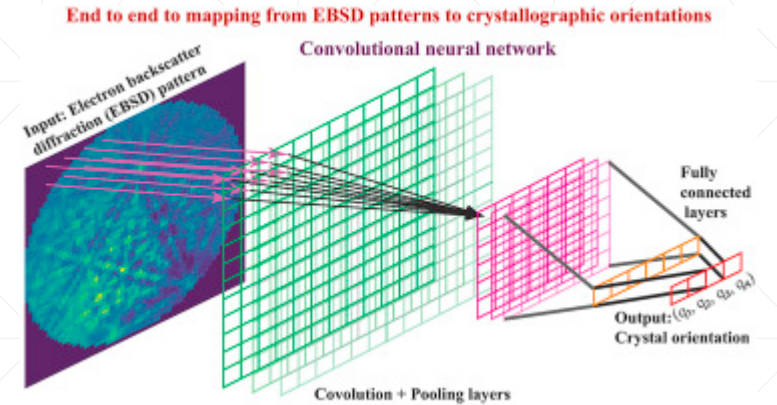


Convolutional Neural Network

~Abhishek Kumar

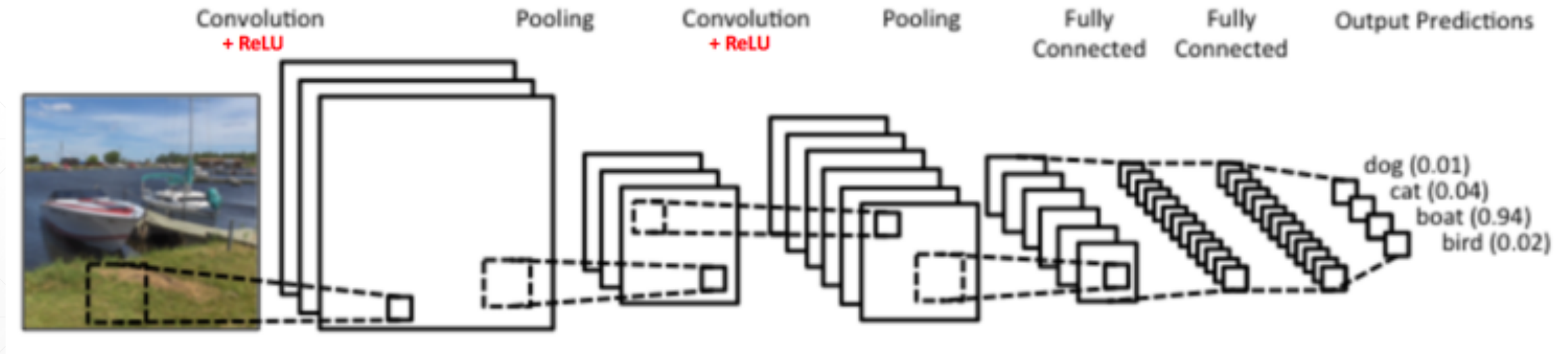


Convolutional Neural Network (CNN)

- Image recognition
 - Classification
 - Natural Language Processing tasks
 - Multi Layer Perceptron's = Fully Connected Layers
-

The LeNet Architecture (1990s)

Character recognition tasks such as reading zip codes, digits,



Operations/ Basic building blocks

- Convolution
 - Non Linearity (ReLU)
 - Pooling or Sub Sampling
 - Classification (Fully Connected Layer)
-

1. Convolution

- Extract features from the input image
 - Convolution preserves the spatial relationship between pixels by learning image features using small squares of input data.
-

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

1	0	1
0	1	0
1	0	1

Filter/
Kernel/
Feature detector



1	1	1	0	0
0	1	1	1	0
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1 _{x0}	1 _{x1}	0 _{x0}
0	1	1 _{x1}	0 _{x0}	0 _{x1}

Image

4	3	4
2	4	3
2	3	4

Convolved
Feature / Activation Map
/ Feature Map

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved
Feature

1	1	1 _{x1}	0 _{x0}	0 _{x1}
0	1	1 _{x0}	1 _{x1}	0 _{x0}
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1	1	0
0	1	1	0	0

Image

4	3	4

Convolved
Feature

1	1	1	0	0
0 _{x1}	1 _{x0}	1 _{x1}	1	0
0 _{x0}	0 _{x1}	1 _{x0}	1	1
0 _{x1}	0 _{x0}	1 _{x1}	1	0
0	1	1	0	0

Image

4	3	4
2		

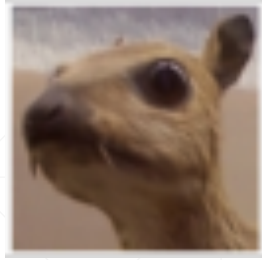
Convolved
Feature

1	1	1	0	0
0	1 _{x1}	1 _{x0}	1 _{x1}	0
0	0 _{x0}	1 _{x1}	1 _{x0}	1
0	0 _{x1}	1 _{x0}	1 _{x1}	0
0	1	1	0	0

Image

4	3	4
2	4	

Convolved
Feature



- Different values of the filter matrix will produce different Feature Maps for the same input image.
- Edge Detection, Sharpen and Blur
- This means that different filters can detect different features from an image, for example edges, curves etc

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

- CNN *learns* the values of these filters on its own during the training process
 - Specify parameters such as number of filters, filter size, architecture of the network
 - More number of filters = More image features get extracted = better our network at recognizing patterns in unseen images.
-

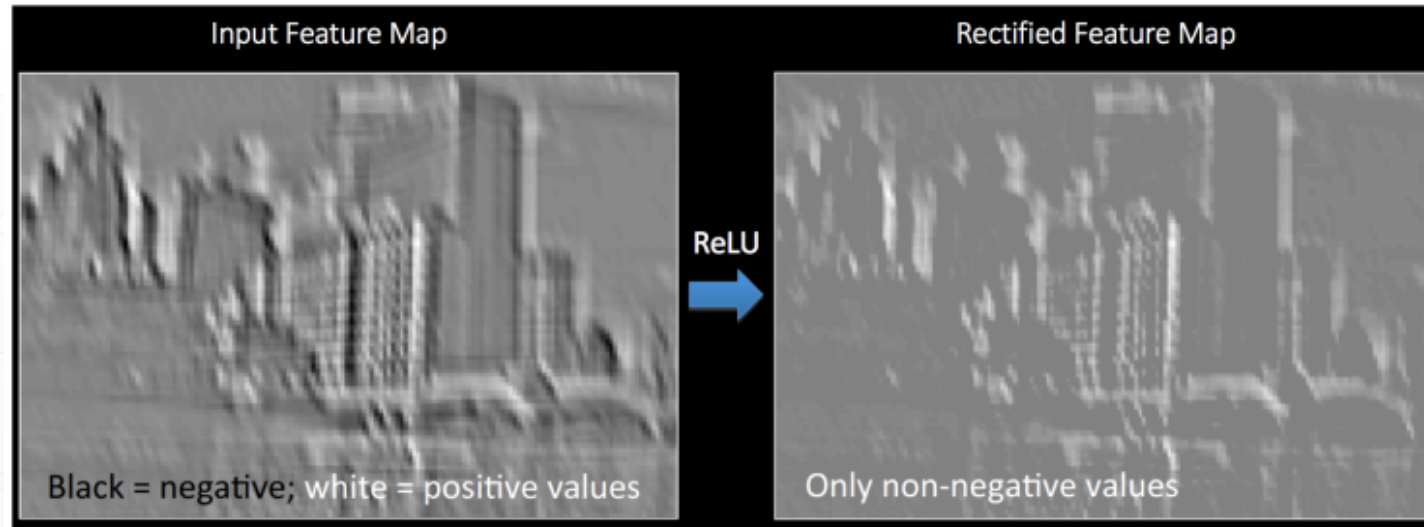
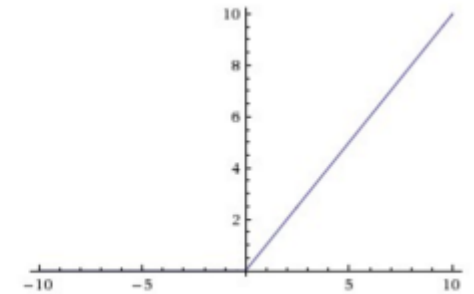
Size of the Feature Map

- **Depth:** number of filters we use for the convolution operation.
 - **Stride:** number of pixels by which we slide our filter matrix over the input matrix.
 - **Zero-padding:** pad the input matrix with zeros around the border, so that we can apply the filter to bordering elements of our input image matrix.
 - Wide Convolution
 - Narrow Convolution
-

2. ReLU (Rectified Linear Unit)

- Introducing Non Linearity
- Perform better in most situations

$$\text{Output} = \text{Max}(\text{zero}, \text{Input})$$

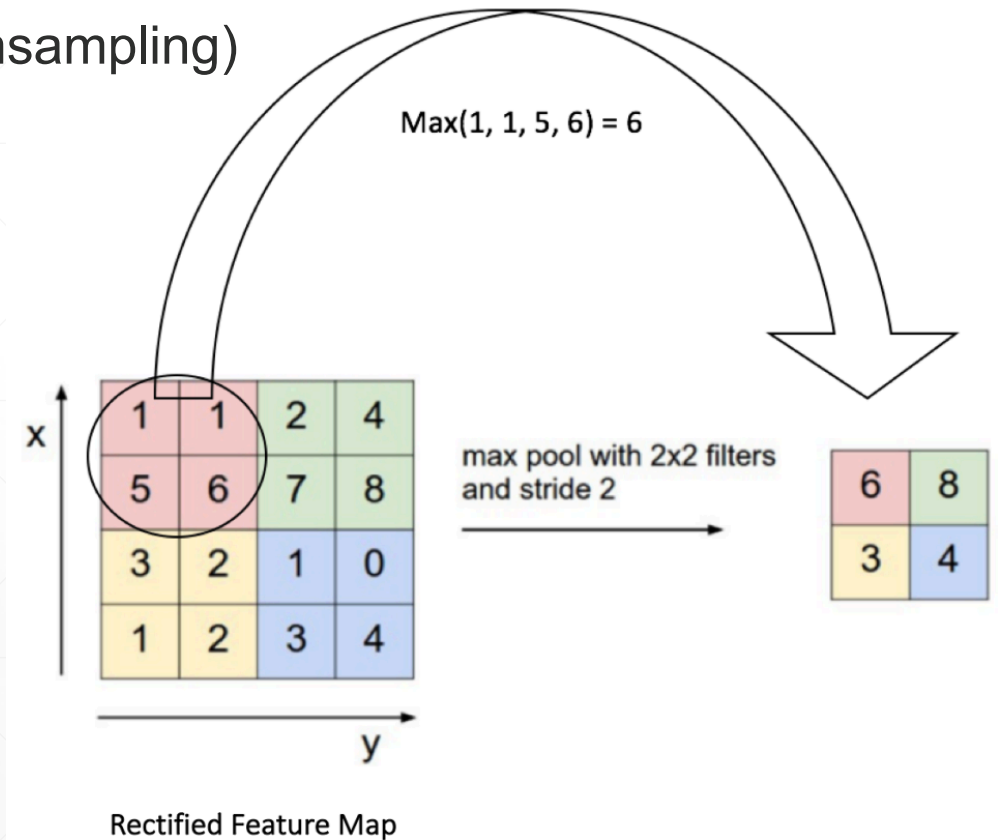


3. Pooling

- Spatial Pooling (also called subsampling or downsampling)

- Reduces the dimensionality of each feature map but retains the most important information. (Like PCA)

- Max
- Sum
- Avg



Convolution
using 3 filters
+ ReLU

Pooling applied
separately on each
feature map





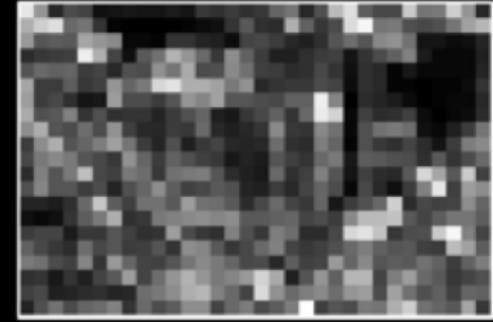
Only non-negative values

Rectified Feature Map

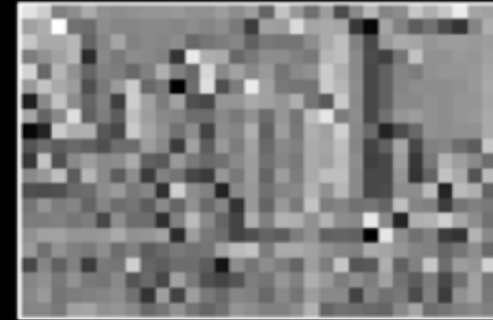
Pooling



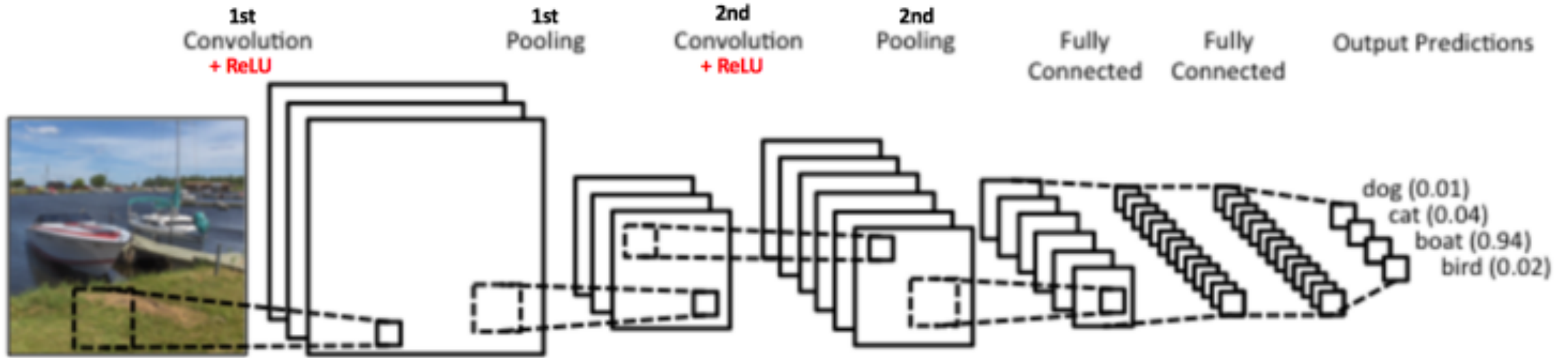
Max



Sum



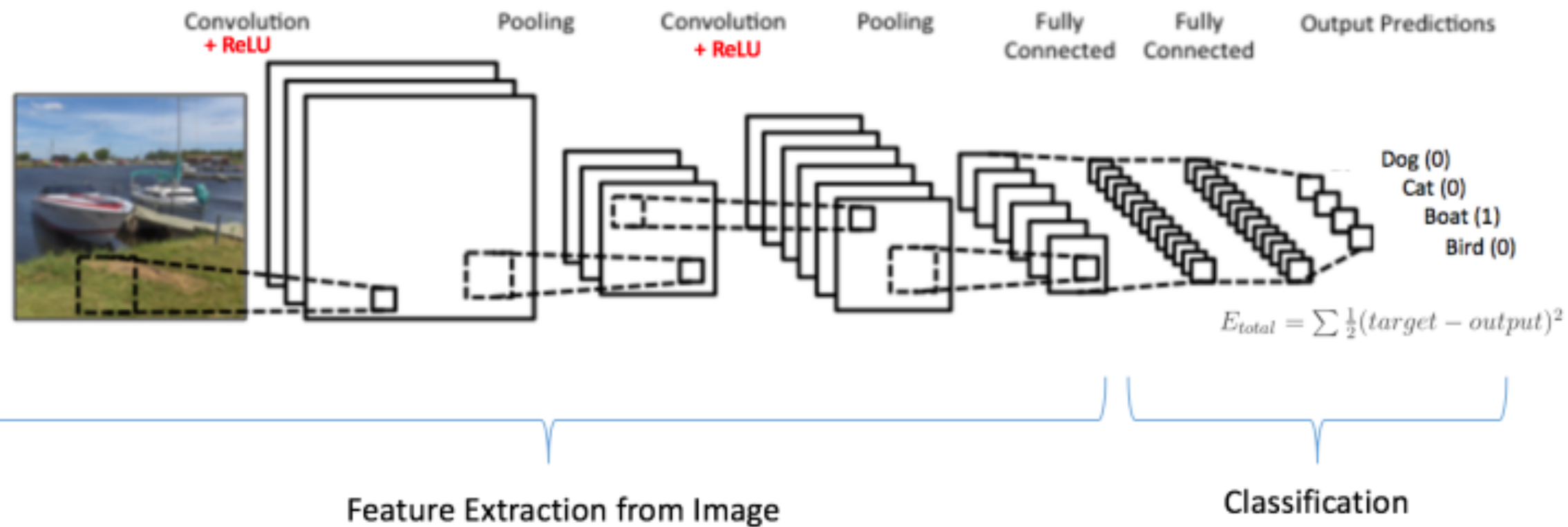
- Makes the input representations (feature dimension) smaller and more manageable
 - reduces the number of parameters and computations in the network, therefore, controlling overfitting.
 - Makes the network invariant to small transformations, distortions and translations in the input image. (since we take the maximum / average value in a local neighborhood).
-



Together these layers extract the useful features from the images, introduce non-linearity in our network and reduce feature dimension while aiming to make the features somewhat equivariant to scale and translation

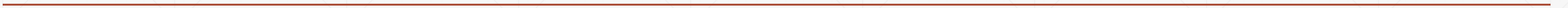
4. Fully Connected Layer

- Output of the 2nd Pooling Layer = Input to the Fully Connected Layer
 - Multi Layer Perceptron that uses a activation function in the output layer (eg. Softmax)
 - The sum of output probabilities from the Fully Connected Layer is 1
-



Step1:

- We initialize all filters and parameters / weights with random values



Step2:

- The network takes a training image as input, goes through the forward propagation step (convolution, ReLU and pooling operations along with forward propagation in the Fully Connected layer) and finds the output probabilities for each class.
 - Lets say the output probabilities for the boat image above are [0.2, 0.4, 0.1, 0.3]
 - Since weights are randomly assigned for the first training example, output probabilities are also random
-

Step3:

- Calculate the total error at the output layer (summation over all 4 classes)
 - **Total Error = $\sum \frac{1}{2} (\text{target probability} - \text{output probability})^2$**
-

Step4:

- Use Backpropagation to calculate the *gradients* of the error with respect to all weights in the network and use *gradient descent* to update all filter values / weights and parameter values to minimize the output error.
 - Parameters like number of filters, filter sizes, architecture of the network etc. have all been fixed before Step 1 and do not change during training process – only the values of the filter matrix and connection weights get updated.
-

Step5:

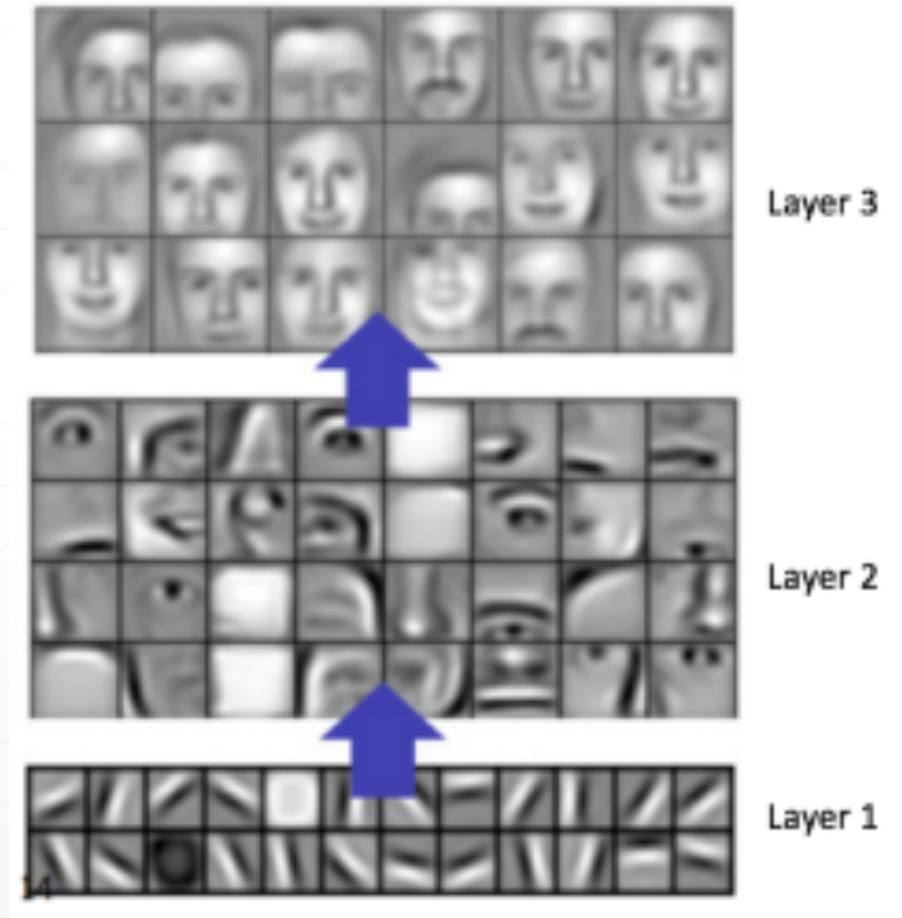
- Repeat steps 2-4 with all images in the training set.



- In general, the more convolution steps we have, the more complicated features our network will be able to learn to recognize.

For example, in Image Classification

- detect edges from raw pixels in the first layer,
- use the edges to detect simple shapes in the second layer,
- and then use these shapes to detect higher-level features, such as facial shapes in higher layers



Adam Harley's example on MNIST database

- <http://scs.ryerson.ca/~aharley/vis/conv/flat.html>



Python code implementation

- Jupyter Notebook





Discussion



Thank you!