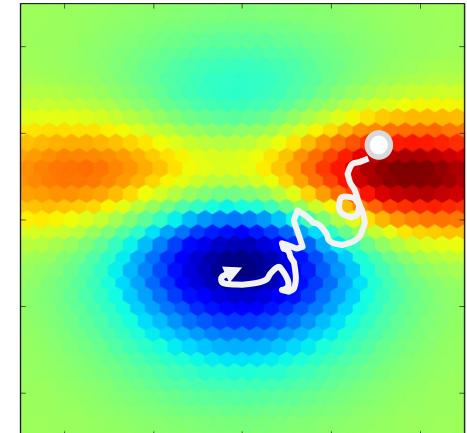
→ Saddle points are the big problem when training neural networks!

And Another Problem

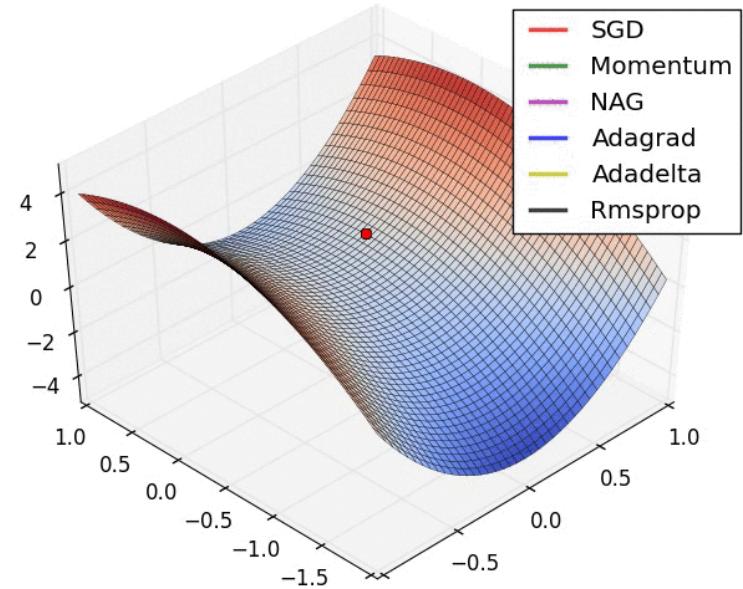
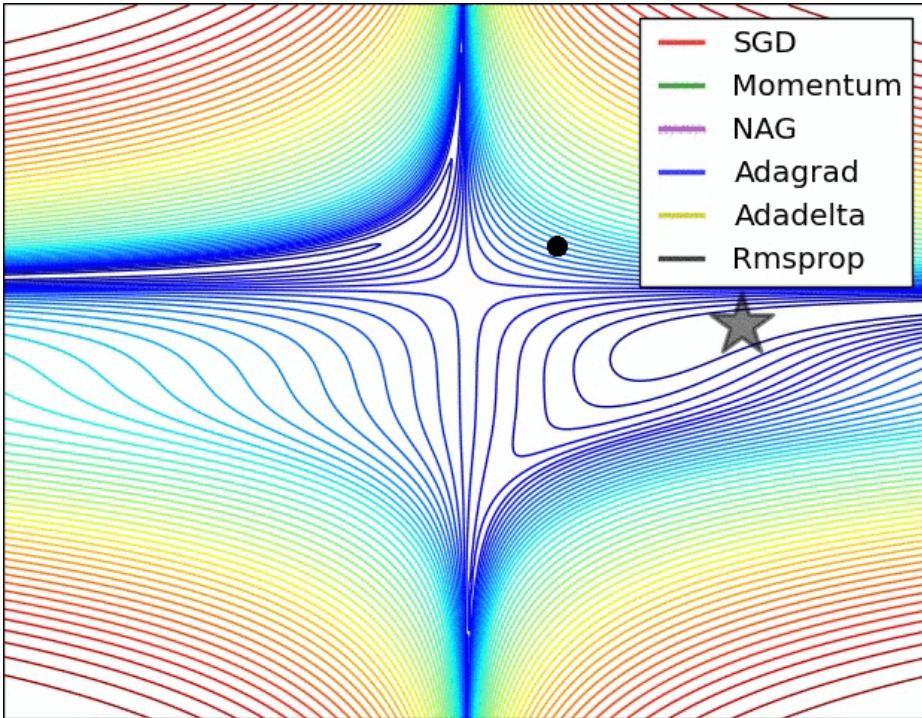
```
while not_converged:  
    data_batch = sample_training_data(data, batch_size)  
    gradient = eval_gradient(loss, data_batch, weights)  
    weights += - step_size * gradient
```

Gradient is computed over mini-batches

- Mini-batches do not necessarily represent the full dataset
- **Gradients can be noisy!**



Optimizers



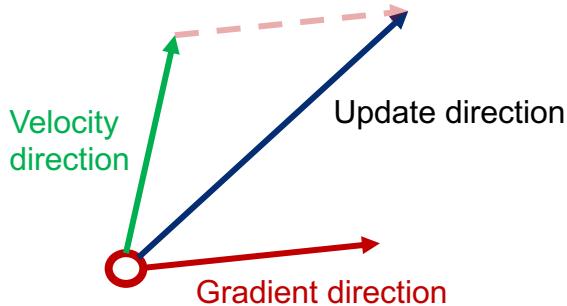
Adding Momentum

Nesterov, Y. E. (1983). A method for solving the convex programming problem with convergence rate O (1/k^2). In *Dokl. Akad. Nauk SSSR* (Vol. 269, pp. 543-547).
[Sutskever, I., Martens, J., Dahl, G., & Hinton, G. \(2013, February\). On the importance of initialization and momentum in deep learning. In International conference on machine learning \(pp. 1139-1147\).](#)

SGD + Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla_W L(W_t)$$
$$W_{t+1} = W_t + v_{t+1}$$

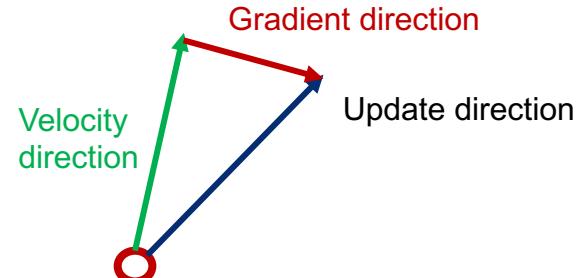
Combine gradient at current point with velocity to get update



Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla_W L(W_t + \rho v_t)$$
$$W_{t+1} = W_t + v_{t+1}$$

Evaluate gradient at where velocity would take us, then mix with velocity



$$g_t = \nabla_W L(W_t)$$

$$S_i = S_i + (g_t)_i^2 \quad \text{with } S_i(t=0) = 0$$

$$(\mathrm{d}W_t)_i = \frac{\alpha}{\sqrt{S_i} + \epsilon} (g_t)_i$$

$$W_{t+1} = W_t - \mathrm{d}W_t$$

1. Compute gradient
2. Compute and accumulate element-wise squared gradient
3. Compute gradient update with **parameter-wise** learning rate
4. Apply gradient update

Adam

[Kingma, D. P., & Ba, J. \(2014\). Adam: A method for stochastic optimization. arXiv:1412.6980.](#)

$$g_t = \nabla_W L(W_t)$$

Bias correction

$$S_i^{(1)} = (\rho_1 S_i^{(1)} + (1 - \rho_1)(g_t)_i) (1 - \rho_1^t)^{(-1)}$$

$$S_i^{(2)} = (\rho_2 S_i^{(2)} + (1 - \rho_2)(g_t)_i^2) (1 - \rho_2^t)^{(-1)}$$

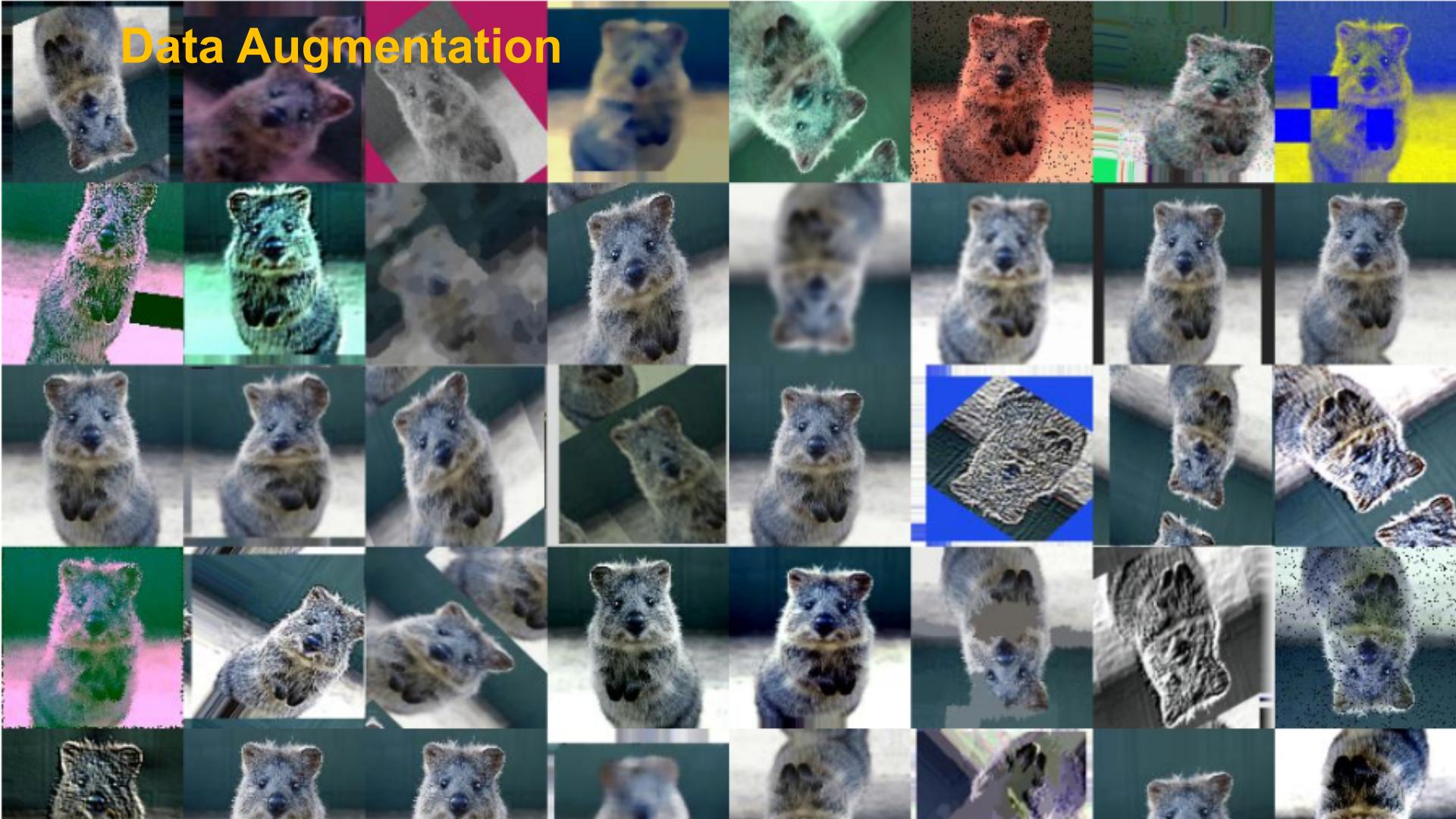
$$(\text{d}W_t)_i = \frac{\alpha}{\sqrt{S_i^{(2)}} + \epsilon} S_i^{(1)}$$

$$\begin{aligned} S_i^{(1)}(t=0) &= 0 \\ S_i^{(2)}(t=0) &= 0 \end{aligned}$$

$$W_{t+1} = W_t - \text{d}W_t$$

1. Compute gradient
2. Compute first momentum (“velocity”)
3. Compute second momentum (parameter-wise normalization)
4. Compute update with **momentum** and **parameter-wise** learning rate
5. Apply update

Data Augmentation



3x3 conv, 64

3x3 conv, 64

pool/2

3x3 conv, 128

3x3 conv, 128

pool/2

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256

pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

pool/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

pool/2

fc 4096

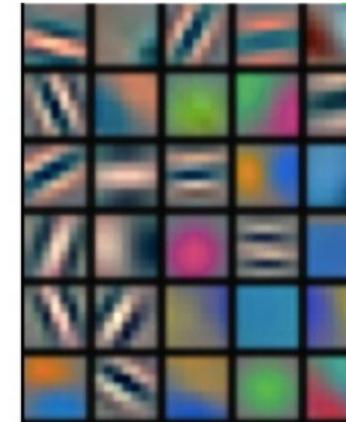
fc 4096

fc 4096

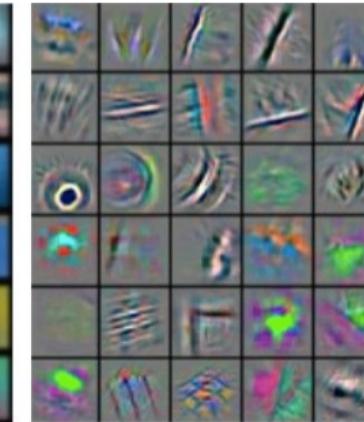
Transfer Learning

Fairly generic

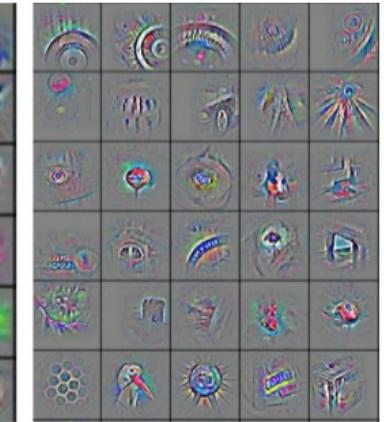
low-level features



mid-level features



high-level features



[Sharif Razavian, A., Azizpour, H., Sullivan, J., & Carlsson, S. \(2014\). CNN features off-the-shelf: an astounding baseline for recognition. CVPR \(pp. 806-813\).](#)

Rather specific



Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
Multilayer Perceptron
Fully Connected (Dense) Layers
Shared Weight Layers,
Convolutional Layers
Pooling Layers

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid

Optimization

Gradient Descent
Backpropagation



Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
Multilayer Perceptron
Fully Connected (Dense) Layers
Shared Weight Layers,
Convolutional Layers
Pooling Layers

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid
Tanh
ReLU
eLu
PreLu
MaxOut

Optimization

Gradient Descent
Backpropagation
Momentum methods
Adaptive gradient methods

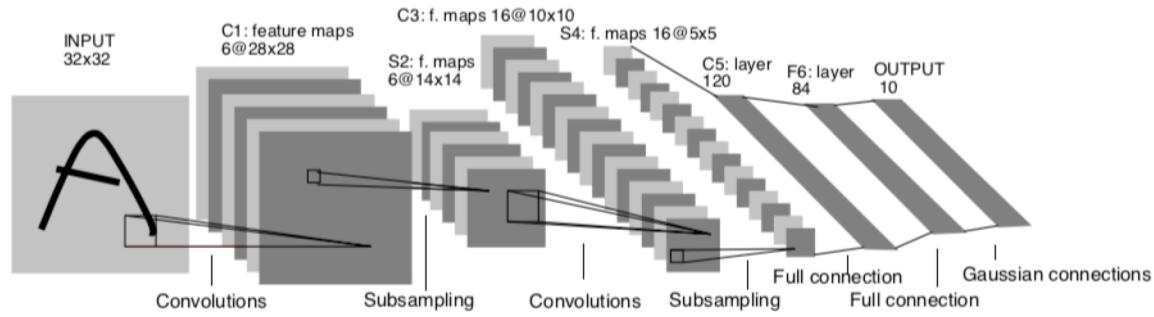
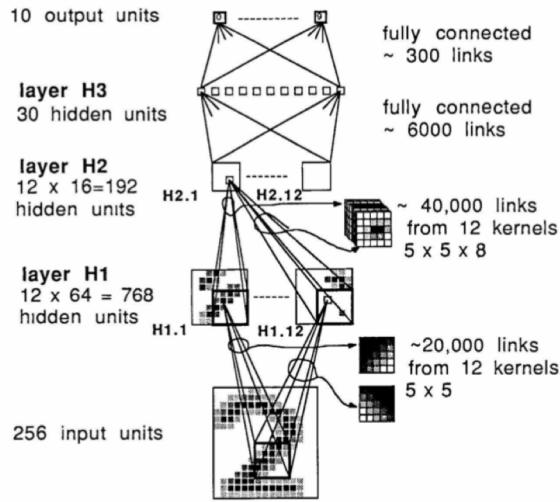
Tricks of the trade

Dropout
Batch Normalization
Initialization
Data Augmentation
Transfer Learning

Recap

Computer Vision Architectures

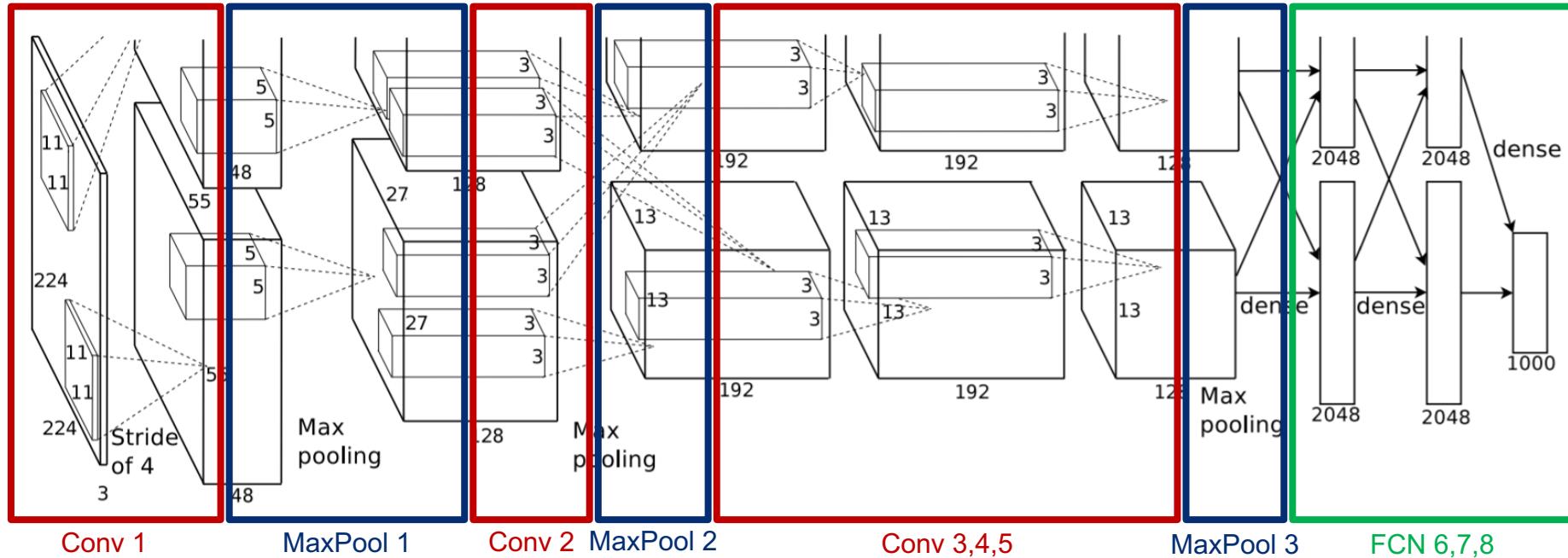
LeNet



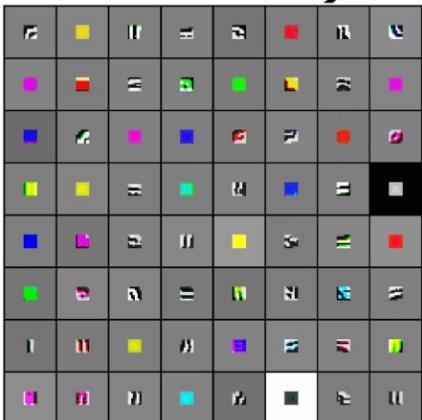
LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., & Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4), 541-551.

LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.

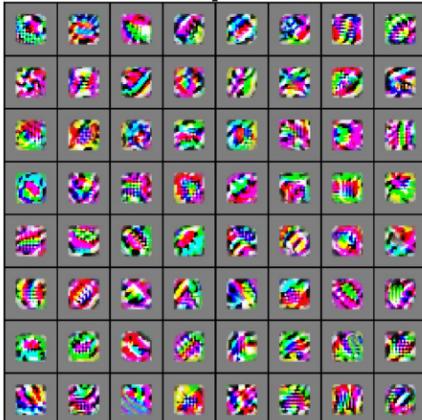
AlexNet



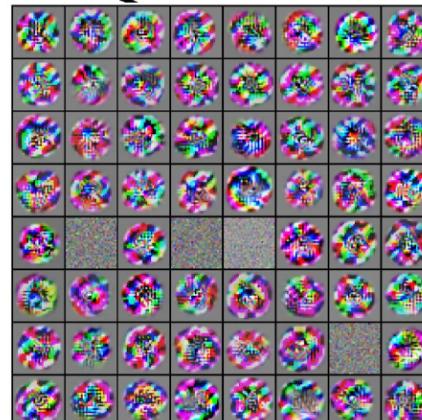
Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. NeurIPS (pp. 1097-1105).



VGG-16 Conv1_1



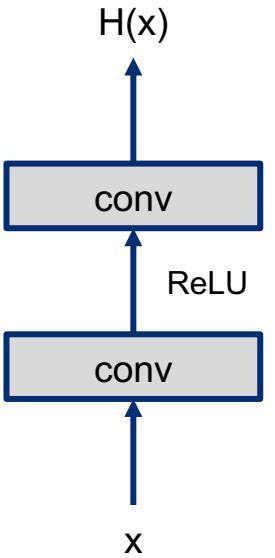
VGG-16 Conv3_2



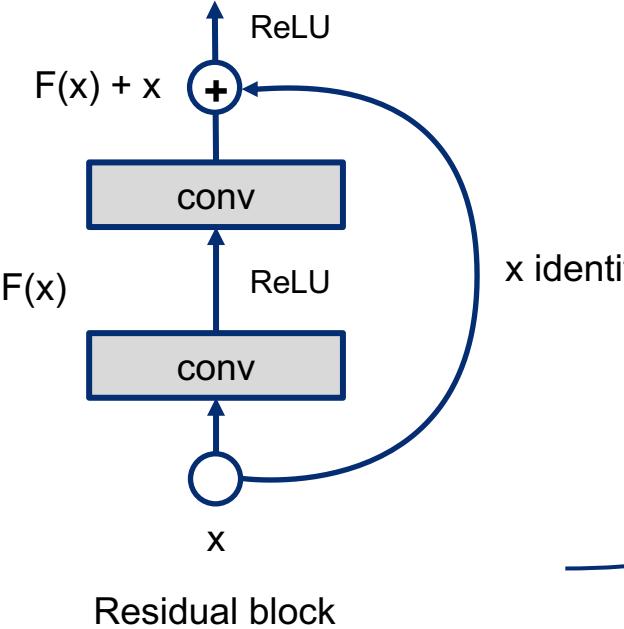
VGG-16 Conv5_3



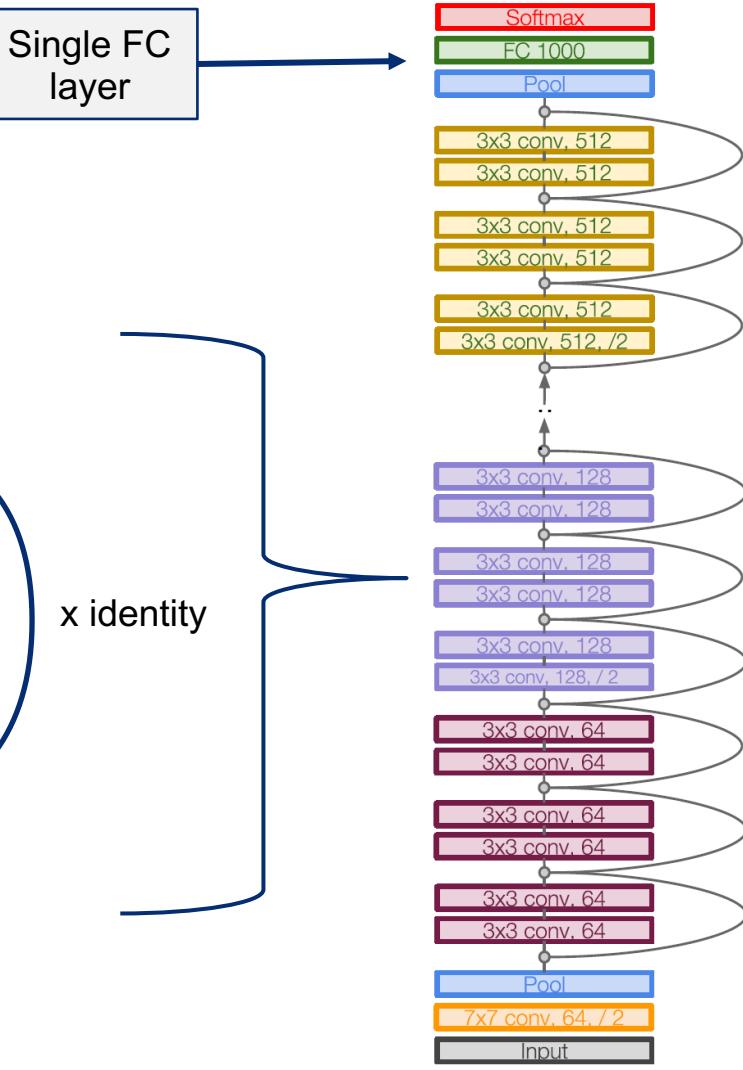
ResNet



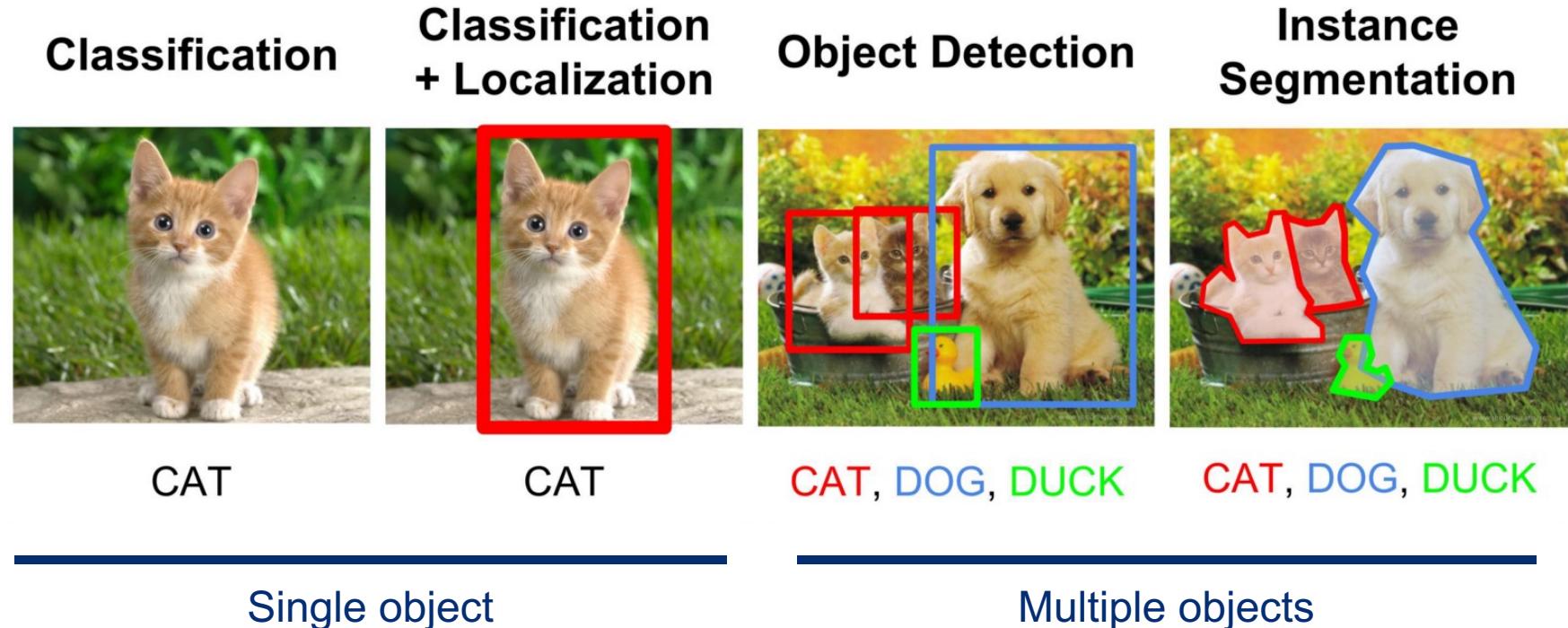
"Plain" conv layers



Residual block



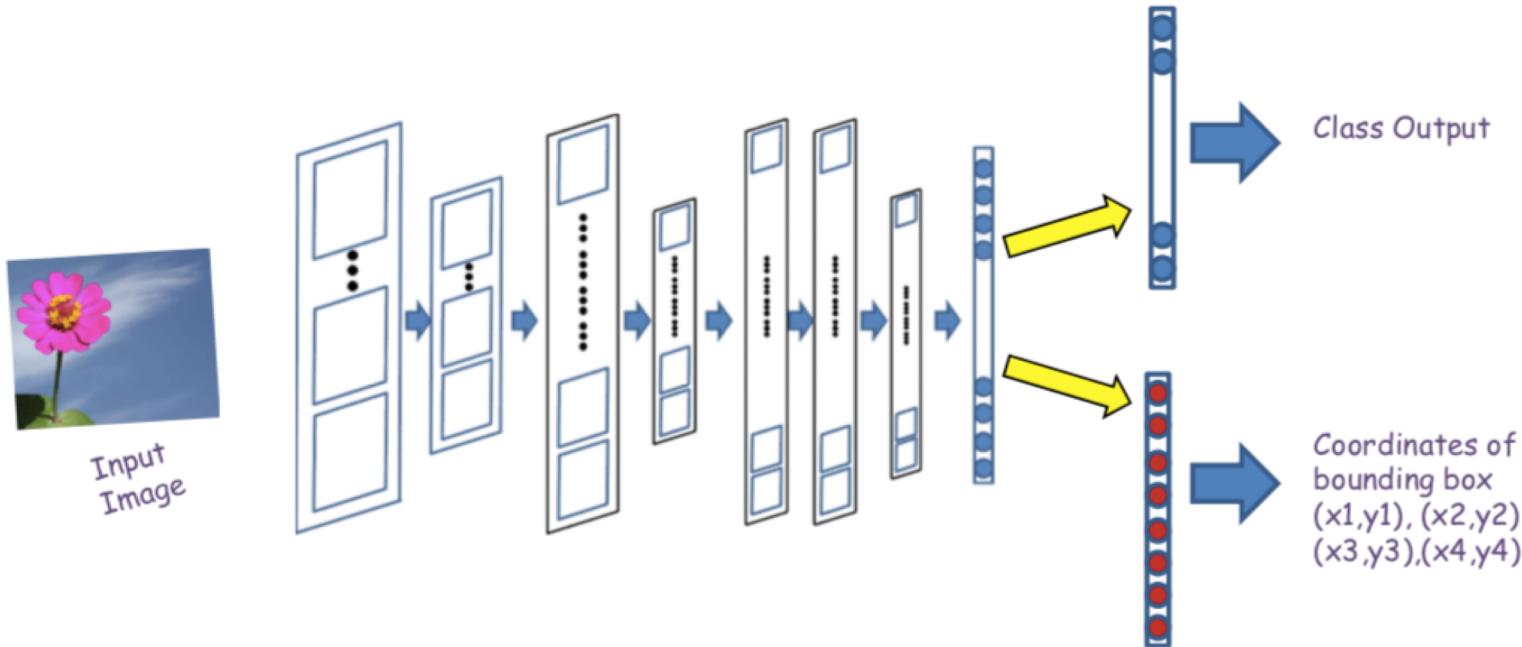
From Classification to Instance Segmentation



[Ouaknine, A. \(2018\) Review of DL for Object Detection Blog Post](#)

Classification + Localization

Learn to predict object **and** bounding box coordinates (Multitask learning)



Pose Estimation/Landmark detection

Learn to predict object **and** landmark coordinates (Multitask learning)

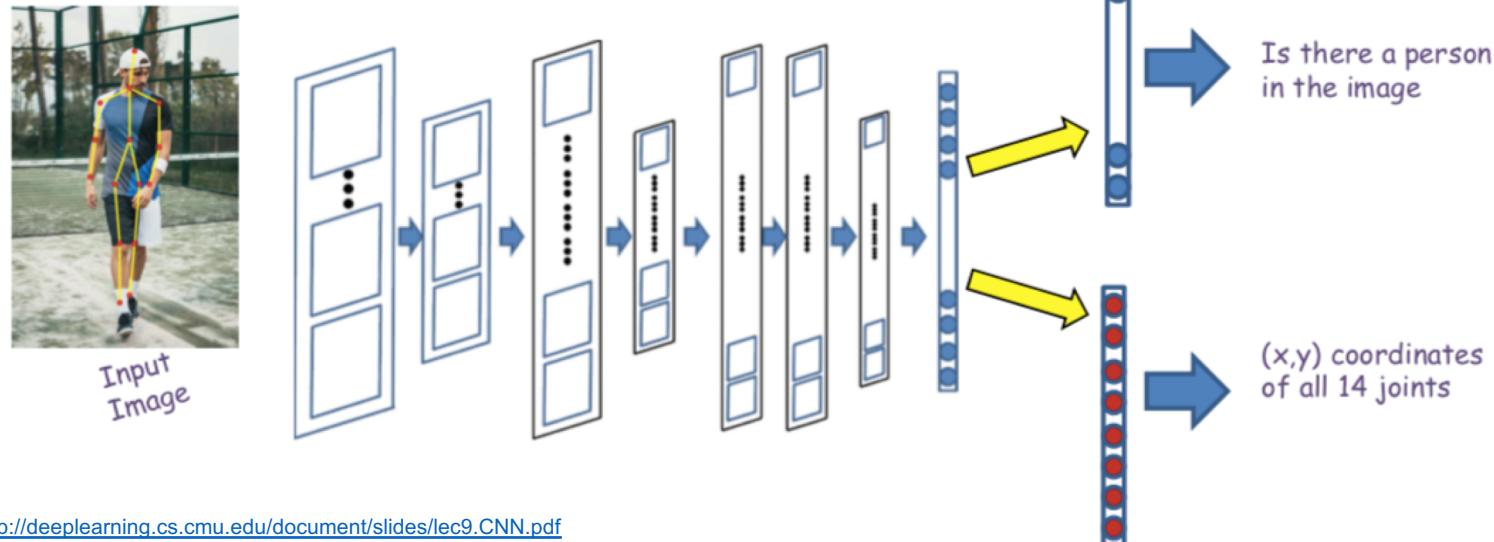
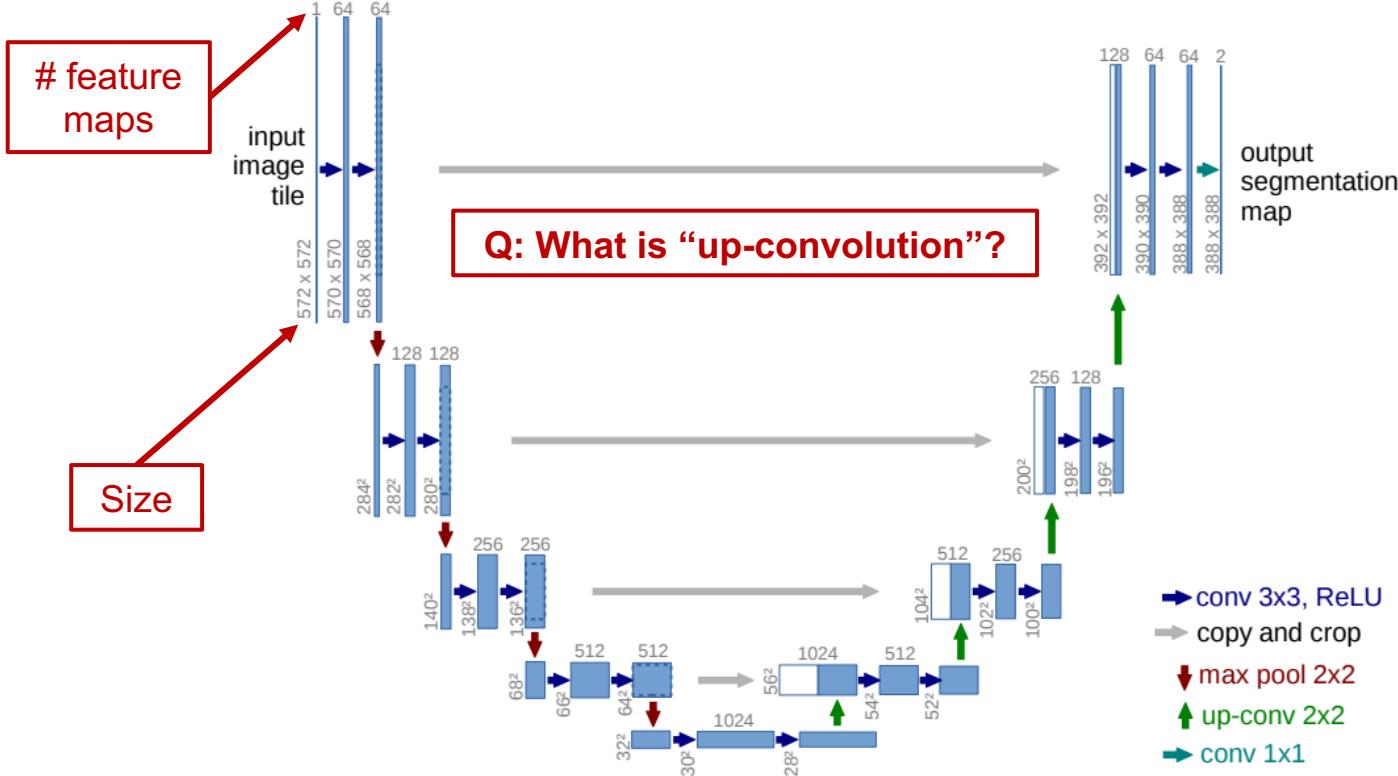


Image Segmentation



Ronneberger, O., Fischer, P., & Brox, T. (2015, October). U-net: Convolutional networks for biomedical image segmentation. MICCAI (pp. 234-241). Springer, Cham.

Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
Multilayer Perceptron
Fully Connected (Dense) Layers
Shared Weight Layers,
Convolutional Layers
Pooling Layers

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid
Tanh
ReLU
eLu
PreLu
MaxOut

Optimization

Gradient Descent
Backpropagation
Momentum methods
Adaptive gradient methods

Tricks of the trade

Dropout
Batch Normalization
Initialization
Data Augmentation
Transfer Learning

Our Deep Learning Toolbox

Components

Linear classifiers $f(\mathbf{x}, \mathbf{W}) = \mathbf{W} \mathbf{x} + \mathbf{b}$
Multilayer Perceptron
Fully Connected (Dense) Layers
Shared Weight Layers,
Convolutional Layers
Pooling Layers
Residual Layers

Loss Functions

Mean Squared Error
Negative Log Likelihood
Hinge Loss
Norm regularization

Activations

Softmax
Sigmoid
Tanh
ReLU
eLu
PreLu
MaxOut

Optimization

Gradient Descent
Backpropagation
Momentum methods
Adaptive gradient methods

Tricks of the trade

Dropout
Batch Normalization
Initialization
Data Augmentation
Transfer Learning
Multitask Learning

CV Architectures

LeNet-5
AlexNet
VGG
RestNet
U-net

Sources

Content adapted from materials created by Mathias Unberath for the
Machine Learning: Deep Learning (Fall 2019) course

Images from Stanford cs231n.

