

MSIS549: ML and AI for Business Applications

Lesson 3



Course Outline

Lesson 1: Introduction to Deep Learning

Lesson 2: Neural Network Fundamentals

Lesson 3: Convolutional Neural Networks

Lesson 4: Recurrent Neural Networks

Lesson 5: Deep Learning Practice

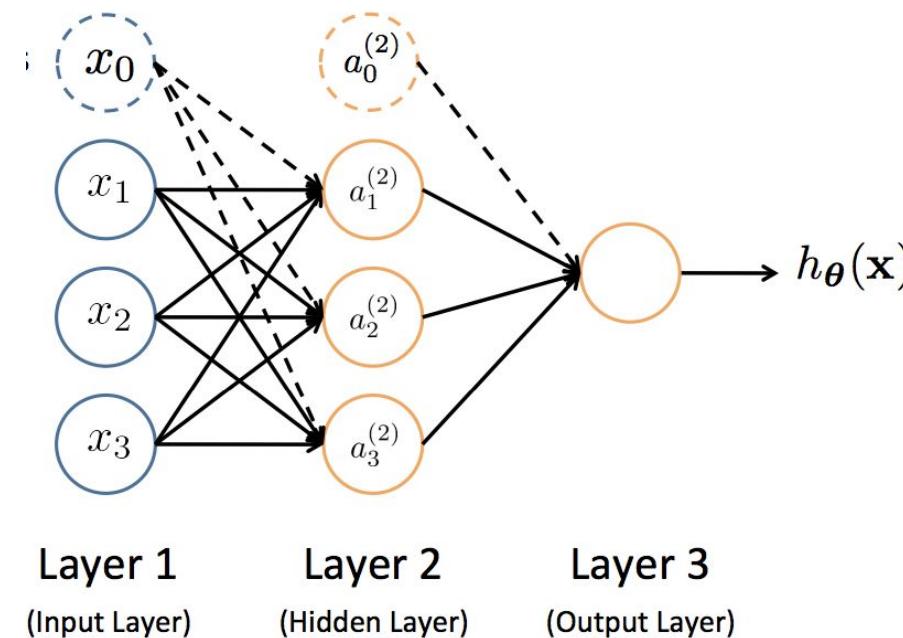
Recap



Neural Networks

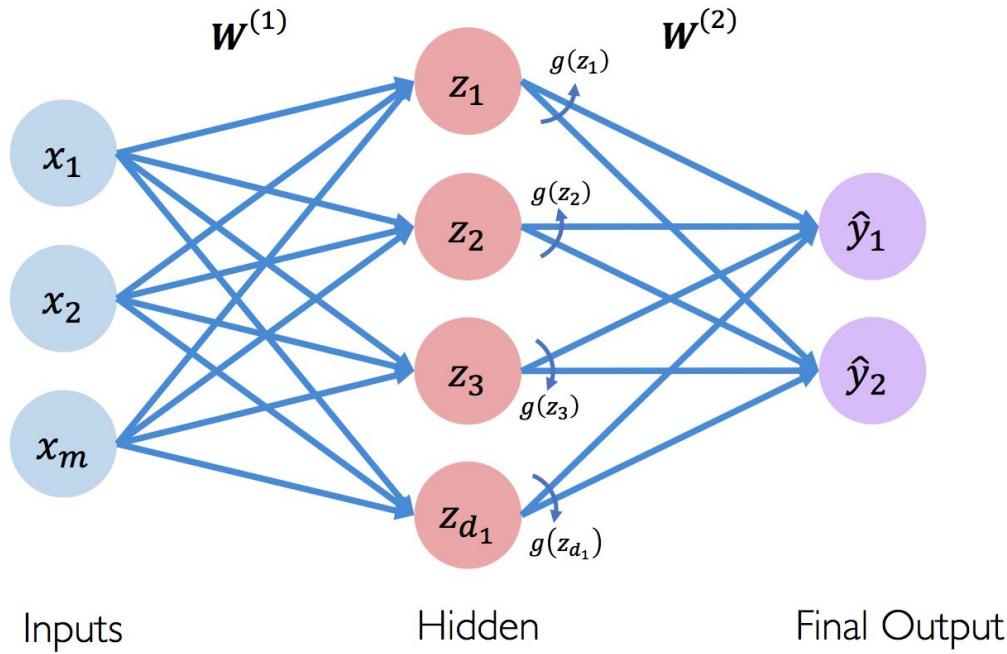
- Origins: Algorithms that try to mimic the brain.
- Very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many applications
- Artificial neural networks are not nearly as complex as the actual brain structure

Universal Approximation Theorem

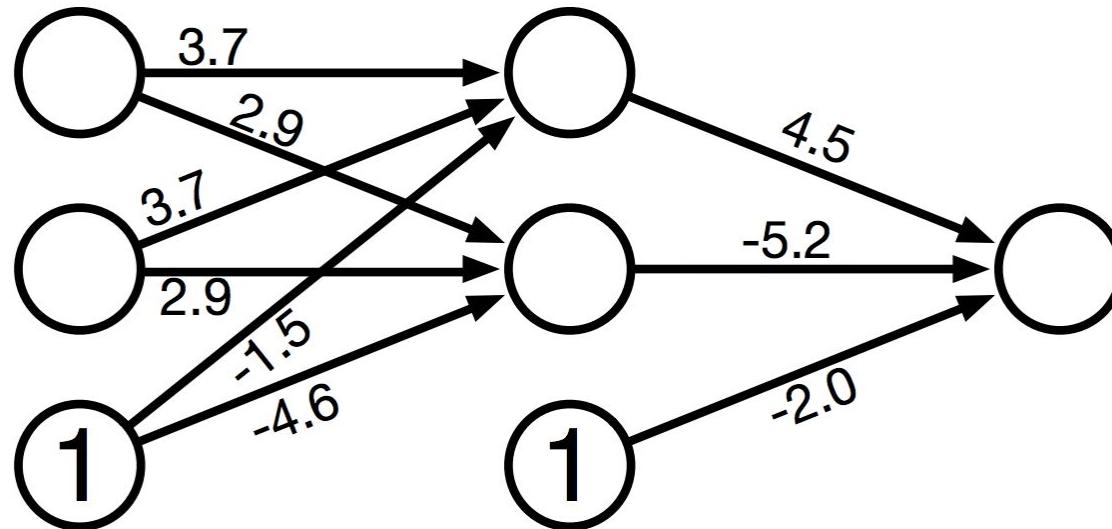


The universal approximation theorem found that a neural network with one hidden layer can approximate any continuous function (George Cybenko, 1989).

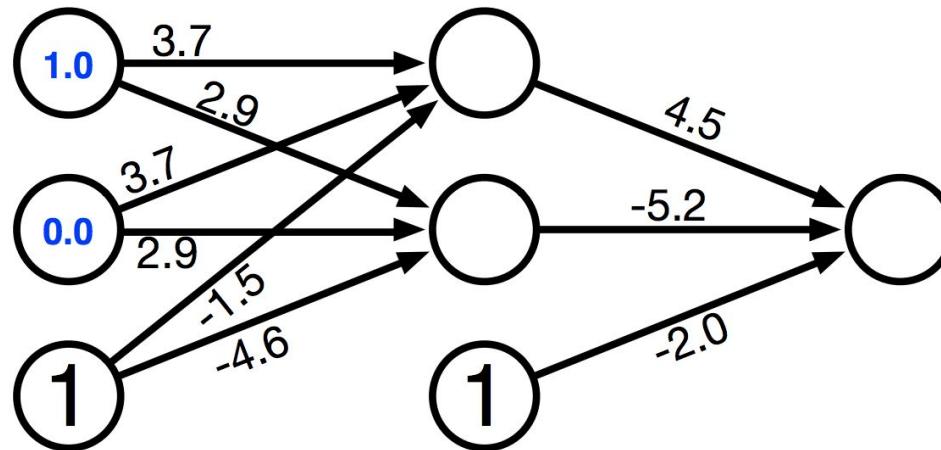
Neural Networks Forward Propagation



Neural Network Example



Neural Network Example

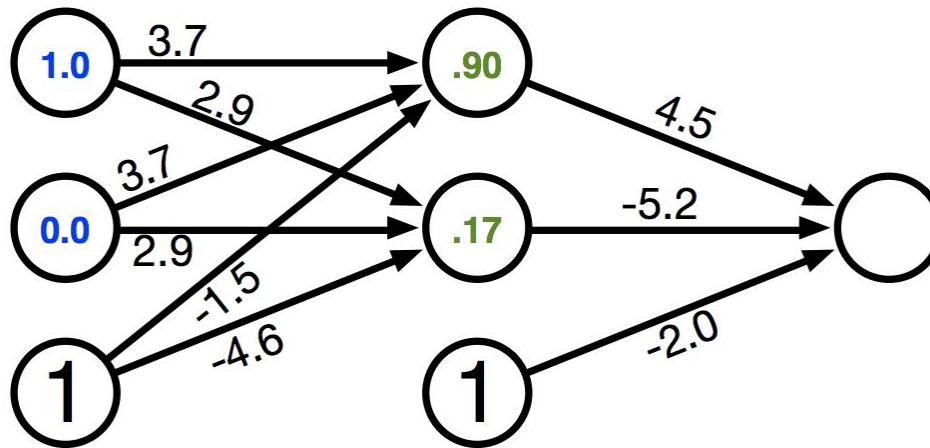


Hidden unit computation

$$\text{sigmoid}(1.0 \times 3.7 + 0.0 \times 3.7 + 1 \times -1.5) = \text{sigmoid}(2.2) = \frac{1}{1 + e^{-2.2}} = 0.90$$

$$\text{sigmoid}(1.0 \times 2.9 + 0.0 \times 2.9 + 1 \times -4.5) = \text{sigmoid}(-1.6) = \frac{1}{1 + e^{1.6}} = 0.17$$

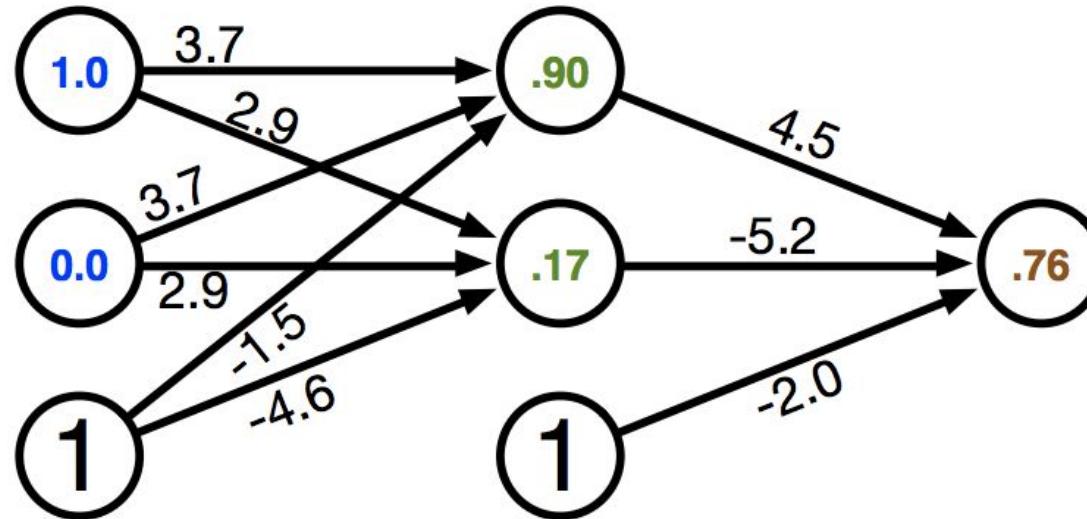
Neural Network Example



Output computation

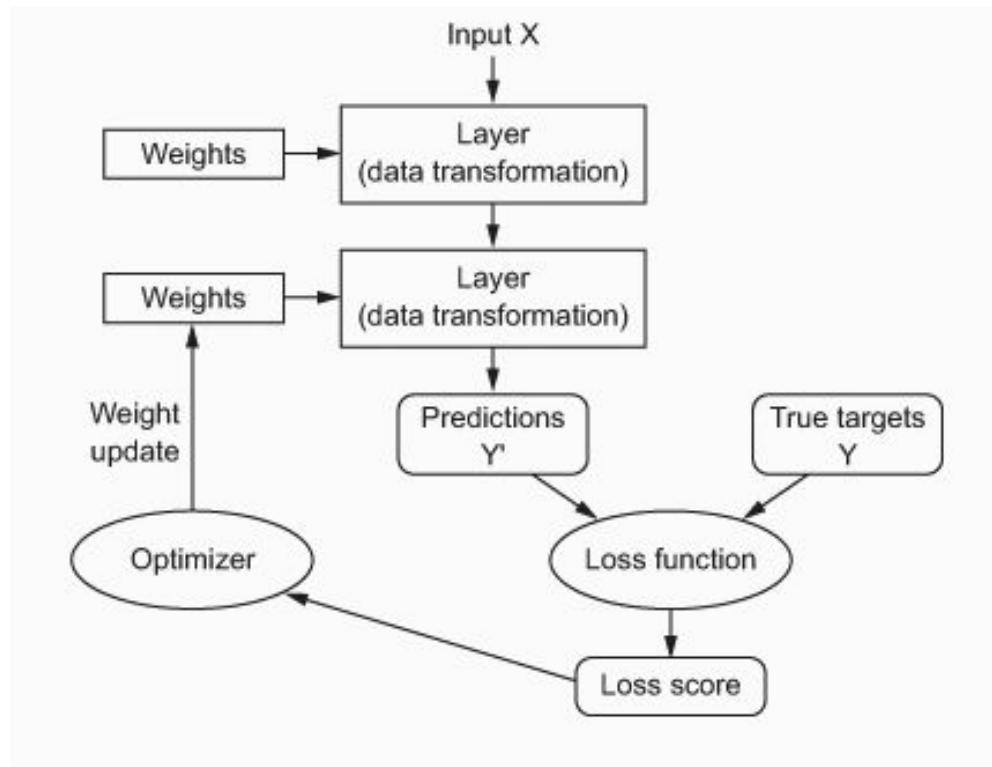
$$\text{sigmoid}(.90 \times 4.5 + .17 \times -5.2 + 1 \times -2.0) = \text{sigmoid}(1.17) = \frac{1}{1 + e^{-1.17}} = 0.76$$

Neural Network Example



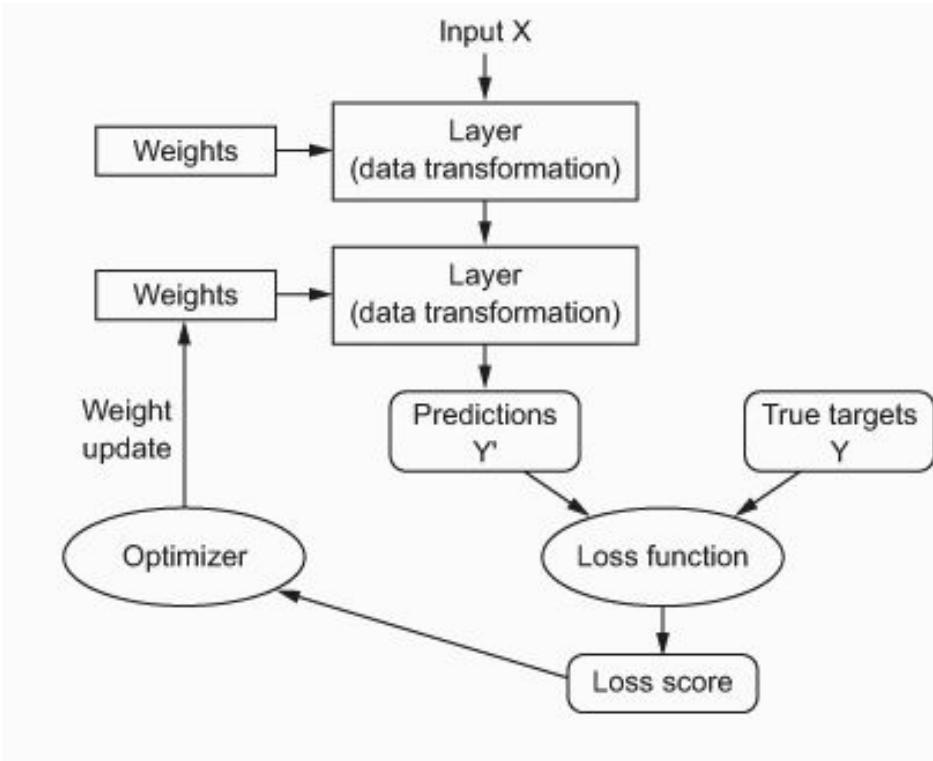
Training Neural Networks

Learning means finding a combination of model parameters that minimizes a loss function for a given set of training data samples and their corresponding targets.



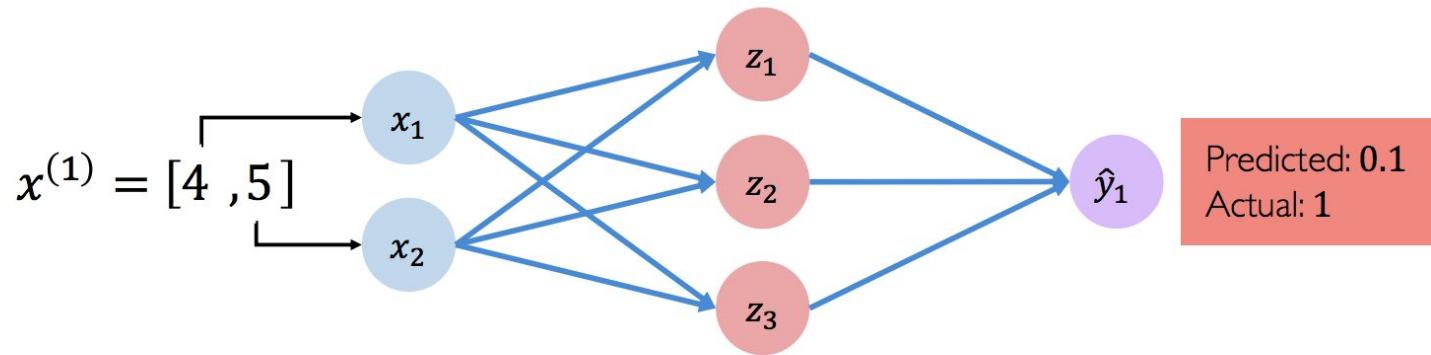
Loss Function

The **loss** is the quantity we attempt to minimize during training, so it should represent a measure of success for the task you're trying to solve.



Quantifying Loss

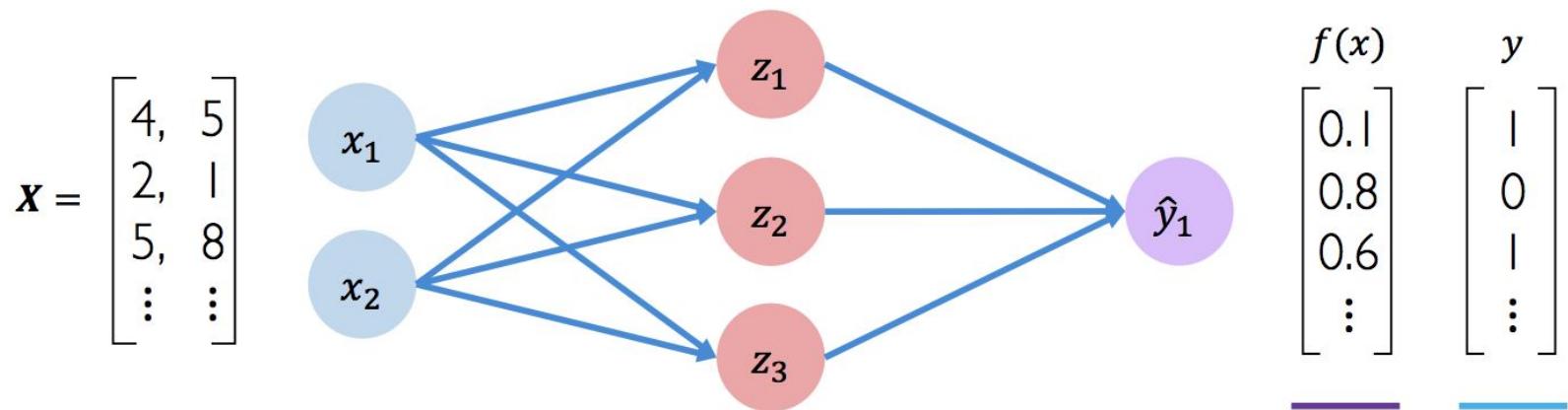
The **loss** of our network measures the cost incurred from incorrect predictions



$$\mathcal{L} \left(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}} \right)$$

Empirical Loss

The **empirical loss** measures the total loss over our entire dataset



Also known as:

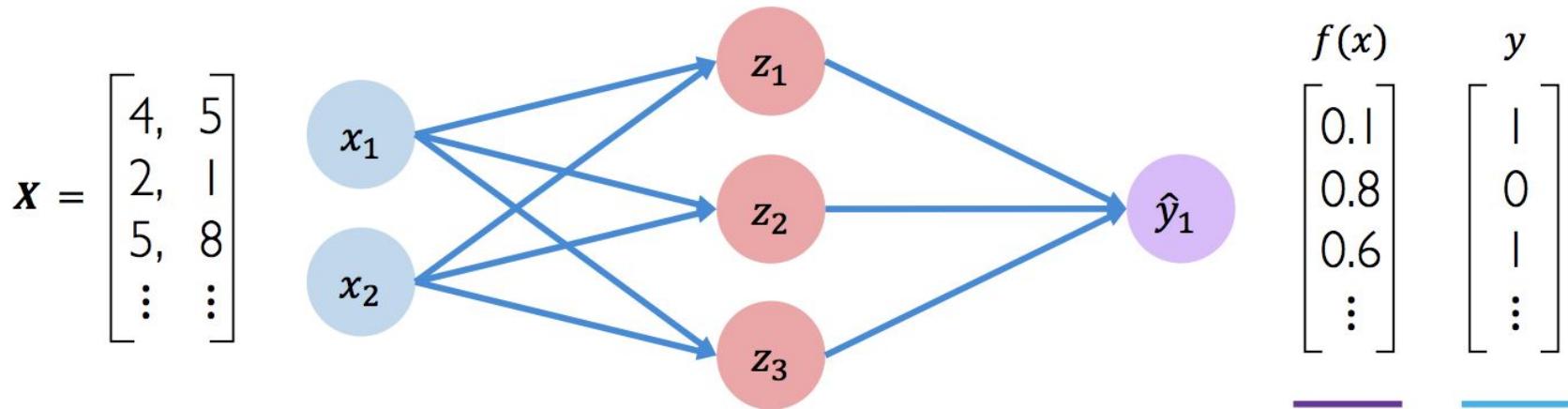
- Objective function
- Cost function
- Empirical Risk

$J(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$

Predicted Actual

Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1

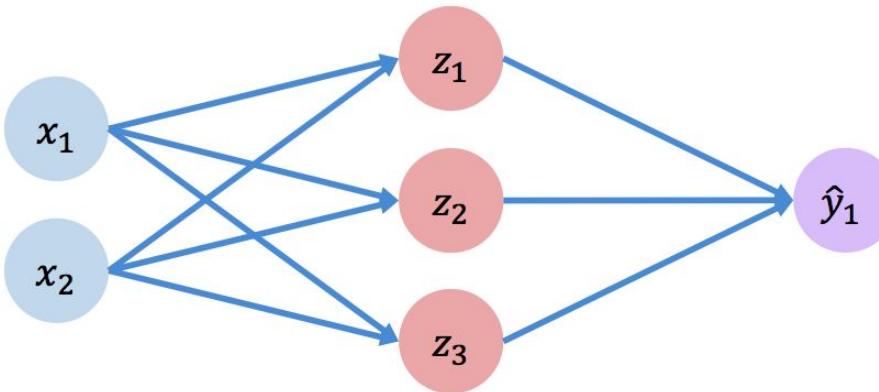


$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} \quad \underbrace{\text{Predicted}}_{\text{Predicted}}$$

Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$\begin{array}{c|c} f(x) & y \\ \hline 30 & 90 \\ 80 & 20 \\ 85 & 95 \\ \vdots & \vdots \end{array}$$

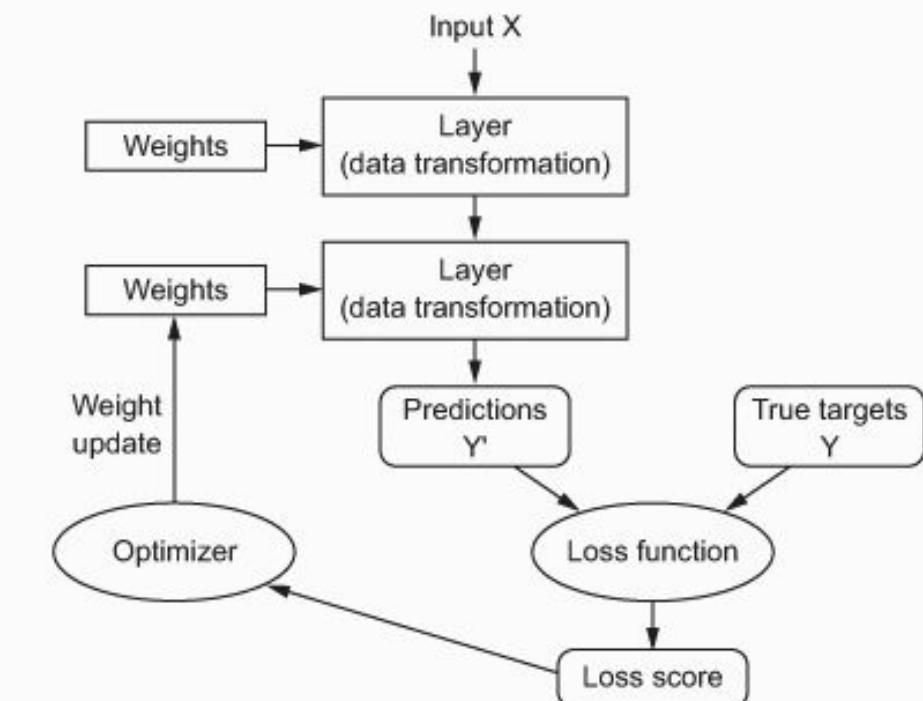
Final Grades
(percentage)

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \left(\underline{y^{(i)}} - \underline{f(x^{(i)}; \mathbf{W})} \right)^2$$

Actual Predicted

Optimizer

The **optimizer** specifies the exact way in which the gradient of the loss will be used to update parameters.



Loss Optimization

We want to find the network weights that achieve the lowest loss.

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:

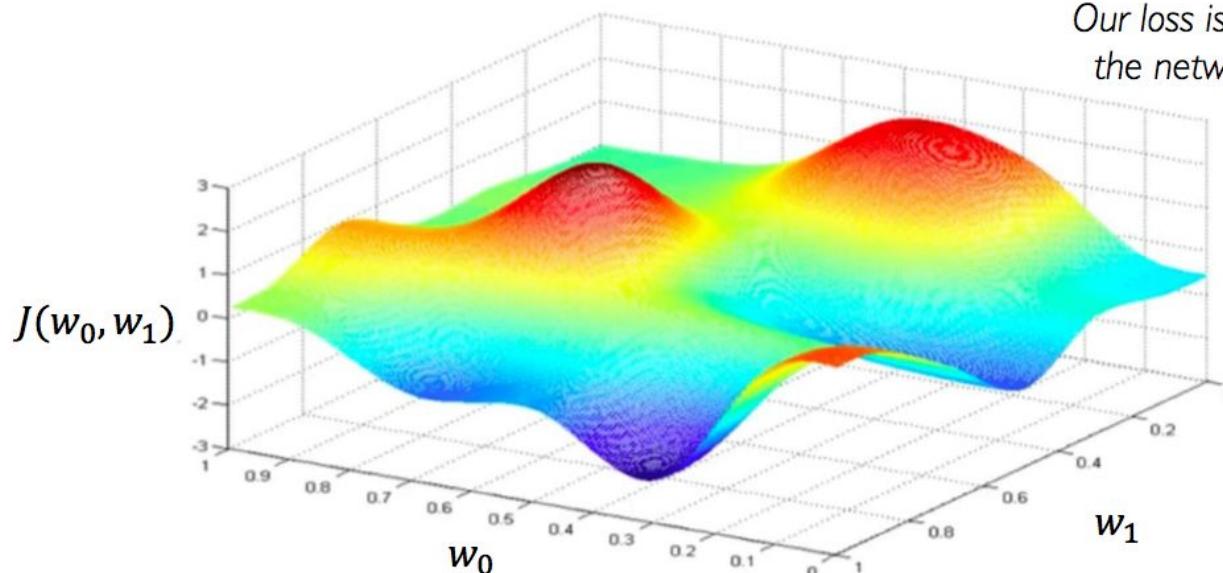
$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

Loss Optimization

Find the network weights that achieve the lowest loss.

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

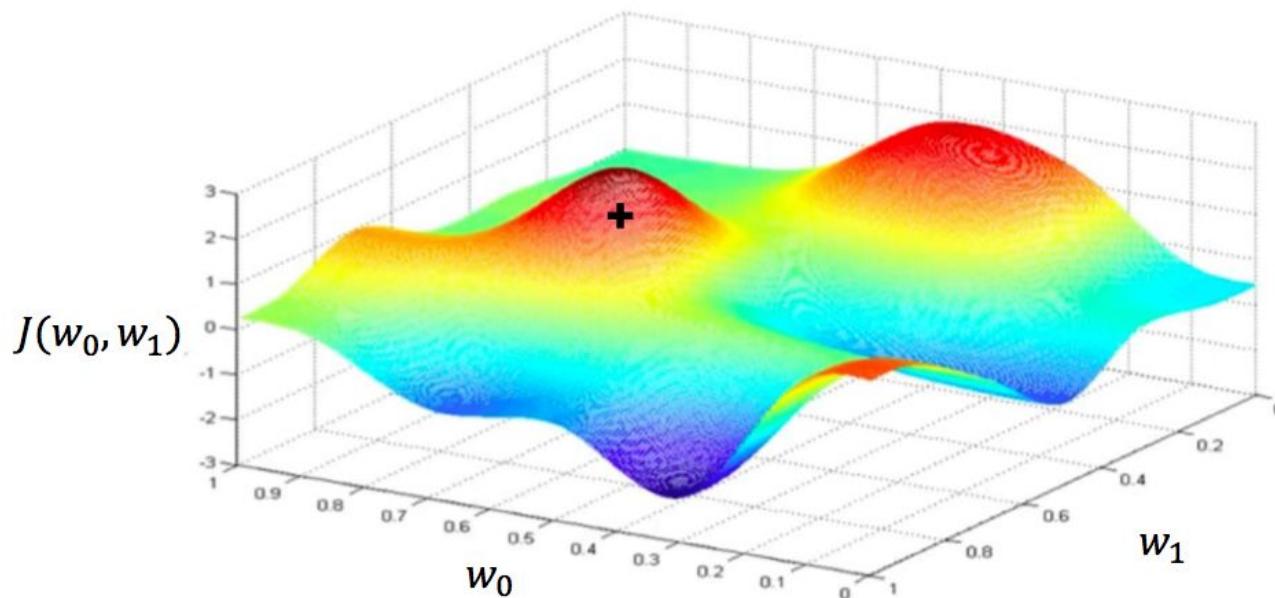
Remember:
Our loss is a function of
the network weights!



Loss Optimization

Find the network weights that achieve the lowest loss.

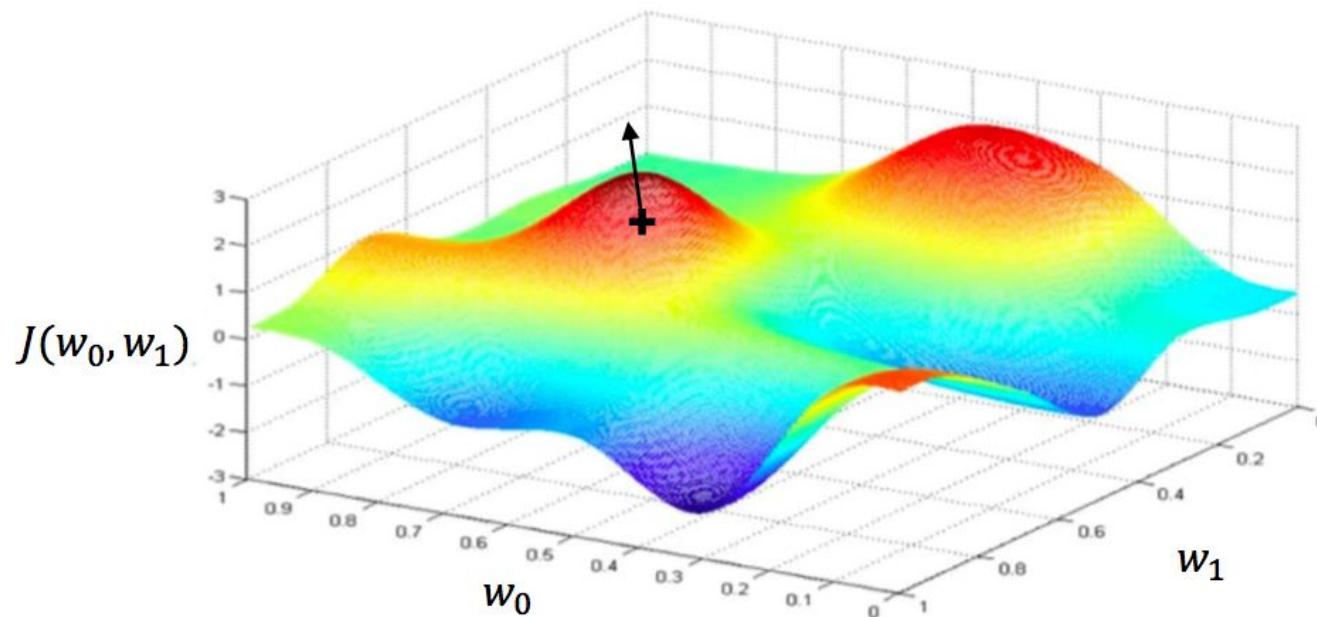
Randomly pick an initial (w_0, w_1)



Loss Optimization

Find the network weights that achieve the lowest loss.

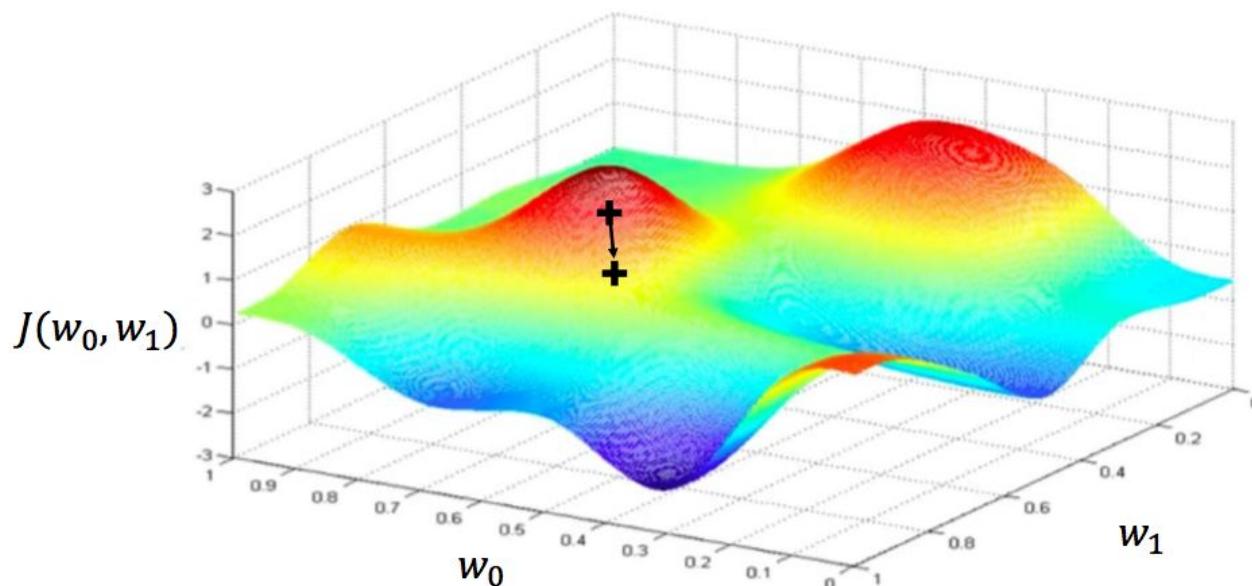
Compute gradient, $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$



Loss Optimization

Find the network weights that achieve the lowest loss.

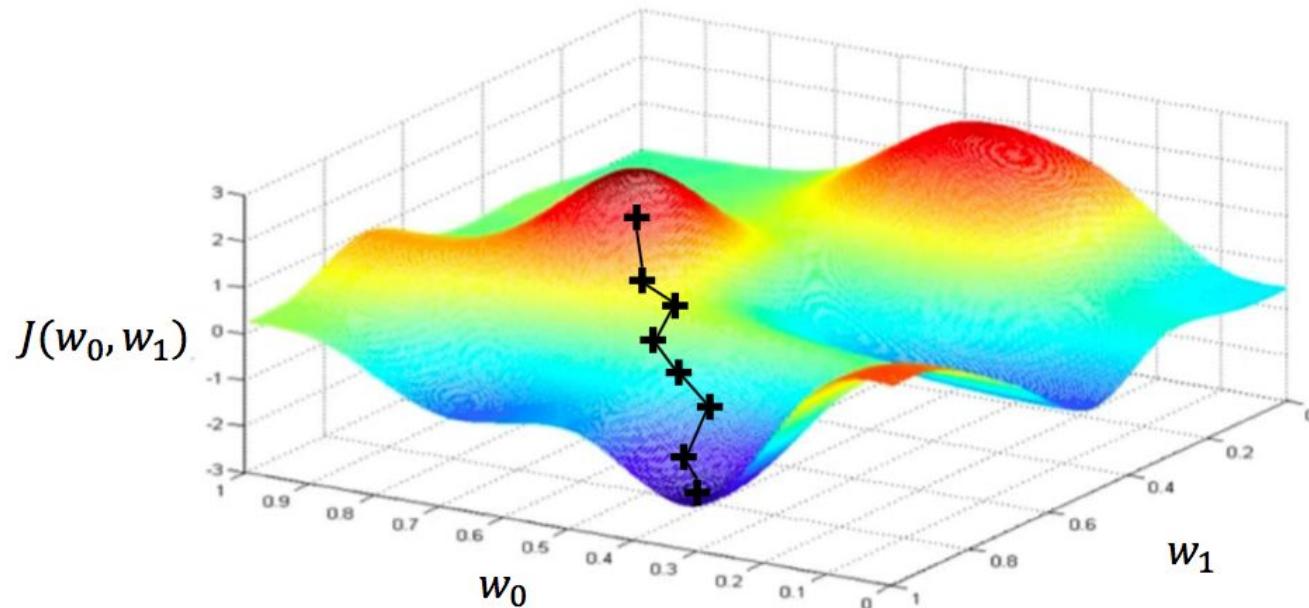
Take small step in opposite direction of gradient



Loss Optimization

Find the network weights that achieve the lowest loss.

Repeat until convergence



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(W)}{\partial W}$
4. Update weights, $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights

```
 weights = tf.random_normal(shape, stddev=sigma)
```

```
 grads = tf.gradients(ys=loss, xs=weights)
```

```
 weights_new = weights.assign(weights - lr * grads)
```

Training Neural Networks: Adaptive learning



Learning Rate

Algorithm

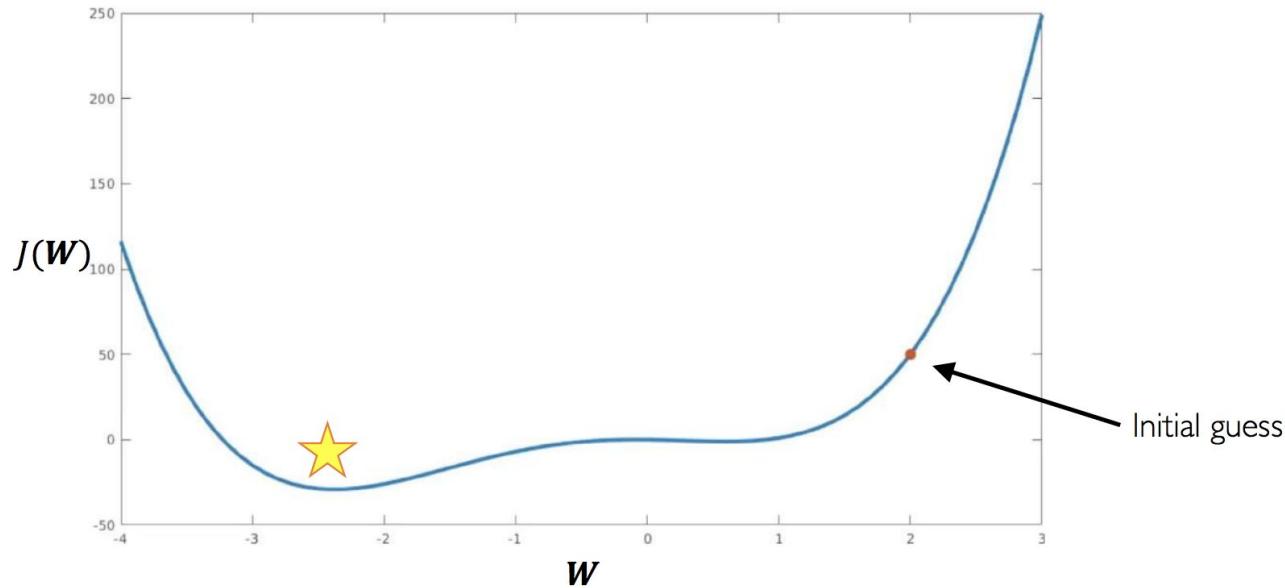
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Optimization through gradient descent

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

How can we set the learning rate?

Setting the Learning Rate



Small learning rate converges slowly and gets stuck in false local minima

Large learning rates overshoot, become unstable and diverge

Stable learning rates converge smoothly and avoid local minima

Adaptive Learning Rate Algorithms

- Momentum
- Adagrad
- Adadelta
- Adam
- RMSProp

 `tf.train.MomentumOptimizer`

 `tf.train.AdagradOptimizer`

 `tf.train.AdadeltaOptimizer`

 `tf.train.AdamOptimizer`

 `tf.train.RMSPropOptimizer`

Qian et al. "On the momentum term in gradient descent learning algorithms." 1999.

Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.

Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.

Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.

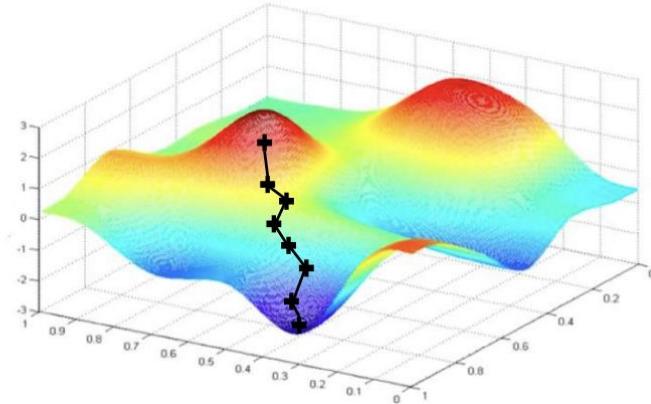
Additional details: <http://ruder.io/optimizing-gradient-descent/>

Batch Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Can be very computational to compute!

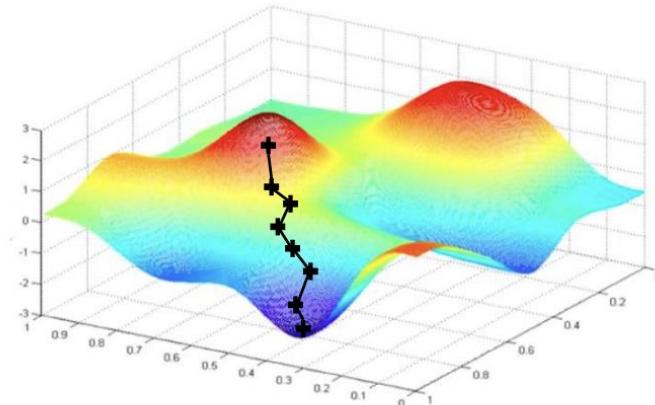


Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

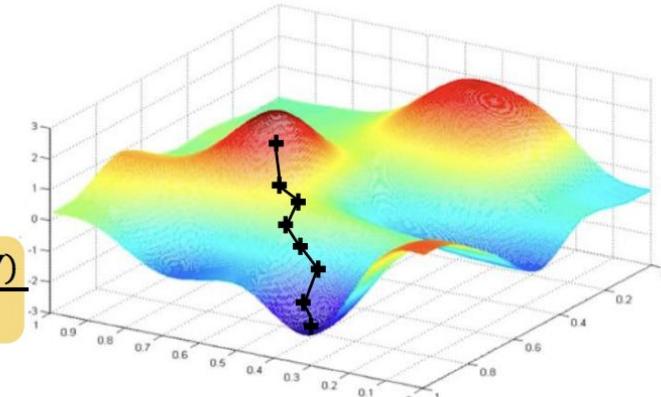
Easy to compute but
very noisy
(stochastic)!



Mini-Batch Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$
5. Update weights,
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$
6. Return weights

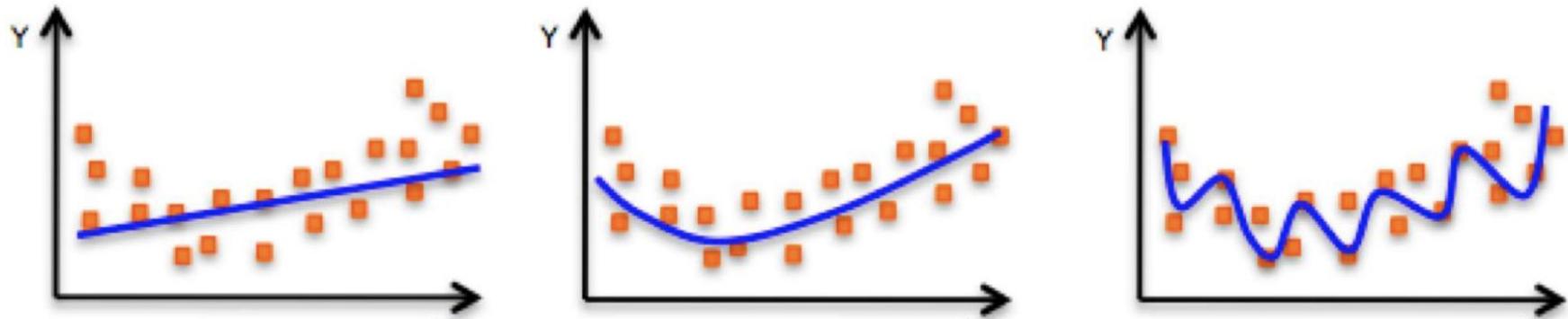


Fast to compute and a much better estimate of the true gradient!

Gradient Descent

Method	Accuracy	Update Speed	Memory Usage	Online Learning
Batch gradient descent	Good	Slow	High	No
Stochastic gradient descent	Good (with annealing)	High	Low	Yes
Mini-batch gradient descent	Good	Medium	Medium	Yes

The Problem of Overfitting



Underfitting

Model does not have capacity
to fully learn the data

Ideal fit

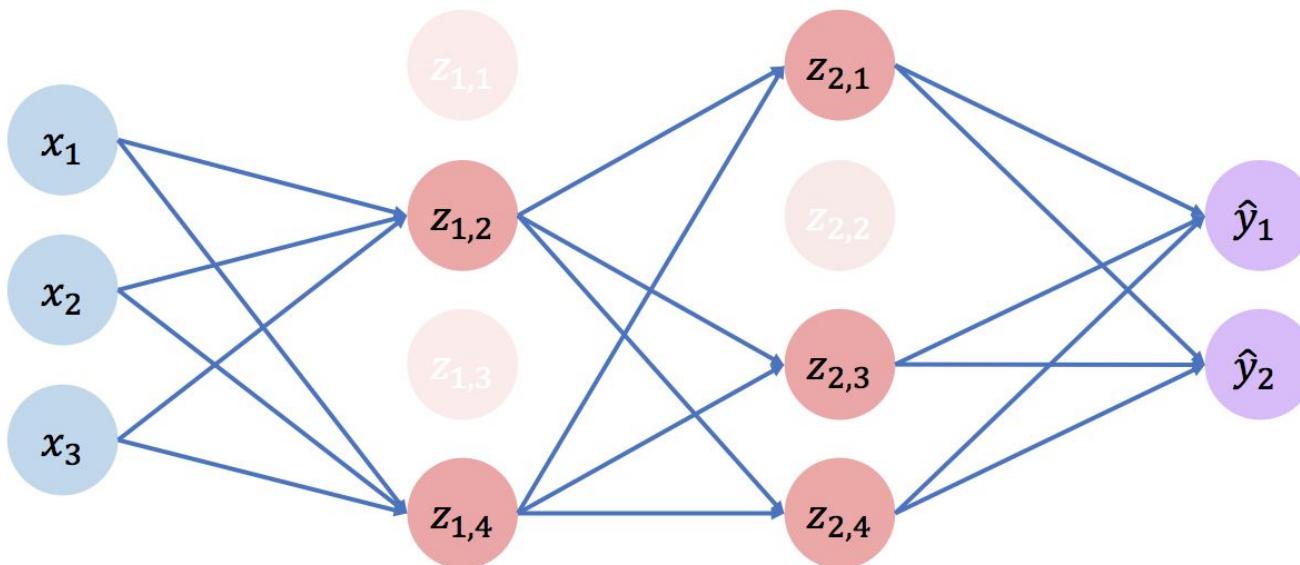
Overfitting

Too complex, extra parameters,
does not generalize well

Regularization: Dropout

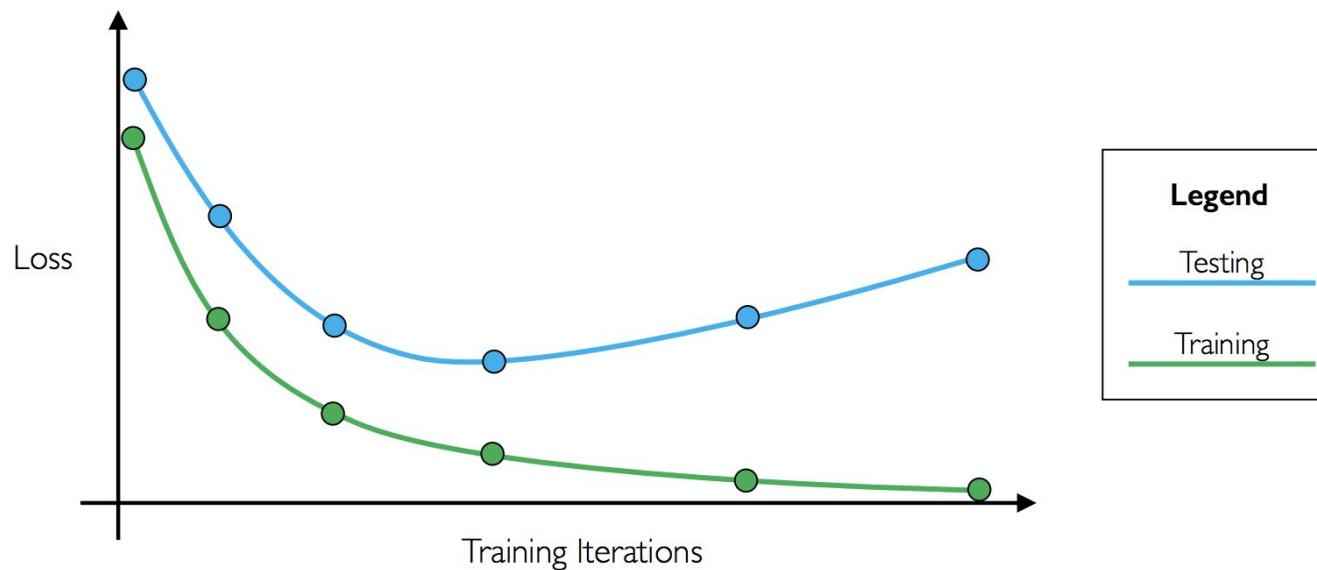
During training, randomly set some activations to 0

- Typically ‘drop’ 50% of activations in layer
- Forces network to not rely on any 1 node



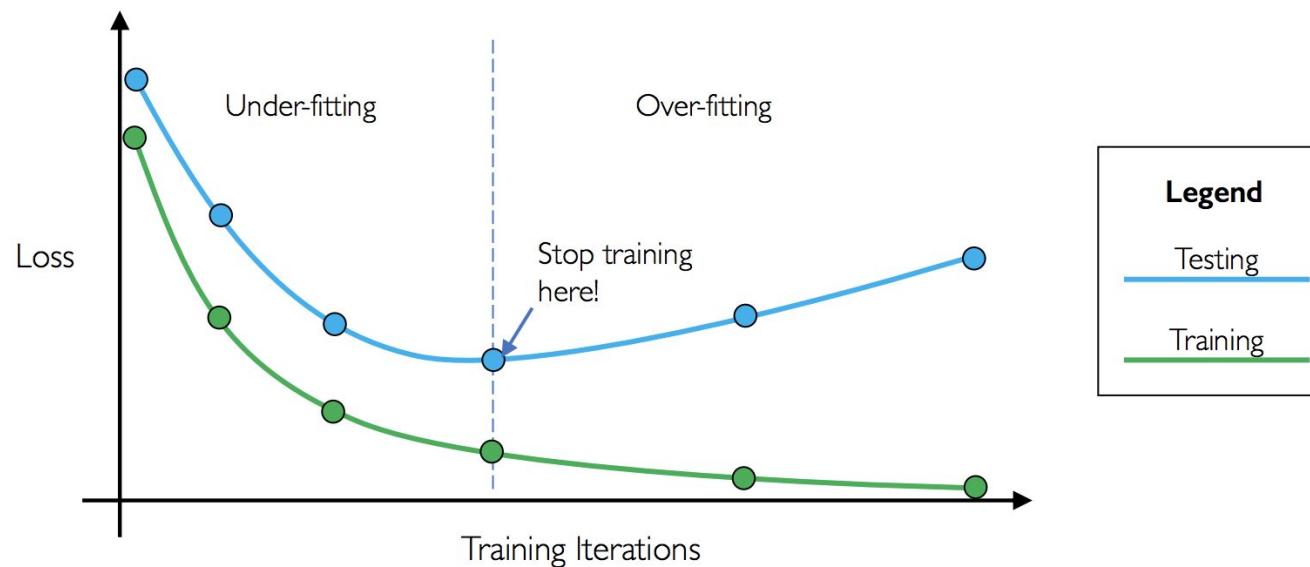
Regularization: Early Stopping

Stop training before we have a chance to overfit



Regularization: Early Stopping

Stop training before we have a chance to overfit



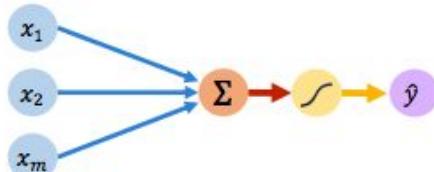
Last-layer Activation and Loss Function

Problem type	Last-layer activation	Loss function
Binary classification	sigmoid	binary_crossentropy
Multiclass, single-label classification	softmax	categorical_crossentropy
Multiclass, multilabel classification	sigmoid	binary_crossentropy
Regression to arbitrary values	None	mse
Regression to values between 0 and 1	sigmoid	mse or binary_crossentropy

Lecture Recap

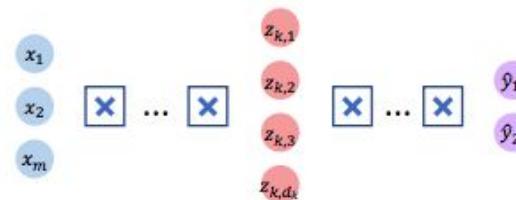
The Perceptron

- Structural building blocks
- Nonlinear activation functions



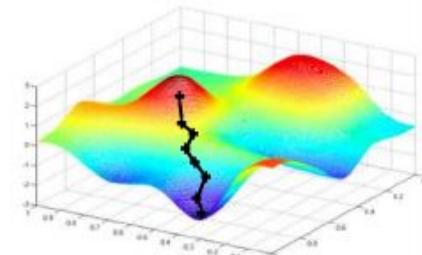
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

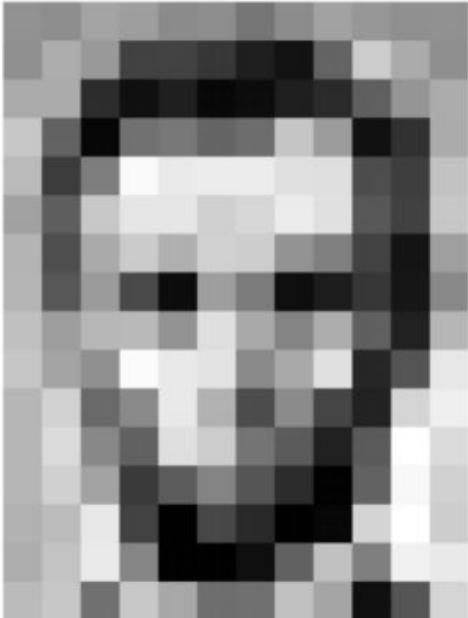
- Adaptive learning
- Batching
- Regularization



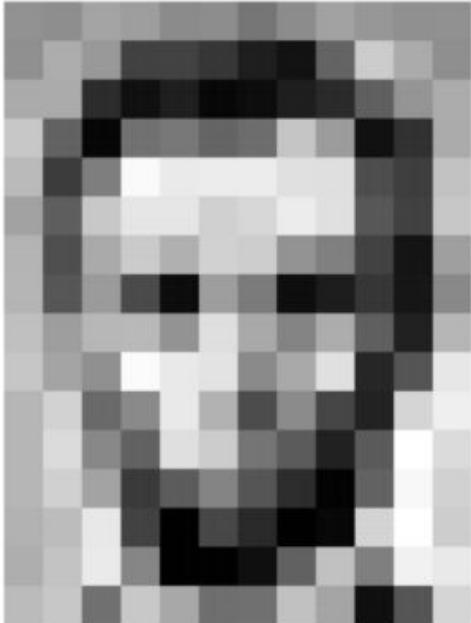
Convolutional Neural Networks



What Computers “See”

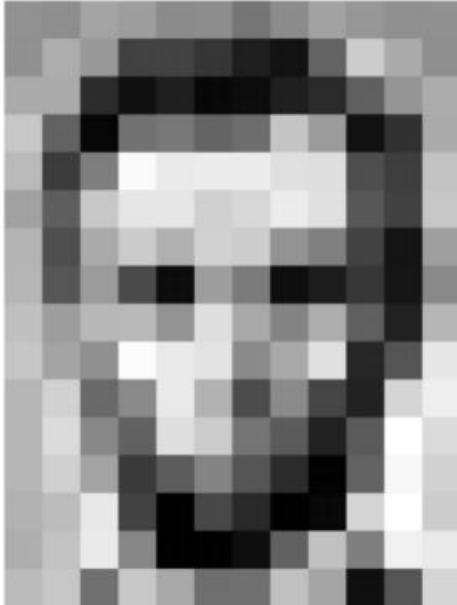


What Computers “See”



157	163	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	93	17	110	210	180	154
180	180	50	14	94	6	10	93	48	106	159	181
206	109	5	124	191	111	120	204	166	15	56	180
194	68	137	251	257	299	299	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	68	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	179	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	106	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	95	218

What Computers “See”



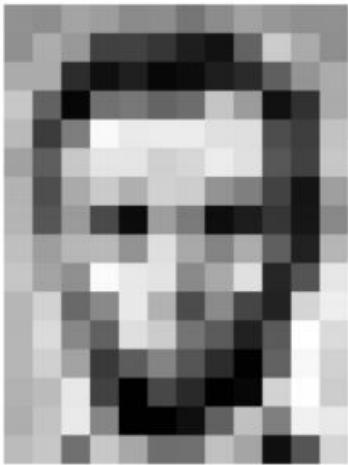
157	153	174	148	150	162	129	151	172	161	165	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	128	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

What the computer sees

157	153	174	168	150	152	129	151	172	161	165	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	128	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	146	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

Tasks in Computer Vision

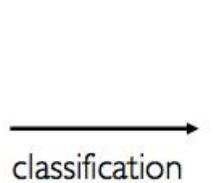


Input Image



167	153	174	168	150	152	129	151	172	161	195	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	56	14	34	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	195	215	211	158	138	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
169	168	191	193	158	227	178	143	182	106	36	190
205	174	195	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	86	190	79	38	218	241
190	224	147	108	227	210	127	102	36	101	295	224
190	214	173	66	103	143	96	50	2	109	240	215
187	196	235	75	1	81	47	0	6	217	295	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Pixel Representation



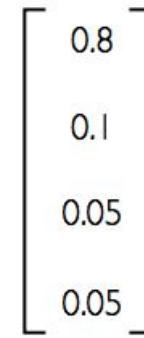
classification

Lincoln

Washington

Jefferson

Obama



- **Regression:** output variable takes continuous value
- **Classification:** output variable takes class label. Can produce probability of belonging to a particular class

High Level Feature Detection

Let's identify key features in each image category



Nose,
Eyes,
Mouth

Wheels,
License Plate,
Headlights

Door,
Windows,
Steps

Manual Feature Extraction

Domain knowledge

Define features

Detect features
to classify



Scale variation

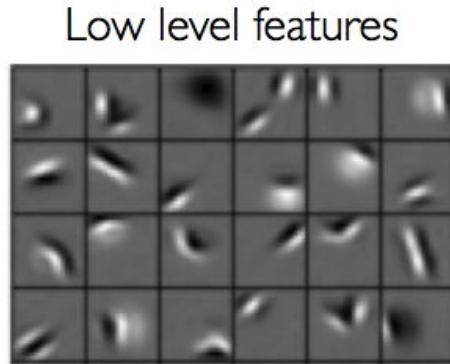


Deformation

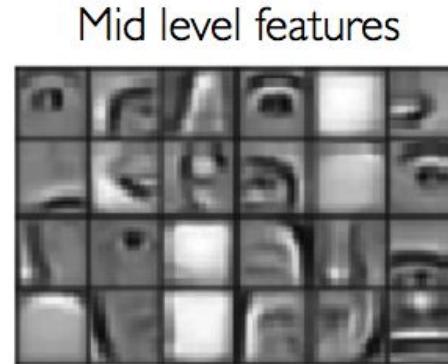


Learning Feature Representations

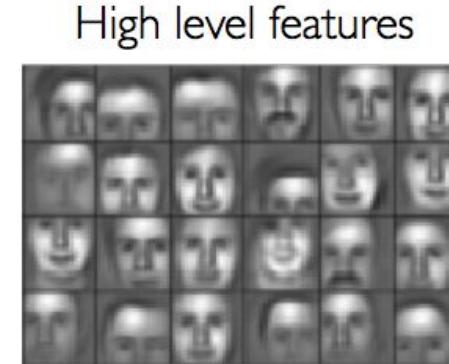
Can we learn a **hierarchy of features** directly from the data instead of hand engineering?



Edges, dark spots



Eyes, ears, nose

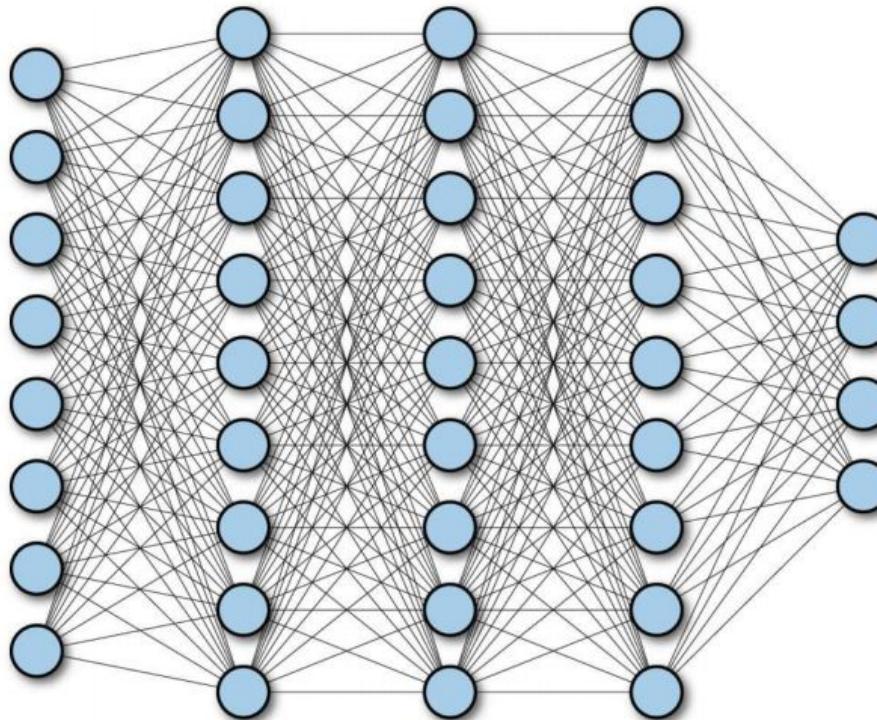


Facial structure

Learning Visual Features



Fully Connected Neural Network



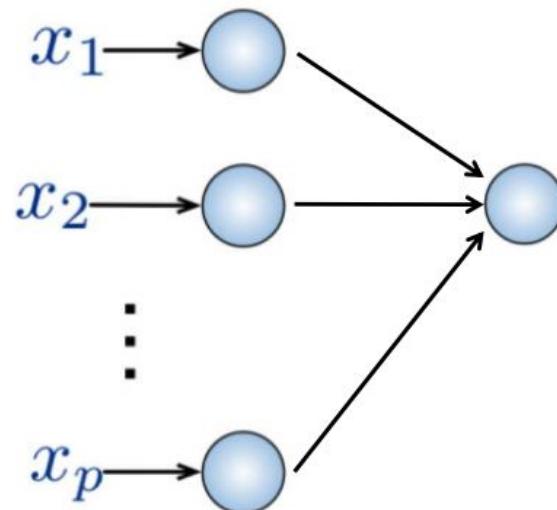


What is the downside of Applying Fully NN to Image Data?

- Too many parameters require a lot of data
- Does not capture the spatial structure

Fully Connected Neural Network

- Input:**
- 2D image
 - Vector of pixel values



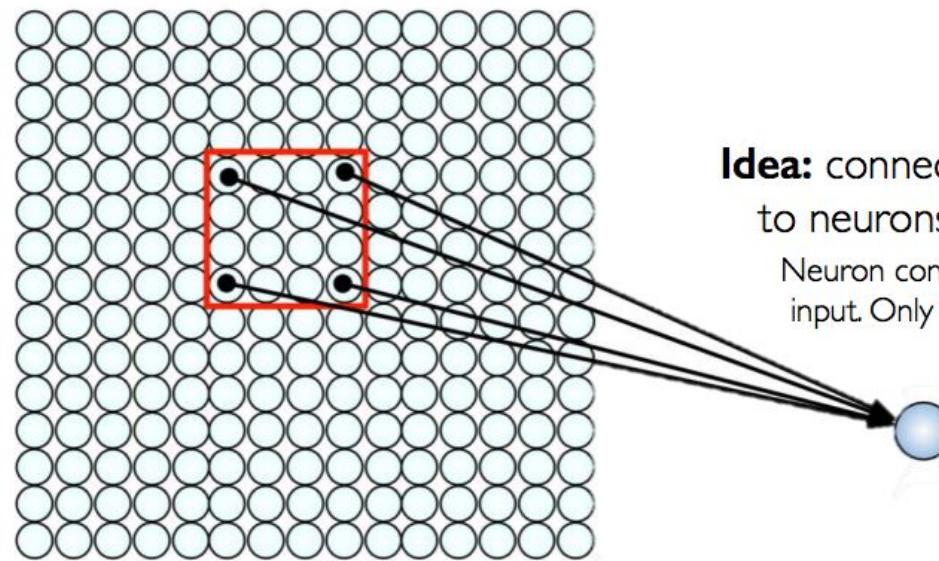
Fully Connected:

- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?

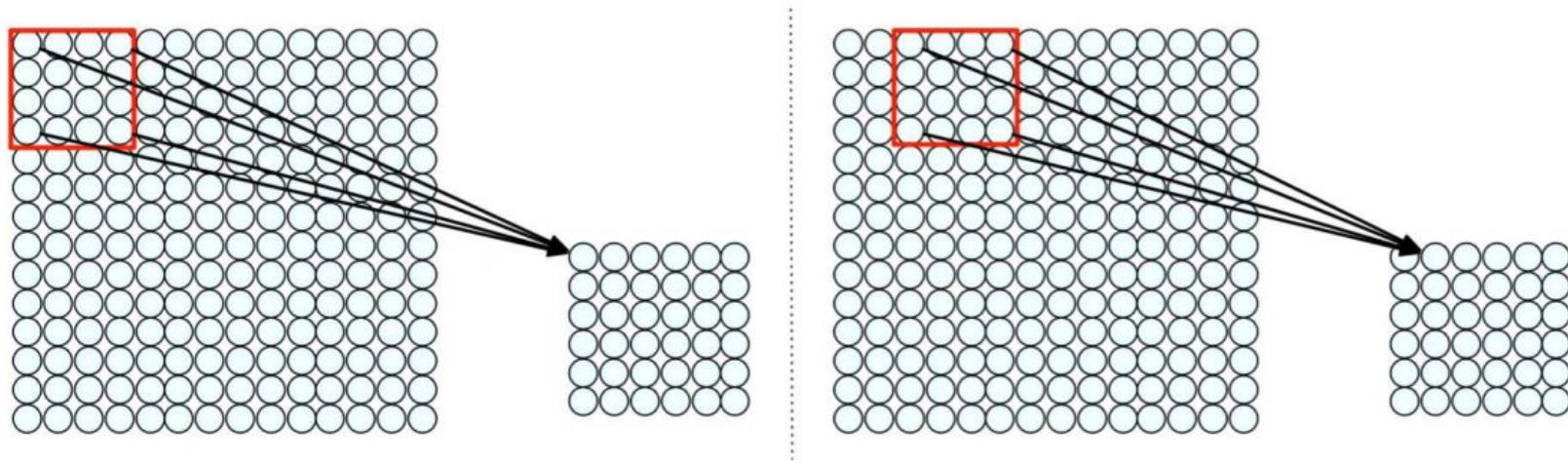
Using Spatial Structure

Input: 2D image.
Array of pixel values



Idea: connect patches of input
to neurons in hidden layer.
Neuron connected to region of
input. Only "sees" these values.

Using Spatial Structure



Connect patch in input layer to a single neuron in subsequent layer.

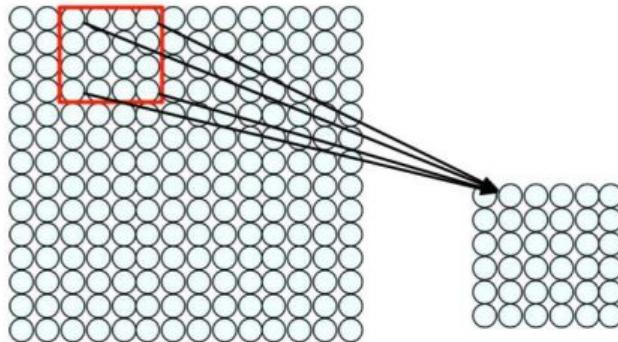
Use a sliding window to define connections.

*How can we **weight** the patch to detect particular features?*

Applying Filters to Extract Features

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) Spatially **share** parameters of each filter
(features that matter in one part of the input should matter elsewhere)

Feature Extraction with Convolution



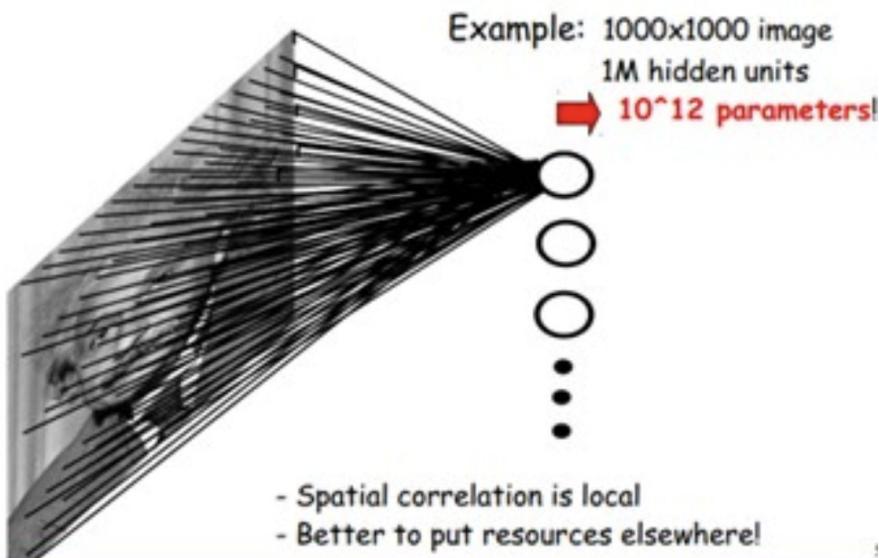
- Filter of size 4×4 : 16 different weights
- Apply this same filter to 4×4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

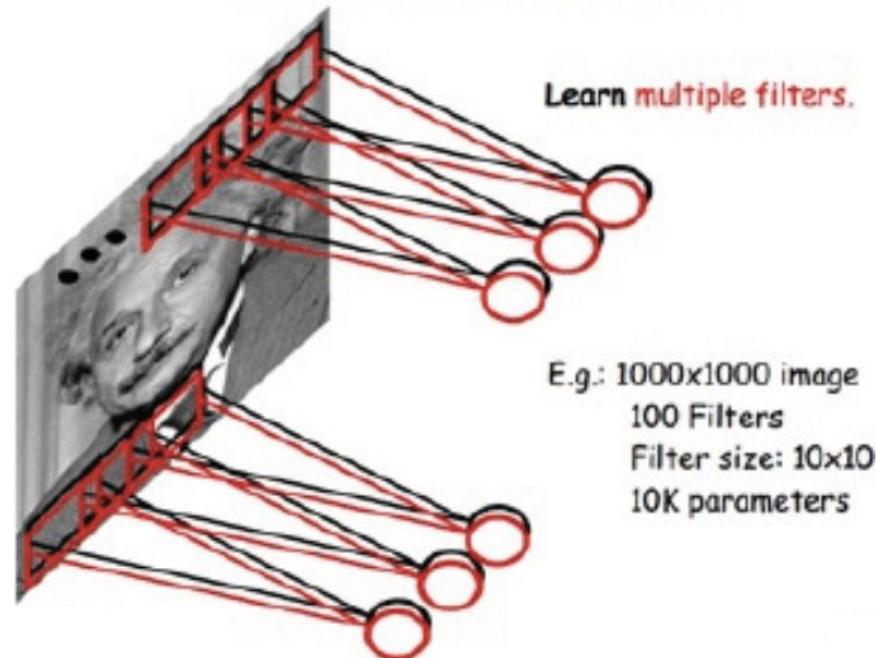
- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Parameter Sharing

FULLY CONNECTED NEURAL NET



CONVOLUTIONAL NET



Feature Extraction and Convolution

A Case Study

X = X?

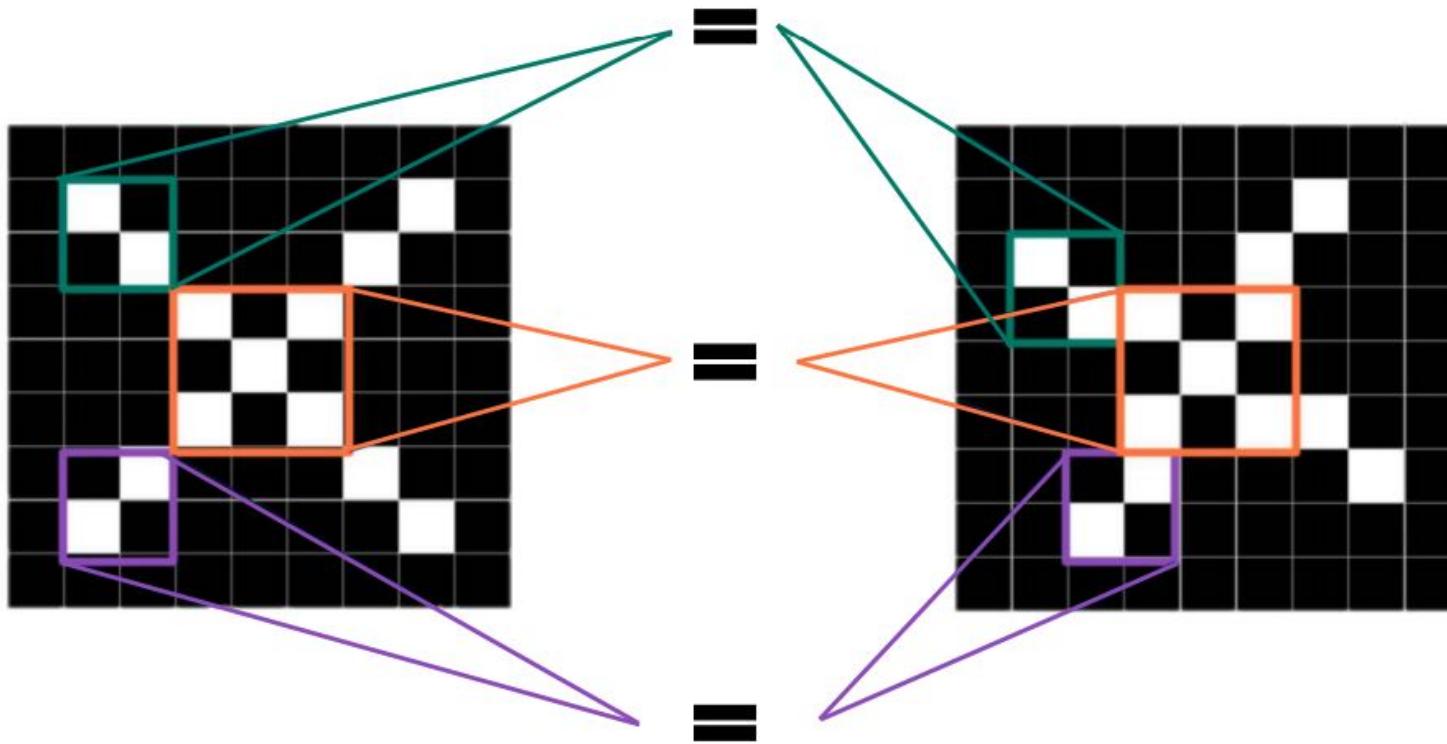
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1



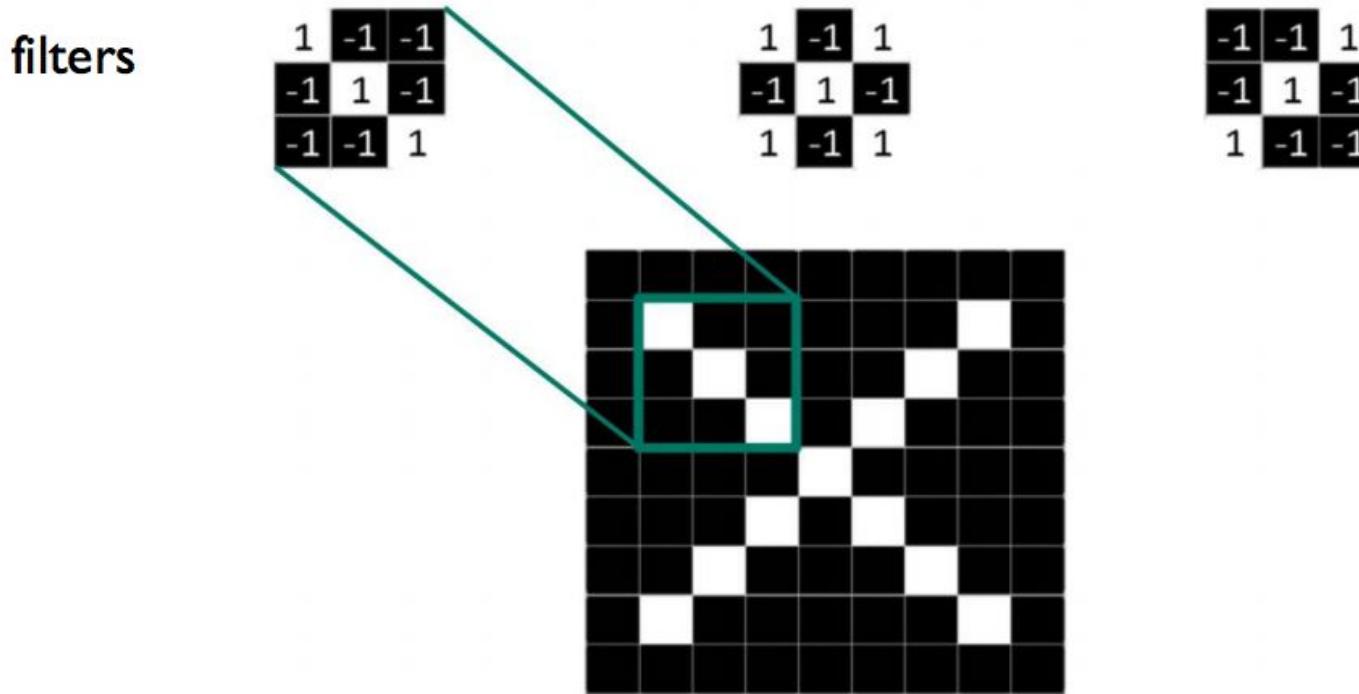
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	1	-1	-1	-1
-1	-1	1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	1	1	-1	-1
-1	-1	-1	1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

Image is represented as matrix of pixel values... and computers are literal!
We want to be able to classify an X as an X even if it's shifted, shrunk, rotated, deformed.

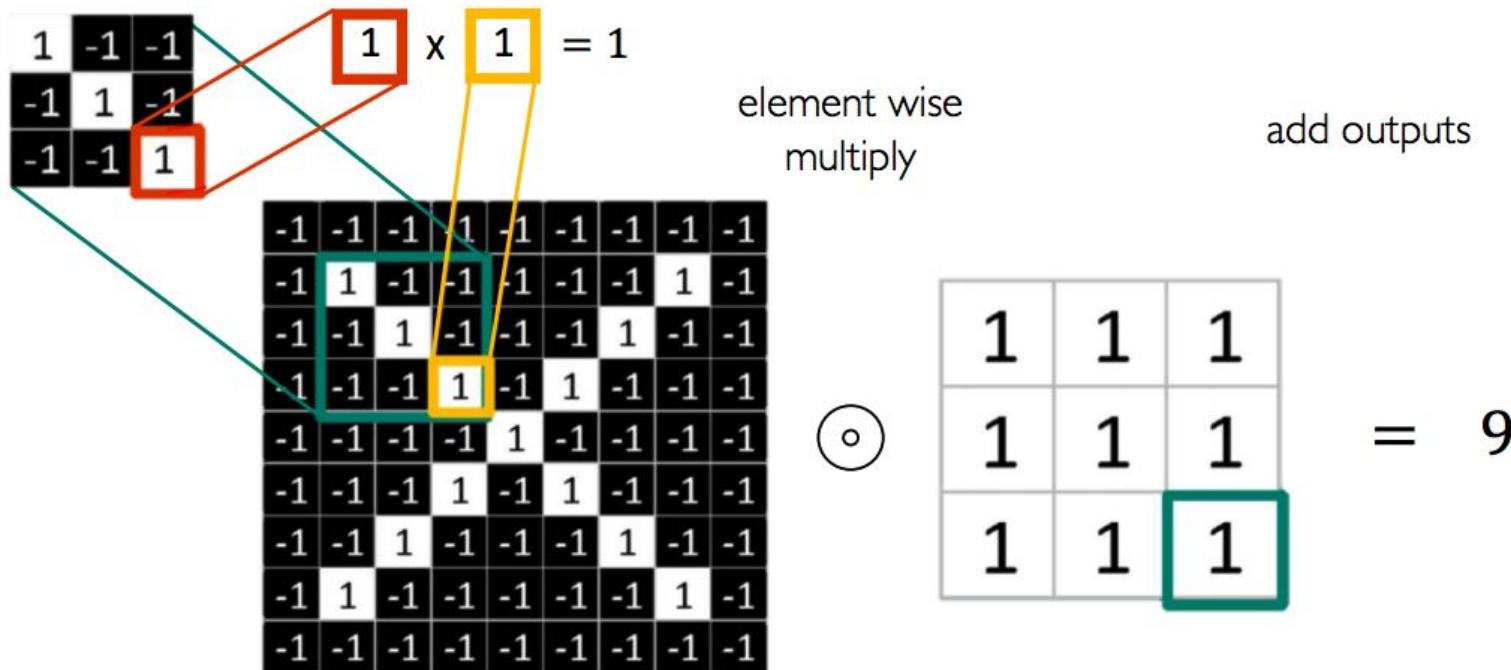
Features of X



Filters to Detect X Features

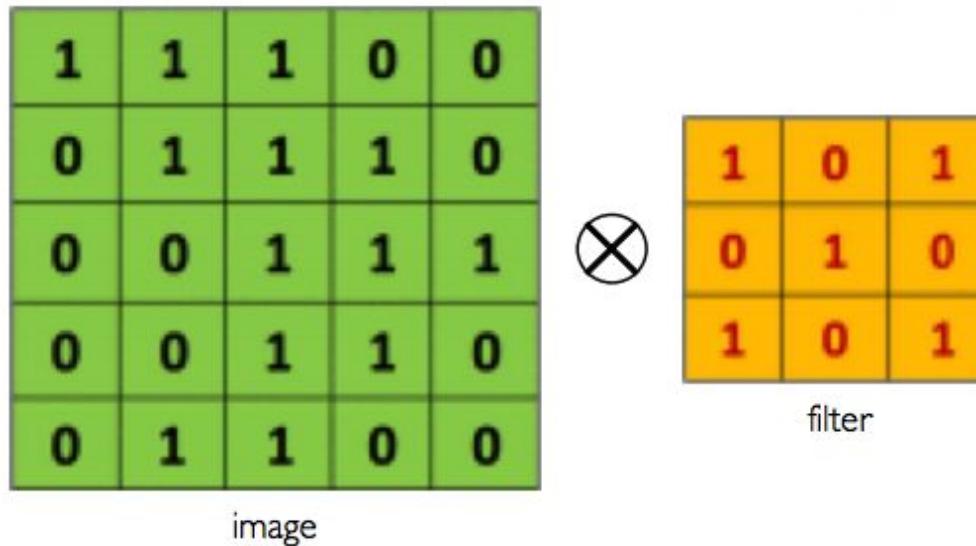


The Convolution Operation



The Convolution Operation

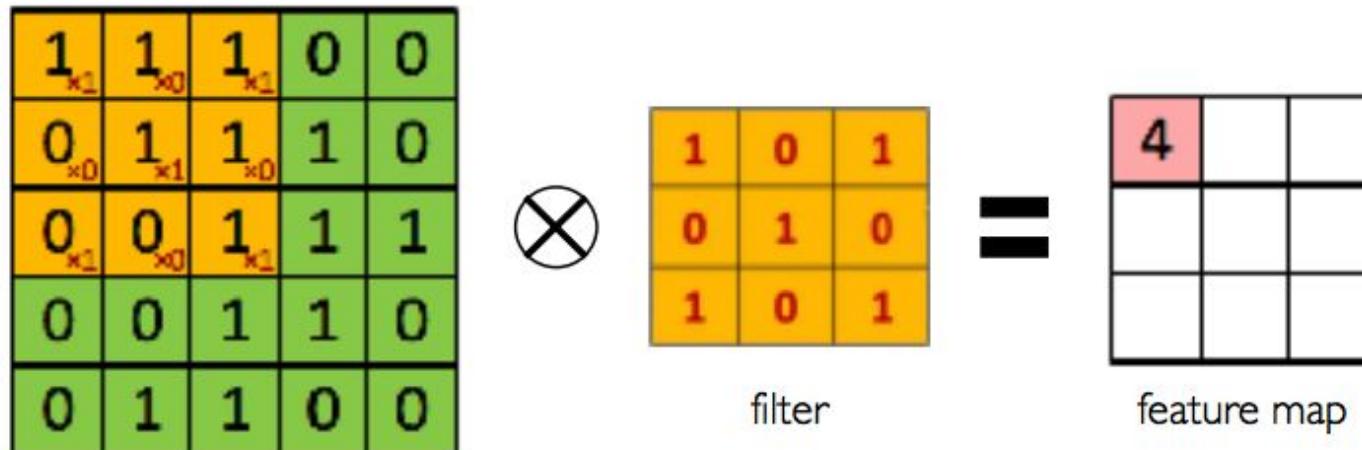
Suppose we want to compute the convolution of a 5x5 image and a 3x3 filter:



We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs...

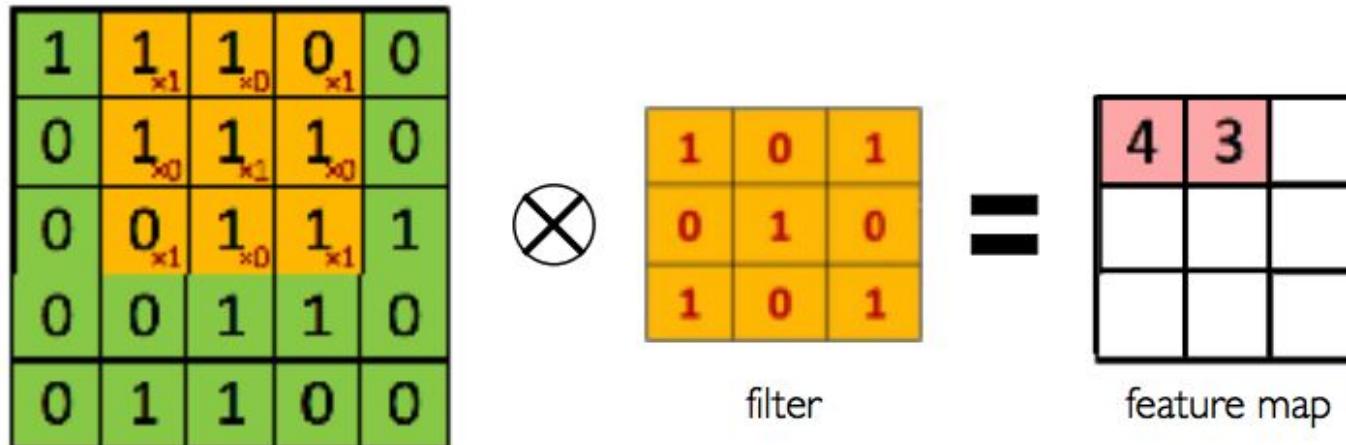
The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



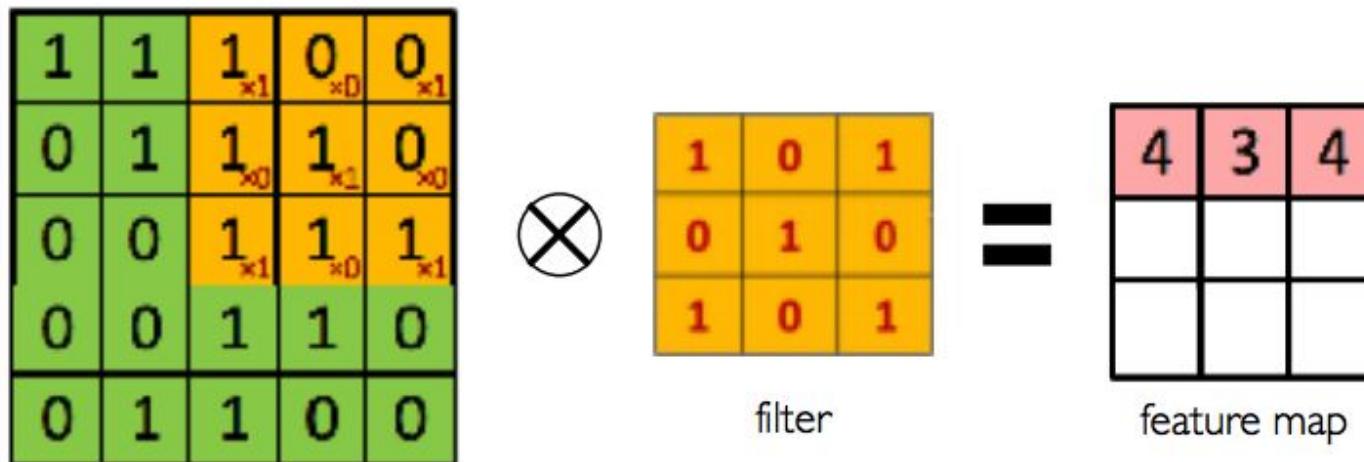
The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



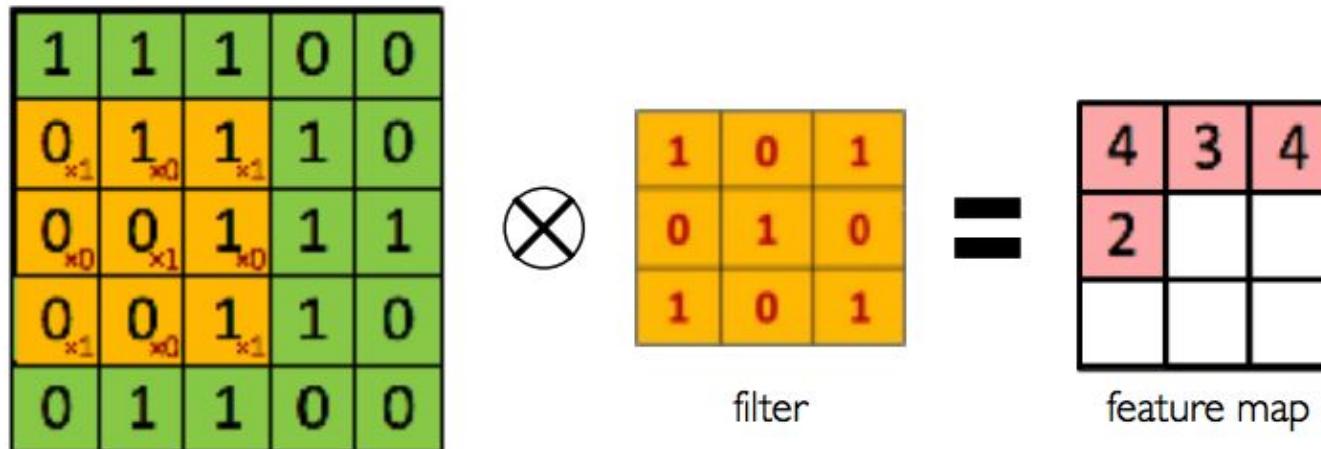
The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



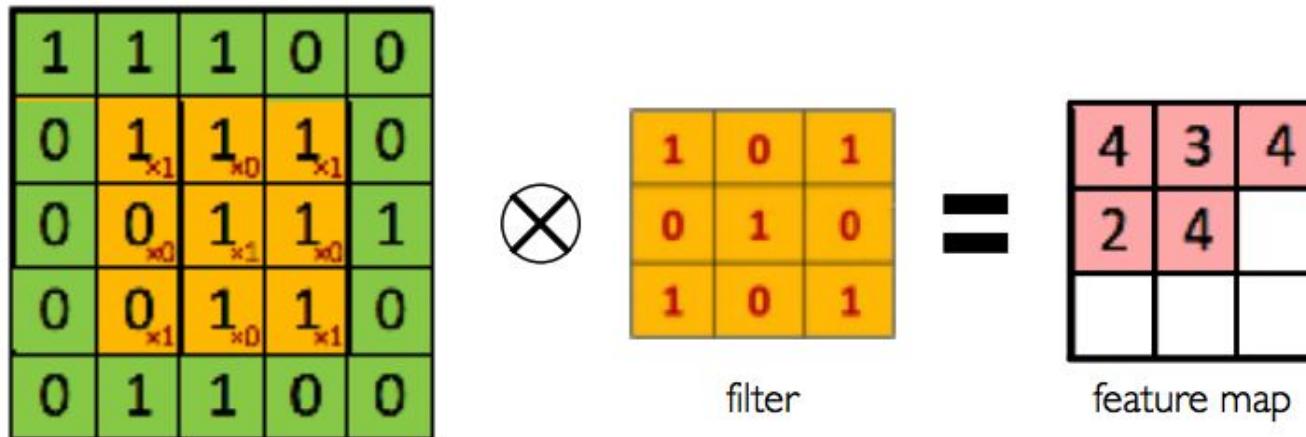
The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:



Quiz

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:

The diagram illustrates the convolution process. On the left is a 5x5 input image with values: Row 1: 1, 1, 1, 0, 0; Row 2: 0, 1, 1, 1, 0; Row 3: 0, 0, 1 (with $\times 1$ below it), 1 (with $\times 0$ below it), 1 (with $\times 1$ below it); Row 4: 0, 0, 1 (with $\times 0$ below it), 1 (with $\times 1$ below it), 0 (with $\times 0$ below it); Row 5: 0, 1, 1 (with $\times 1$ below it), 0 (with $\times 0$ below it), 0 (with $\times 1$ below it). In the center is a 3x3 filter with values: Row 1: 1, 0, 1; Row 2: 0, 1, 0; Row 3: 1, 0, 1. To the right of the filter is an equals sign (=). Below the filter is the label "filter". To the right of the equals sign is a 3x3 feature map with values: Row 1: 4, 3, 4; Row 2: 2, 4, 3; Row 3: 2, 3, 4. Below the feature map is the label "feature map".

1	1	1	0	0
0	1	1	1	0
0	0	1 _{×1}	1 _{×0}	1 _{×1}
0	0	1 _{×0}	1 _{×1}	0 _{×0}
0	1	1 _{×1}	0 _{×0}	0 _{×1}

\otimes

1	0	1
0	1	0
1	0	1

=

4	3	4
2	4	3
2	3	4

filter

feature map

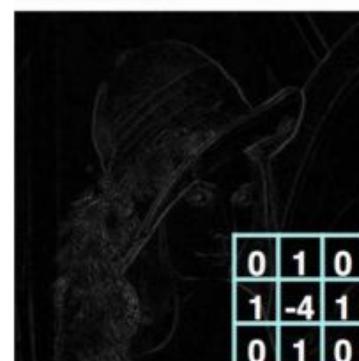
Producing Feature Maps



Original



Sharpen

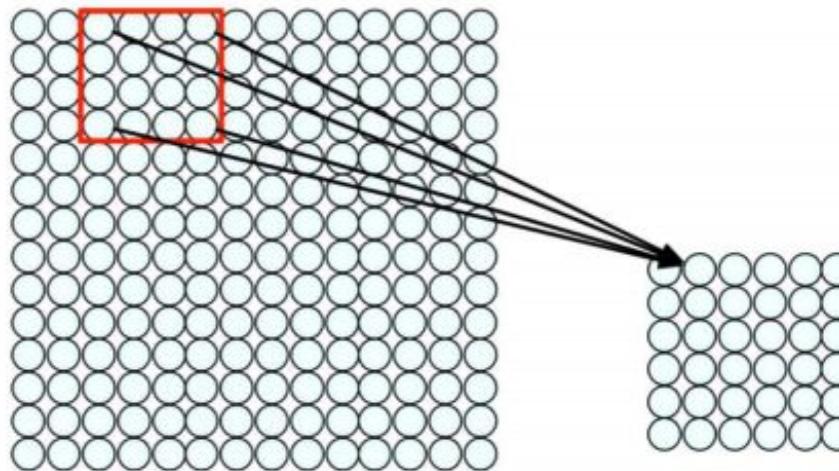


Edge Detect



"Strong" Edge
Detect

Convolution Recap

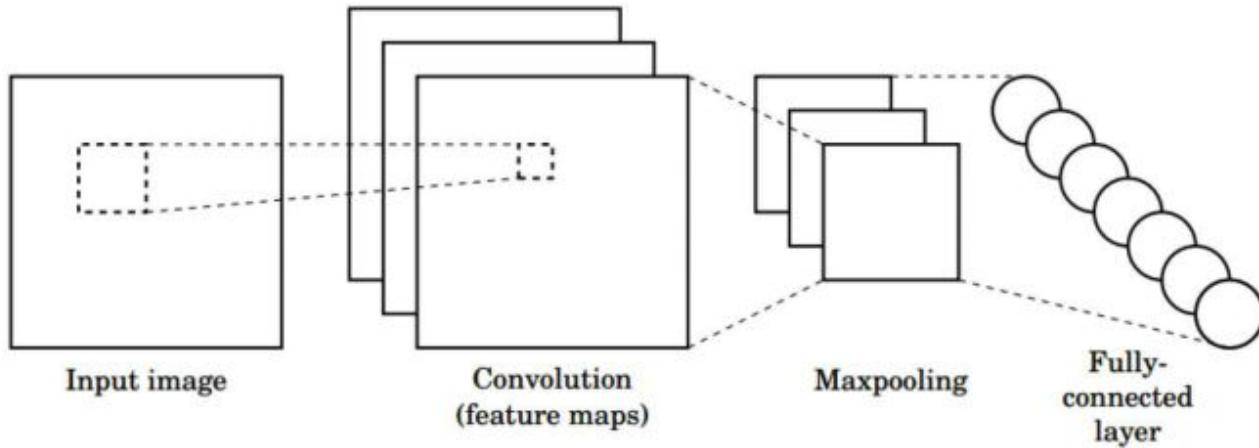


- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Convolutional Neural Networks (CNNs)



CNNs for Classification

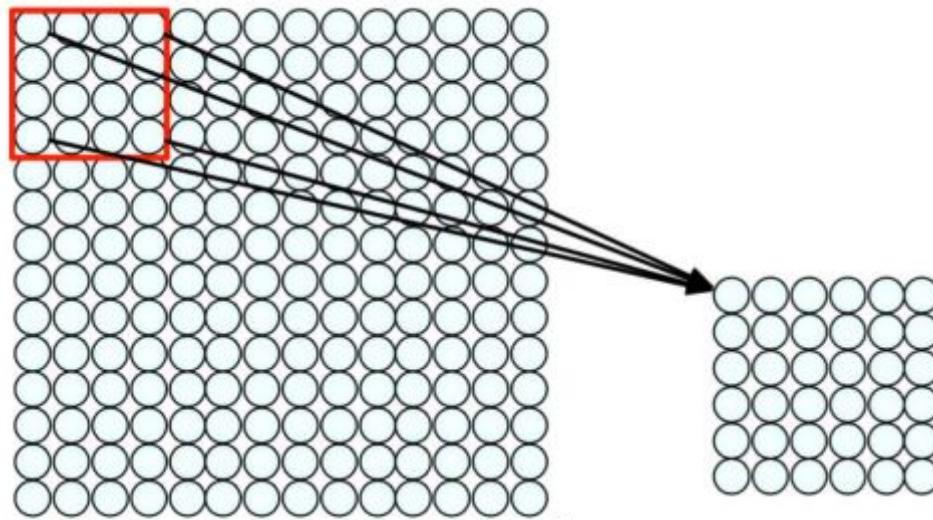


- 1. Convolution:** Apply filters with learned weights to generate feature maps.
- 2. Non-linearity:** Often ReLU.
- 3. Pooling:** Downsampling operation on each feature map.

Train model with image data.

Learn weights of filters in convolutional layers.

Convolutional Layers: Local Connectivity



$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

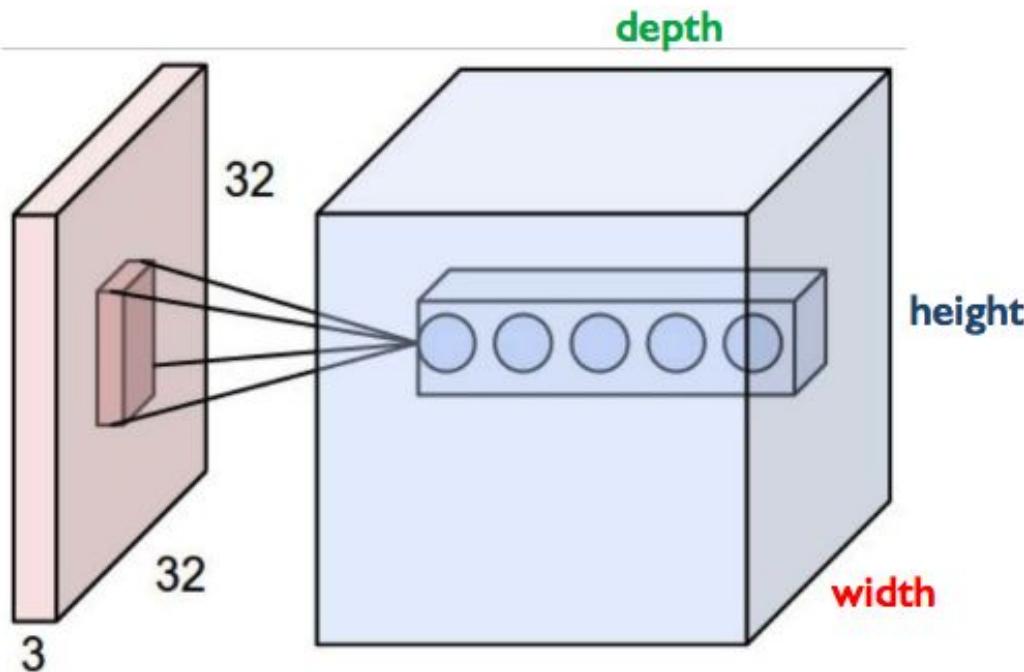
for neuron (p,q) in hidden layer

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

- 1) applying a window of weights
- 2) computing linear combinations
- 3) activating with non-linear function

CNNs: Spatial Arrangement of Output Volume



Layer Dimensions:

$$h \times w \times d$$

where h and w are spatial dimensions
d (depth) = number of filters

Stride:

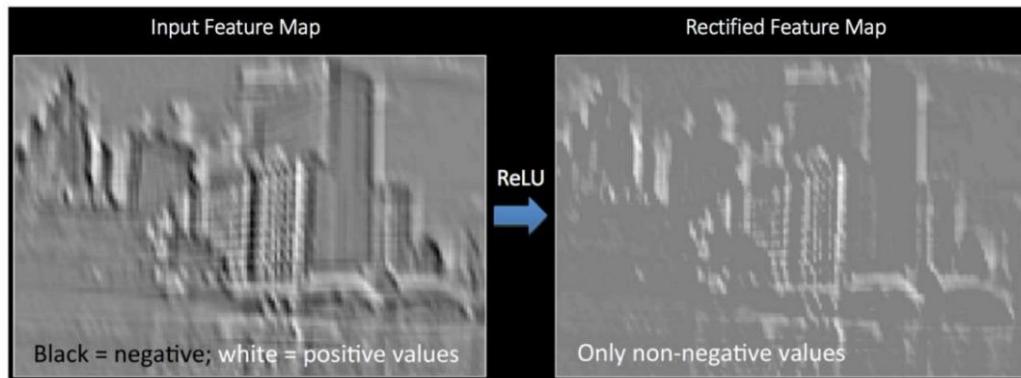
Filter step size

Receptive Field:

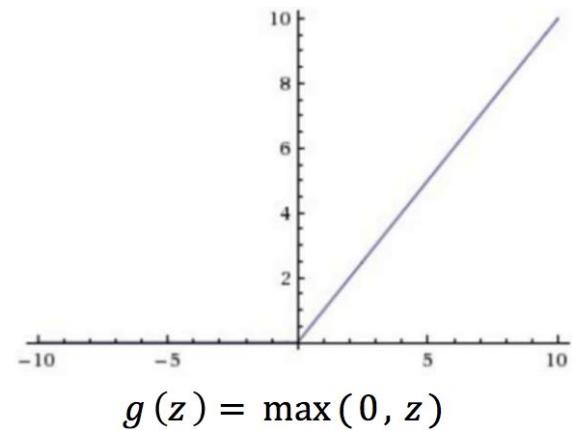
Locations in input image that
a node is path connected to

Introducing Non-Linearity

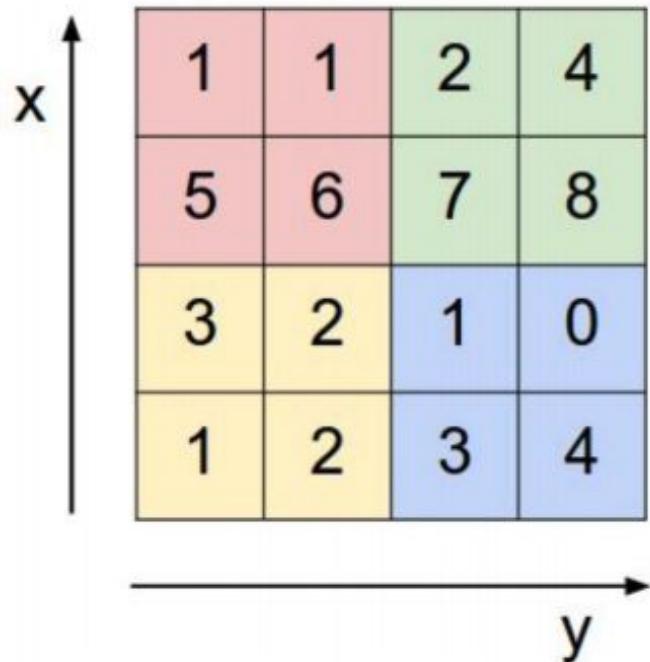
- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**



Rectified Linear Unit (ReLU)



Max Pooling



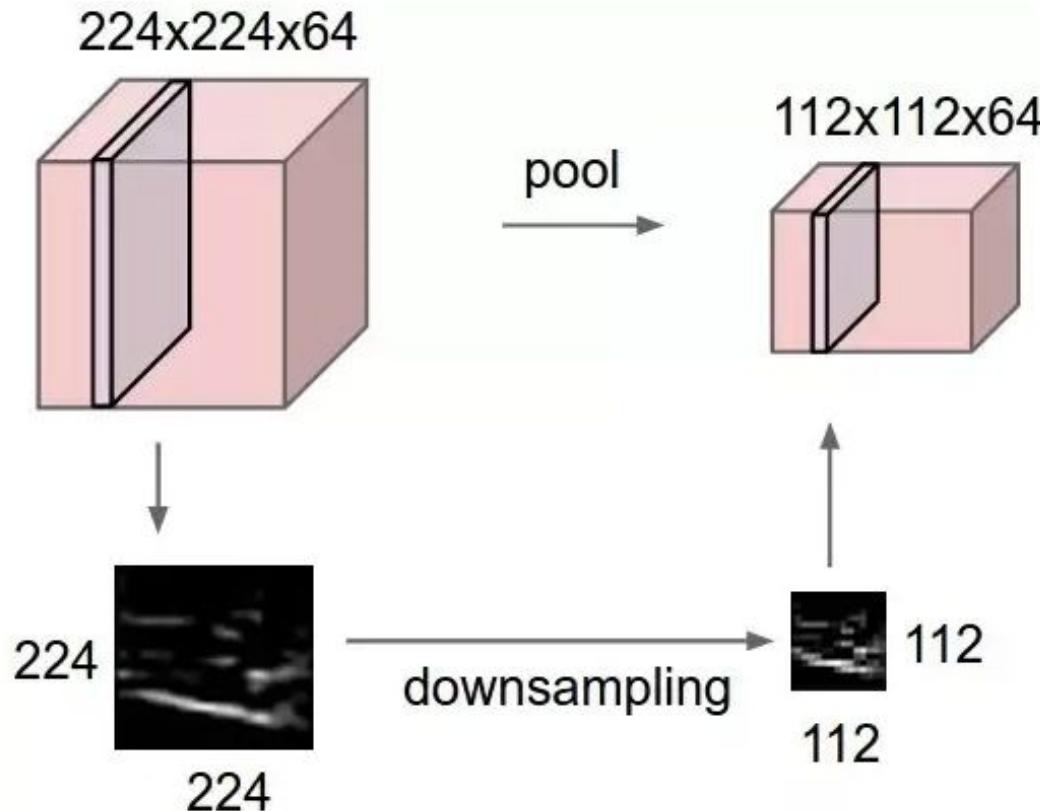
max pool with 2x2 filters
and stride 2



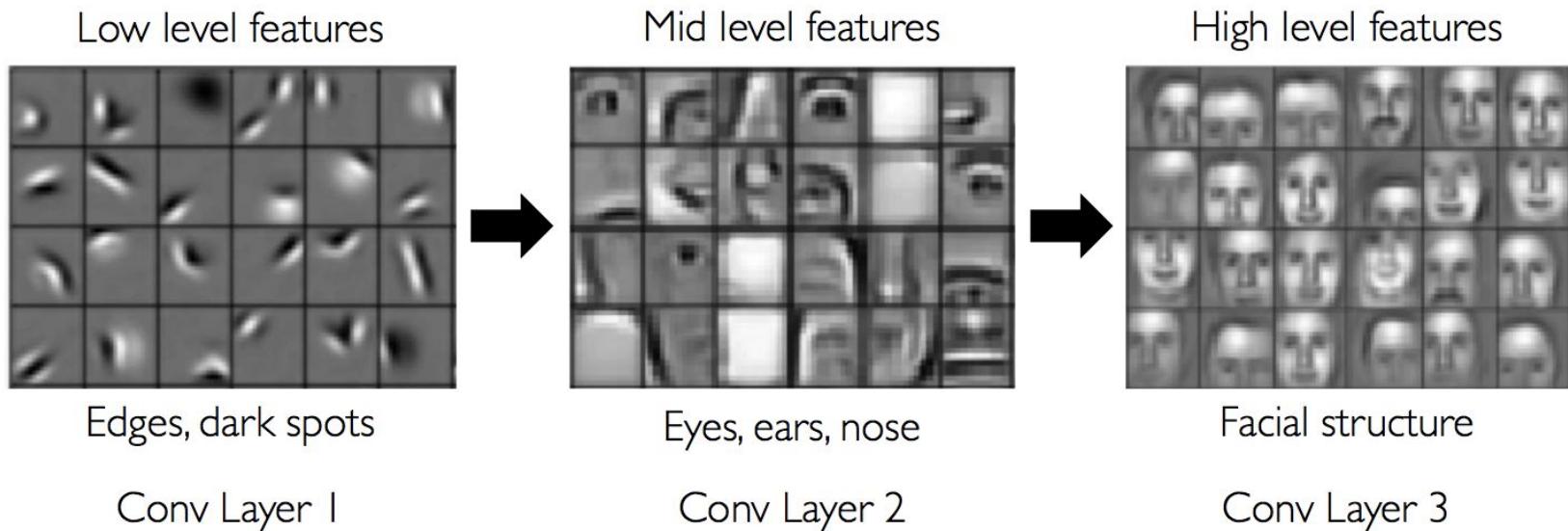
6	8
3	4

- 1) Reduced dimensionality
- 2) Spatial invariance

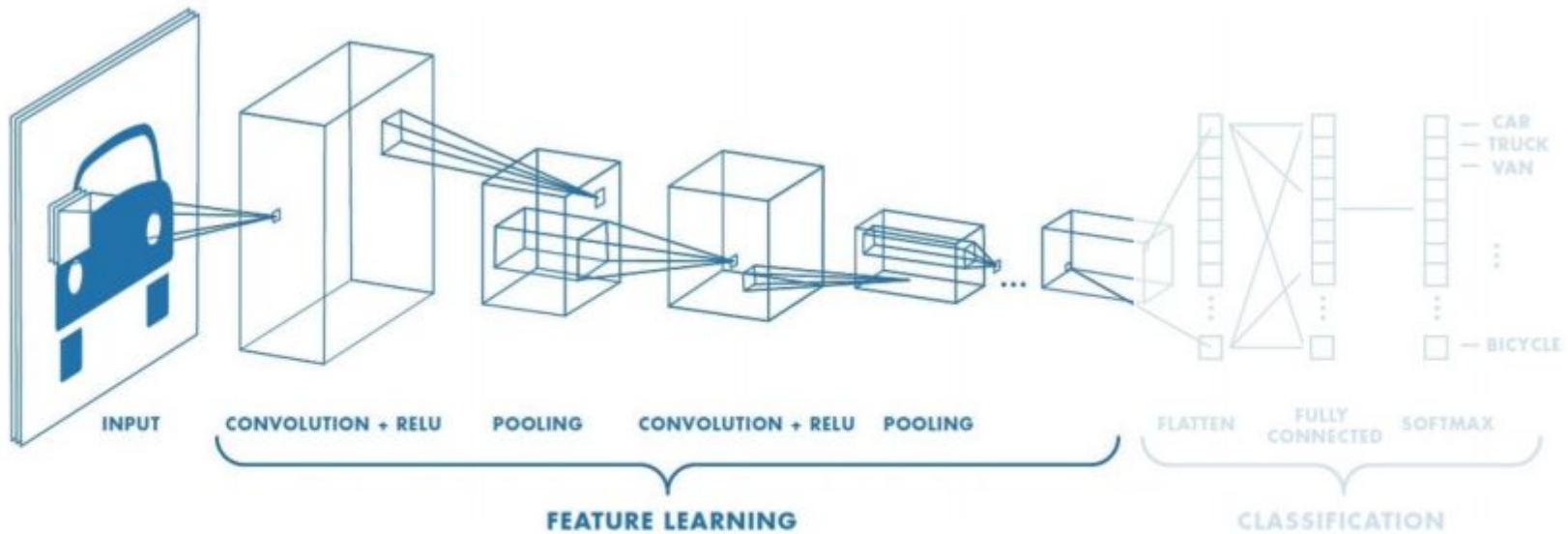
Max Pooling Example



Representation Learning in Deep CNNs

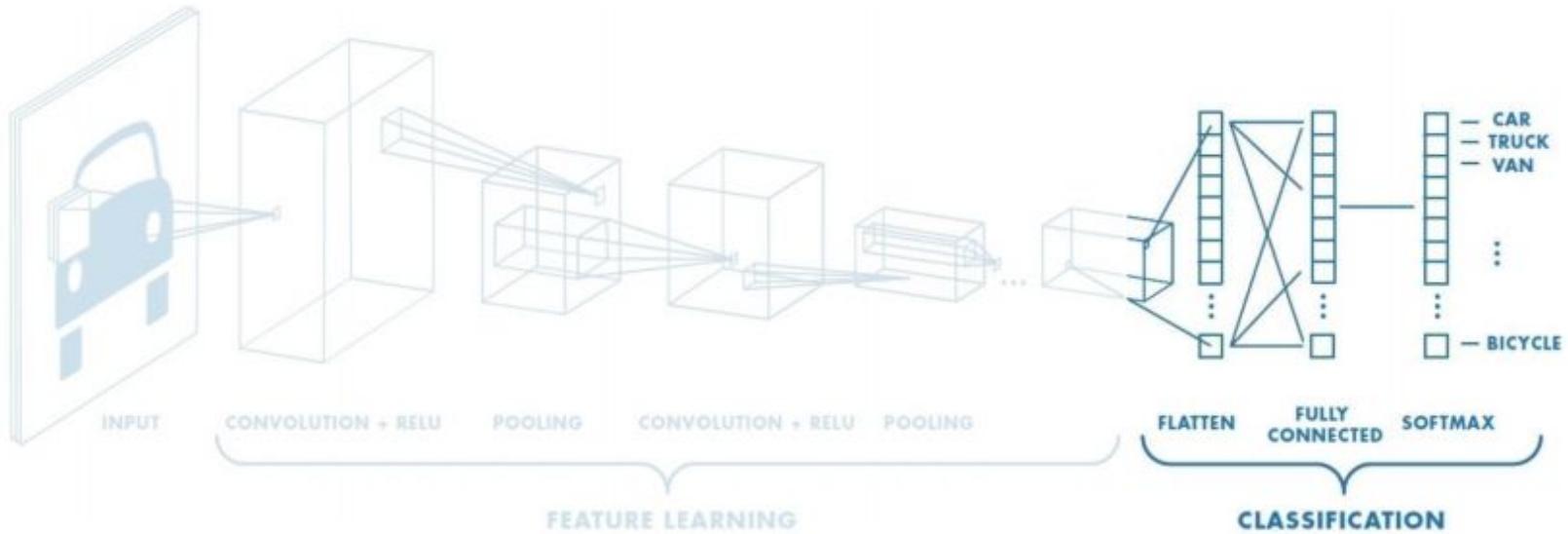


CNNs for Classification: Feature Learning



1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**

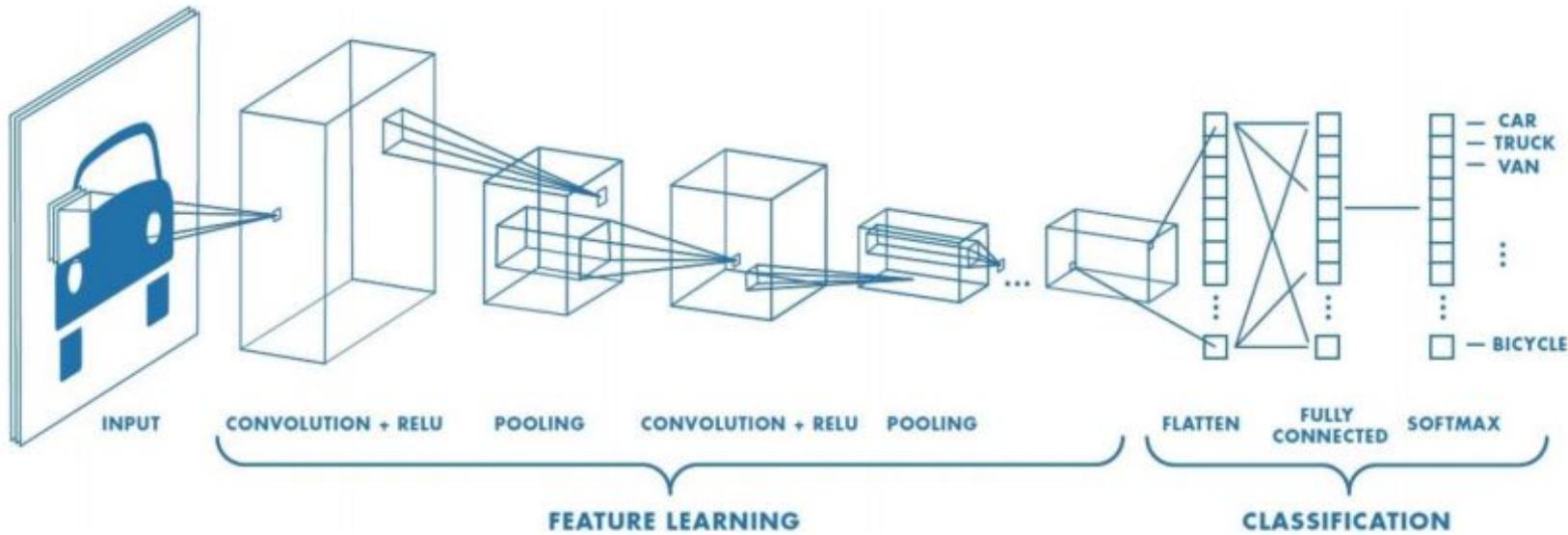
CNNs for Classification: Class Probabilities



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

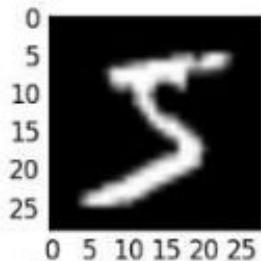
CNNs: Training with Backpropagation



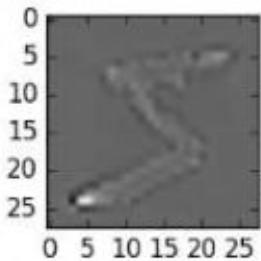
Learn weights for convolutional filters and fully connected layers

Backpropagation: cross-entropy loss

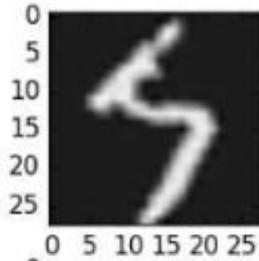
Data Augmentation



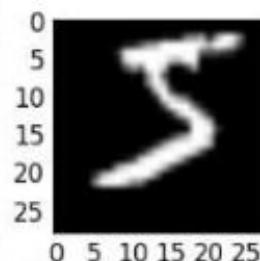
Example MNIST
images



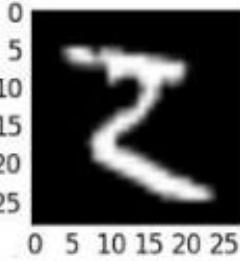
ZCA Whitening



Random Rotations



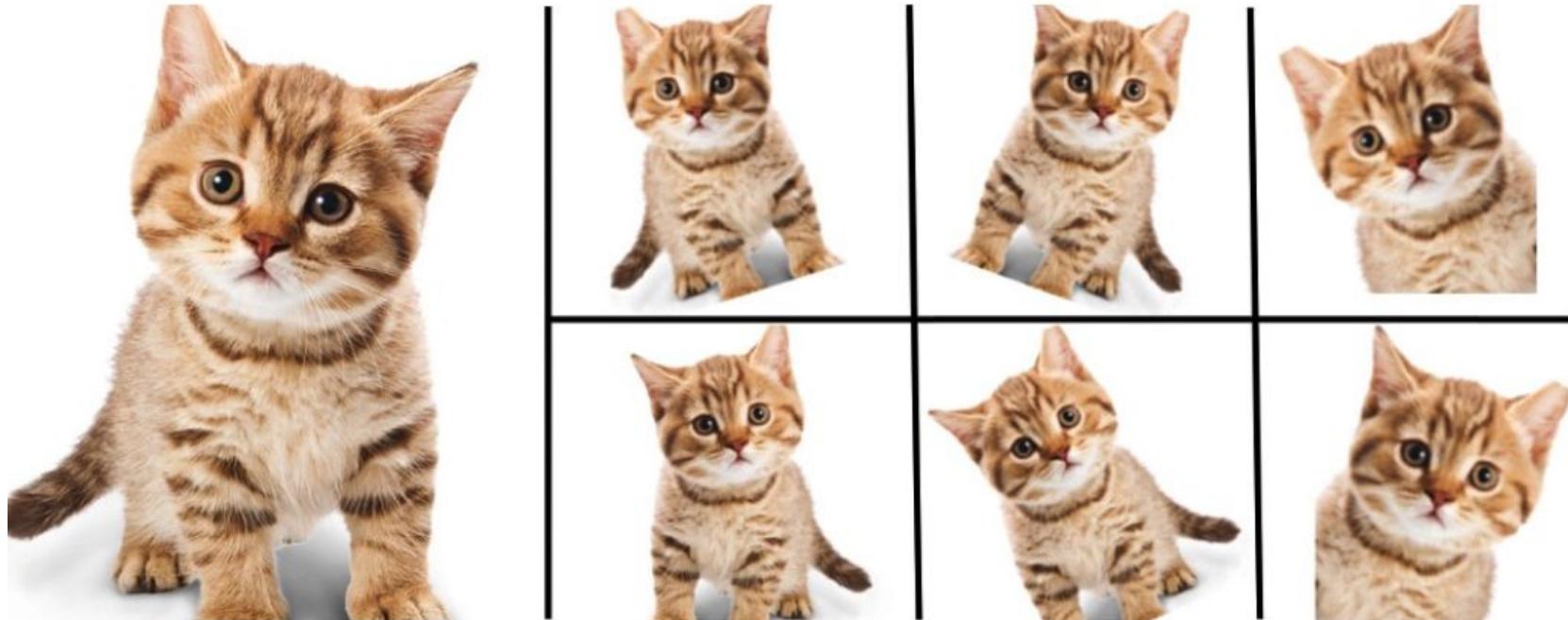
Random shift



Random flip

- Approaches that alter the training data in ways that change the array representation while keeping the label the same.
- Some popular augmentations people use are gray scales , horizontal flips, vertical flips, random crops, color jitters, translations, rotations, and much more.

Data Augmentation

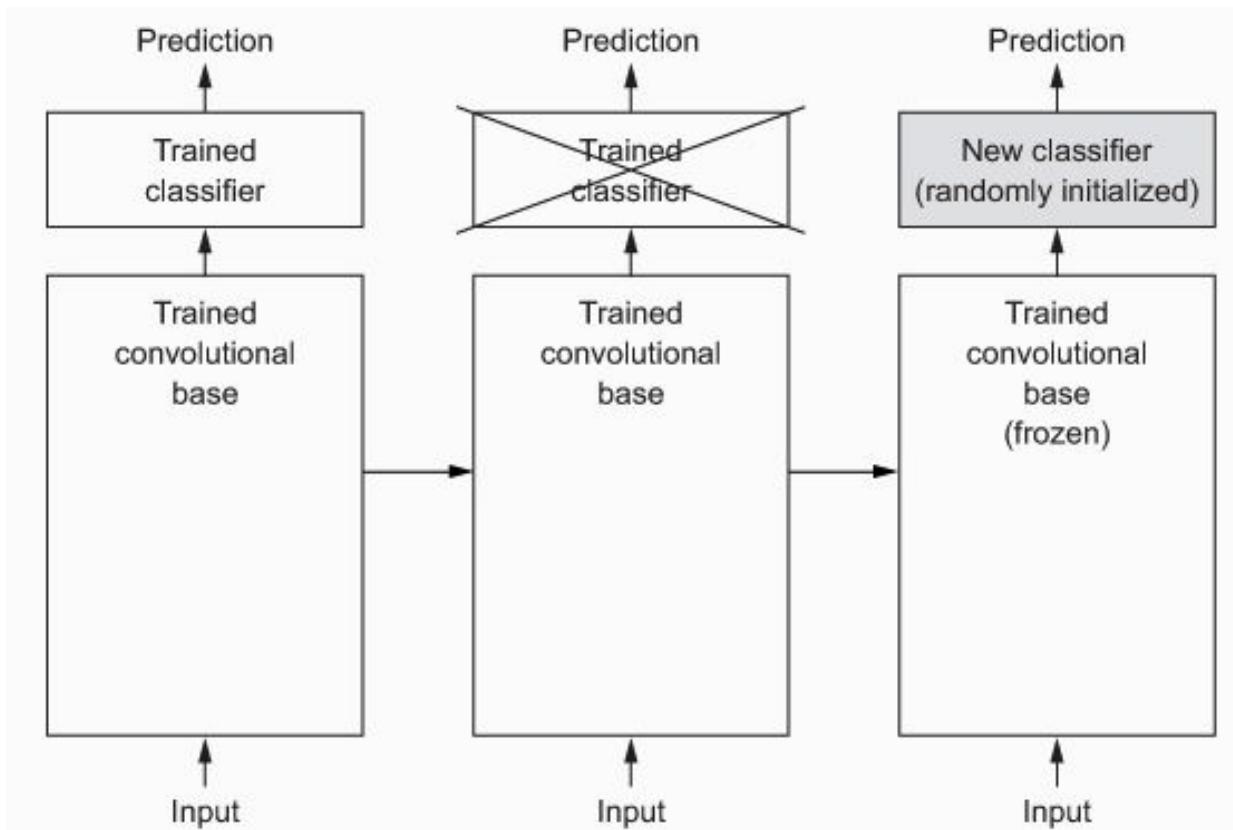


Transfer Learning

Transfer learning is the process of taking a pre-trained model (the weights and parameters of a network that has been trained on a large dataset by somebody else) and “fine-tuning” the model with your own dataset.

- Pre-trained model will act as a feature extractor.
- Freeze the weights of all the other layers .
- Remove the last layer of the network and replace it with your own classifier.
- Train the network normally.

Transfer Learning



CNNs on ImageNet



ImageNet Dataset

Dataset of over 14 million images across 21,841 categories

"Elongated crescent-shaped yellow fruit with soft sweet flesh"



1409 pictures of bananas.

ImageNet Challenge



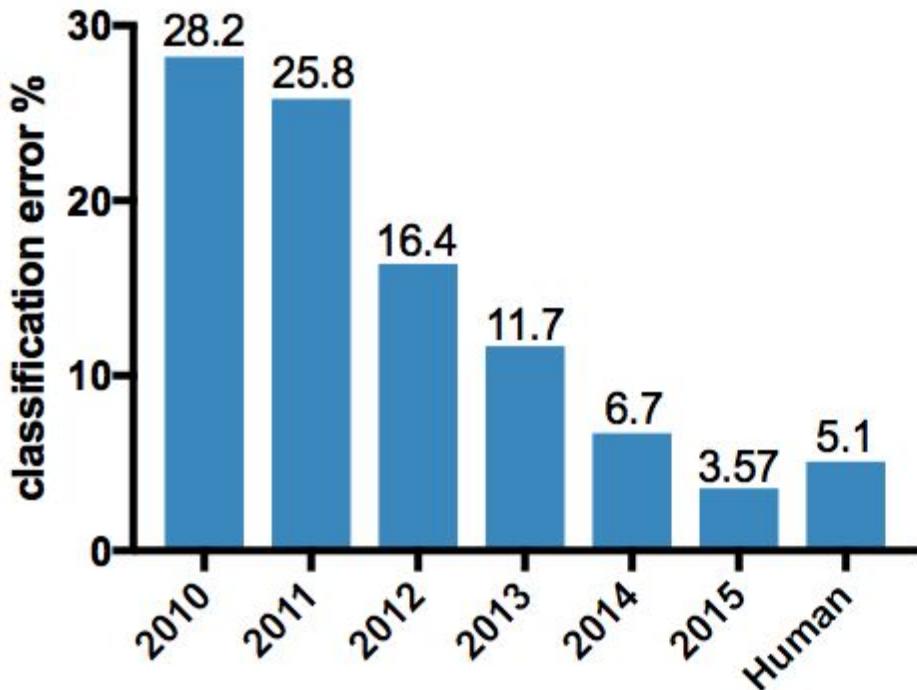
ImageNet Large Scale Visual Recognition Challenges



Classification task: produce a list of object categories present in image. 1000 categories.

“Top 5 error”: rate at which the model does not output correct label in top 5 predictions

ImageNet Challenge: Classification Task



2012: AlexNet. First CNN to win.

- 8 layers, 61 million parameters

2013: ZFNet

- 8 layers, more filters

2014: VGG

- 19 layers

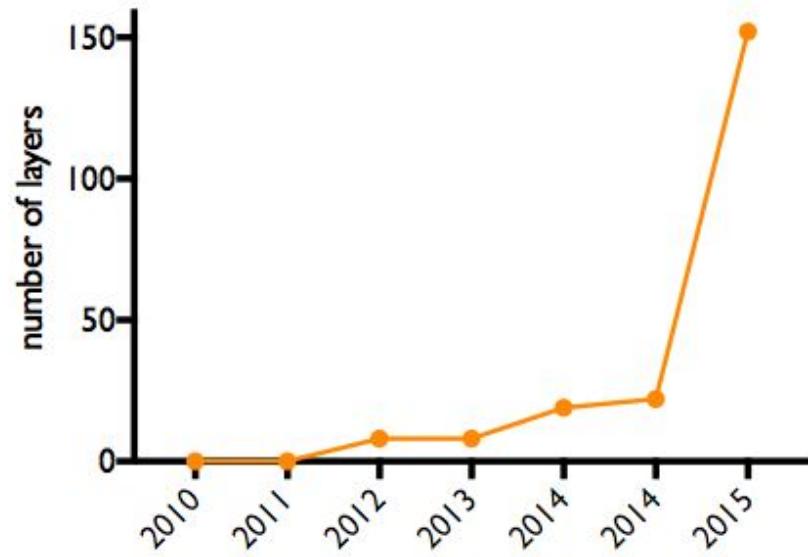
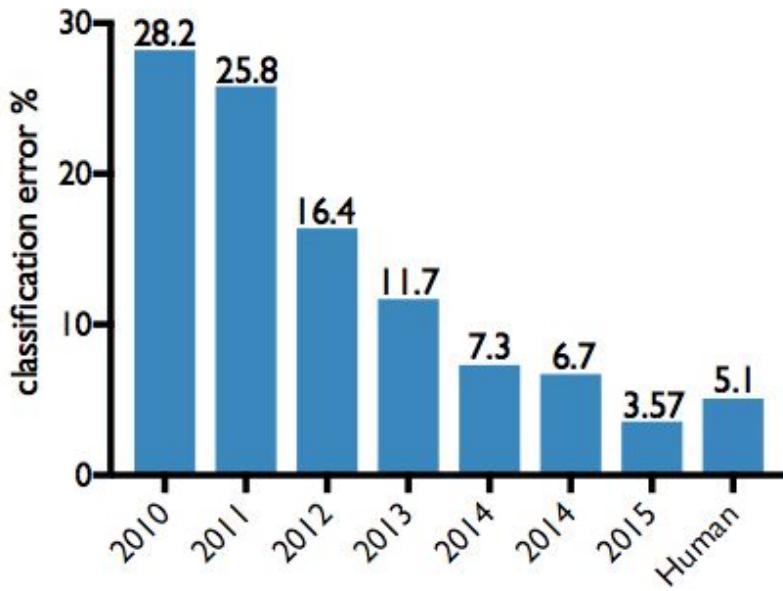
2014: GoogLeNet

- "Inception" modules
- 22 layers, 5million parameters

2015: ResNet

- 152 layers

ImageNet Challenge: Classification Task



Lab

