Lecture Notes for

# Machine Learning in Python

[ 👨‍🏫 , 👨‍💻 , 🐍 , 👨‍🔬 ]

## Professor Eric Larson

# Convolutional Neural Networks

# Logistics and Agenda

- Logistics
  - Wide/Deep due soon!
  - late turn in… reminder
- Agenda
  - Basic CNN architectures and Gradient

# Class Overview, by topic

**Table Data Visualization**

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

**Dimension Reduction and Image Processing**

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

**Linear and Logistic Regression**

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

**Neural Networks and Back Prop.**

Numpy
Detailed mathematics for NN operations

**Wide and Deep Networks**

**Convolutional Networks**

**Recurrent Networks**

Keras, Tensorflow
Intuition, Detailed implement.

**Ethics in Language Models**

ConceptNet
Case studies

# Convolutional Neural Networks

$$\sum \left( \mathbf{I} \left[ i \pm \frac{r}{2}, j \pm \frac{c}{2} \right] \odot \mathbf{f} \right) = \mathbf{O}[i,j]$$

output image
at pixel i,j

input image at $r$ x $c$ range of
pixels centered in $i,j$

kernel of size, $r$ x $c$
usually $r=c$

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 12 | 9 | 8 | 0 |
| 0 | 5 | 2 | 3 | 4 | 12 | 9 | 8 | 0 |
| 0 | 5 | 2 | 1 | 4 | 10 | 9 | 8 | 0 |
| 0 | 7 | 2 | 1 | 4 | 12 | 7 | 8 | 0 |
| 0 | 7 | 2 | 1 | 4 | 14 | 9 | 8 | 0 |
| 0 | 5 | 2 | 3 | 4 | 12 | 7 | 8 | 0 |
| 0 | 5 | 2 | 1 | 4 | 12 | 9 | 8 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

input image, $\mathbf{I}$

| 1 | 2 | 1 |
|---|---|---|
| 2 | 4 | 2 |
| 1 | 2 | 1 |

kernel
filter, $\mathbf{f}$
3x3

| 20 | 21 | 36 | … | … | … | … |
|----|----|----|---|---|---|---|
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |
| … | … | … | … | … | … | … |

output image, $\mathbf{O}$

# Breaking Apart Convolution Operations



$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

Lecture Notes for Machine Learning in Python    |    Professor Eric C. Larson

Convolutional Filters + Statistics of Filters (Edge Magnitude)

Convolutional Filters (Orientation)

Convolutional Filters (Local Gaussian)

Statistics of Filter Outputs (Histograms)

take normalized histogram at point $u,v$
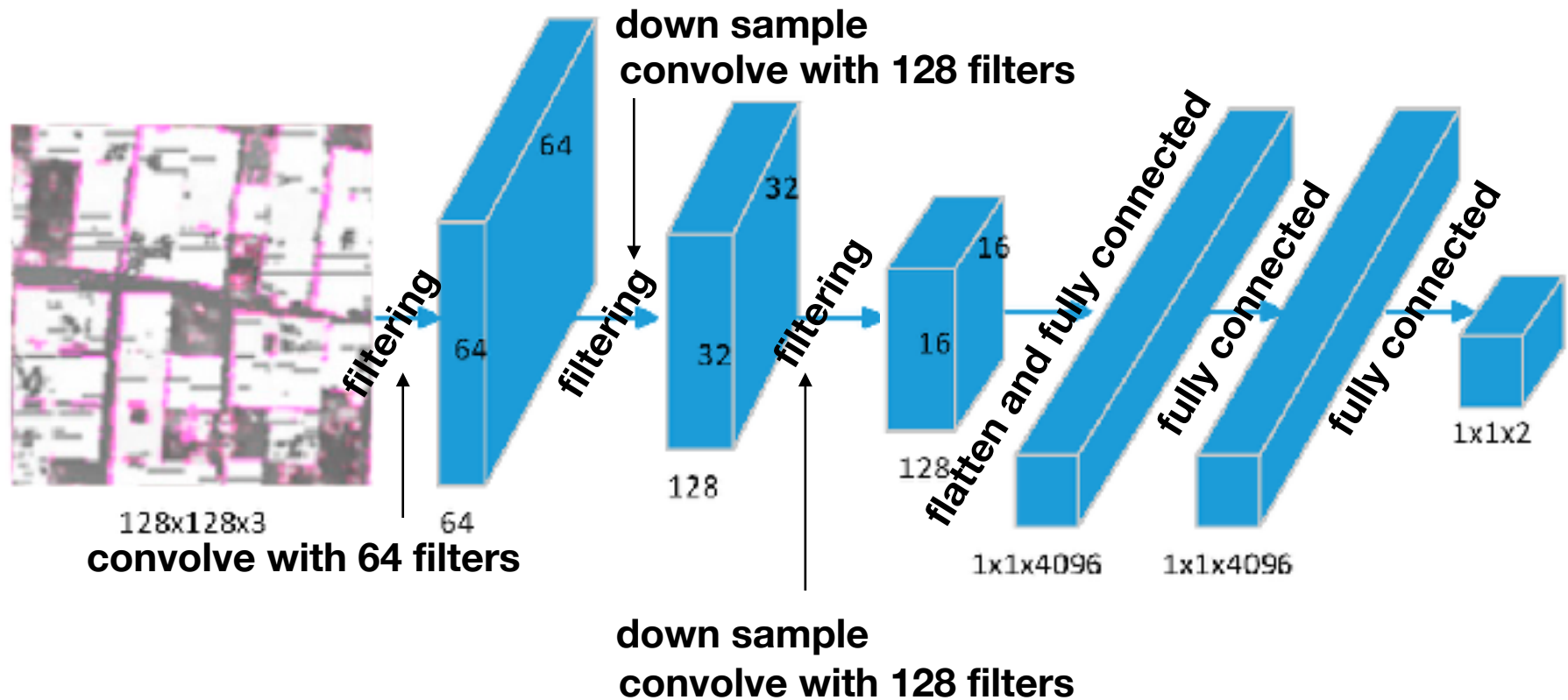
$$\widetilde{\mathbf{h}}_\Sigma(u, v) = \left\| [\mathbf{G}_1^\Sigma(u, v), \ldots, \mathbf{G}_H^\Sigma(u, v)]^\top \right\|$$

$$\mathcal{D}(u_0, v_0) =$$

$$\left[ \widetilde{\mathbf{h}}_{\Sigma_1}^\top(u_0, v_0), \right.$$

$$\widetilde{\mathbf{h}}_{\Sigma_1}^\top(\mathbf{l}_1(u_0, v_0, R_1)), \cdots, \widetilde{\mathbf{h}}_{\Sigma_1}^\top(\mathbf{l}_T(u_0, v_0, R_1)),$$

$$\widetilde{\mathbf{h}}_{\Sigma_2}^\top(\mathbf{l}_1(u_0, v_0, R_2)), \cdots, \widetilde{\mathbf{h}}_{\Sigma_2}^\top(\mathbf{l}_T(u_0, v_0, R_2)),$$

Tola et al. "Daisy: An efficient dense descriptor applied to wide- baseline stereo." Pattern Analysis and Machine Intelligence, IEEE Transactions
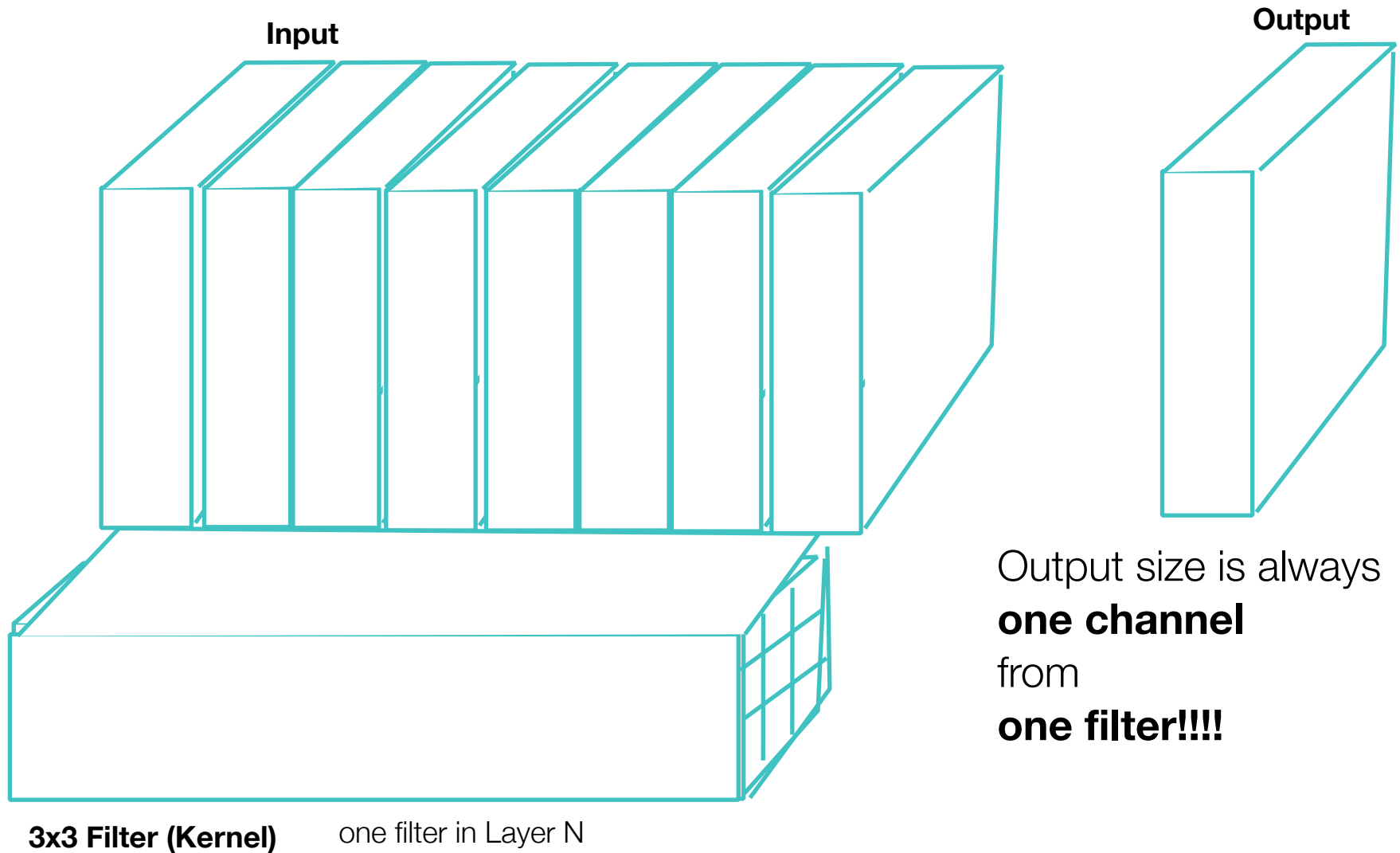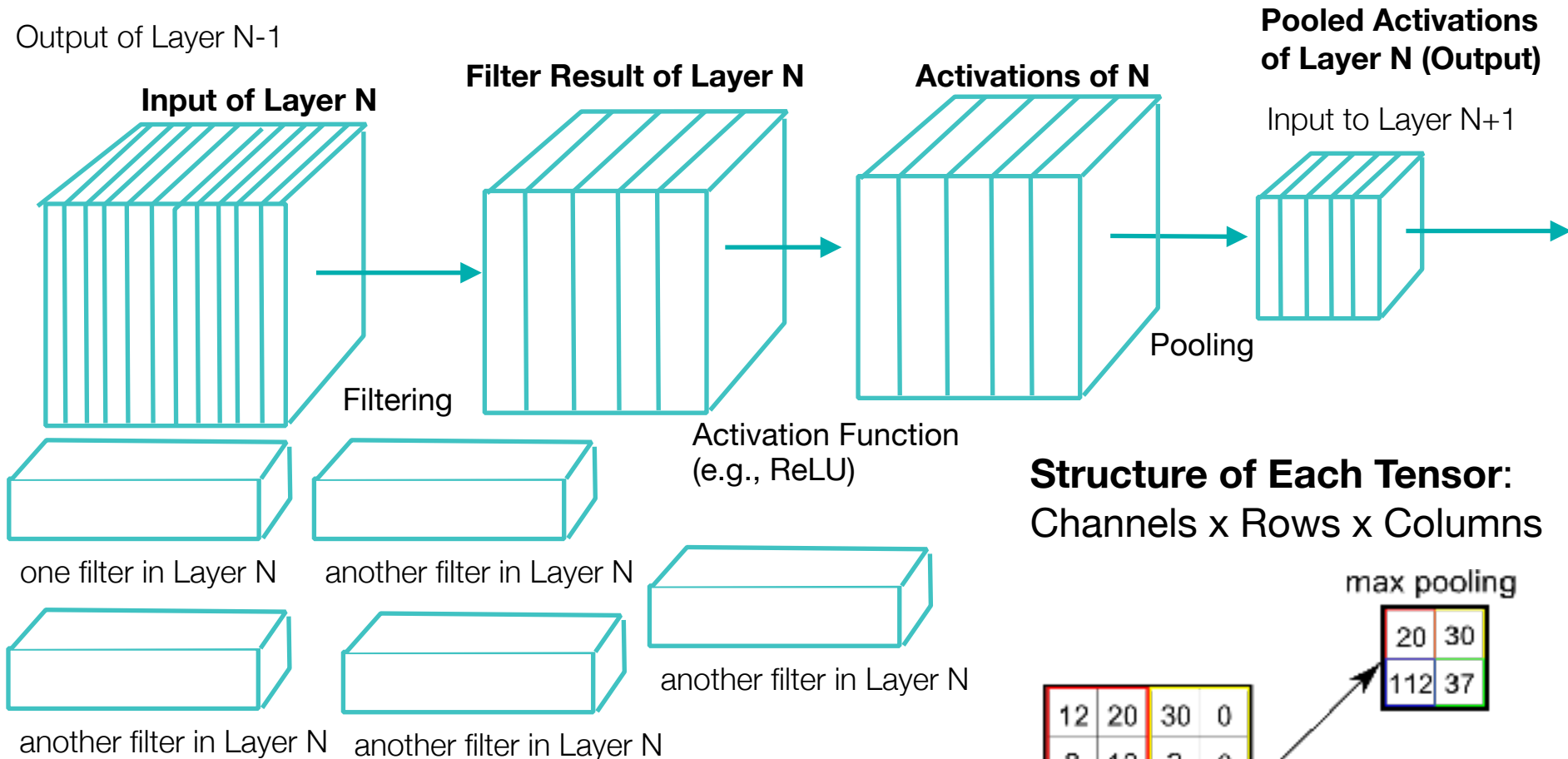
down sample
convolve with 128 filters

64

64

32

16

filtering

64

filtering

32

16

filtering

128

flatten and fully connected

fully connected

fully connected

1x1x2

128x128x3
**convolve with 64 filters**

64

128

128

1x1x4096

1x1x4096

**down sample
convolve with 128 filters**

**Blue Tensors**: Outputs tensors of Each Layer

**Learned Params**: Weights in Each Arrow

**Input**

**Output**

**3x3 Filter (Kernel)**     one filter in Layer N

Output size is always
**one channel**
from
**one filter!!!!**

Output of Layer N-1

**Pooled Activations of Layer N (Output)**

**Filter Result of Layer N**

**Activations of N**

**Input of Layer N**

Input to Layer N+1

Filtering

Pooling

Activation Function (e.g., ReLU)

one filter in Layer N

another filter in Layer N

another filter in Layer N

another filter in Layer N

another filter in Layer N

**Structure of Each Tensor**:
Channels x Rows x Columns

max pooling

| 20 | 30 |
|----|----|
| 112 | 37 |

| 12 | 20 | 30 | 0 |
|----|----|----|----|
| 8 | 12 | 2 | 0 |
| 34 | 70 | 37 | 4 |
| 112 | 100 | 25 | 12 |

average pooling

| 13 | 8 |
|----|----|
| 79 | 20 |

One channel input

One channel output

**Self Test:** What are the learned parameters?
A. Activations
B. Pooling Weights
C. Filter Weights        D. All of these

# CNN Overview



- Conv. layer(s):
  - filtering
  - activation
  - pooling
  - Each pooling layer *can* make the input image "smaller"
    - allows for "Information Distillation"
    - less dependence on exact pixels
- Final layers are densely connected
  - typically multi-layer perceptrons

$\mathbf{a}^{(1)}$  $N \times 1$

$\mathbf{W}^{(1)}$  $N \times S^1$

$\mathbf{z}^{(1)}$  $S^1 \times 1$

$\varphi$

$\mathbf{a}^{(2)}$  $S^1 \times 1$

$\mathbf{W}^{(2)}$  $S^1 \times S^2$

$\mathbf{z}^{(2)}$  $S^2 \times 1$

$\varphi$

$\mathbf{a}^{(3)}$  $S^2 \times 1$

0
1
2
3
4
5
6
7
8
9

what we did in the past

If image is 9x9, and each fully connected layer is 20 hidden neurons wide, how many parameters are in this NN (ignore bias)?

for 9x9,     $9^2 \times 20 + (20 \times 10) = 1{,}820$ parameters

$$(K^2 \times 20) + (20 \times 10) = 200 + 20\,K^2$$

12

three filters

3x3 kernels

**(3x3) x3 kernels = 27 weights**

three outputs

three activations

three pooled activations

flatten

max pool

max pool

max pool

output convolutions

If image is 9x9, and each fully connected layer is 20 hidden neurons wide, how many parameters are in this NN (ignore bias)?

3x3x3 = 27 inputs

$\mathbf{a}^{(1)}$   $N \times 1$

$\mathbf{W}^{(1)}$   $N \times S^1$

$\mathbf{z}^{(1)}$   $S^1 \times 1$

$\varphi$

$\mathbf{a}^{(2)}$   $S^1 \times 1$

$\mathbf{W}^{(2)}$   $S^1 \times S^2$

$\mathbf{z}^{(2)}$   $S^2 \times 1$

$\varphi$

$\mathbf{a}^{(3)}$   $S^2 \times 1$

$27 + (27 \times 20) + (20 \times 10) = 767$

0
1
2
3
4
5
6
7
8
9

three outputs

three activations

three pooled activations

flatten

φ

max pool

φ

max pool

3x3 kernels

φ

max pool

**kernel size = k x k**

**num kernels = $N_k$**

**k x k x $N_k$ weights**

output convolutions

pooling outputs

**convolutional params**

**$N_k$ x $k^2$**

filter dimension

num filters

**if stride = k**

**(K/k) x (K/k) x $N_k$**

**Input to MLP**

**$N_k$ x ($K^2/k^2$)**

image dimension

$\mathbf{a}^{(1)}$   $\mathbf{z}^{(1)}$   $\mathbf{a}^{(2)}$   $\mathbf{z}^{(2)}$   $\mathbf{a}^{(3)}$

$\mathbf{W}^{(1)}$   $\mathbf{W}^{(2)}$

N x 1   $S^1$ x 1   $S^1$ x1   $S^2$ x 1   $S^2$ x 1

N x $S^1$   $S^1$ x $S^2$

φ   φ

three outputs

three activations

three pooled activations

flatten



**X**

3x3 kernels

**F$_{a,b,c}$**

output convolutions **O**

$\sigma$

max pool

$\sigma$

max pool

$\sigma$

max pool

Derivative of max pool is easy to compute:

$$\frac{\partial}{\partial L_i}\text{pool}(L_i) = 1, \text{ if } L_i \text{ is max}$$

$$0, \text{ else}$$

output activations

**L** = $\sigma(\mathbf{O})$

**X$^{(Next)}$** = pool( $\sigma(\mathbf{O})$ )

Derivative of convolution is more involved:

$$\begin{pmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{pmatrix} = \text{Convolution}\left( \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{pmatrix}, \begin{pmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{pmatrix} \right)$$

Output O    Input X    Filter F

**Back propagate to previous Layer**

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial X}$$

X

O

$$\frac{\partial L}{\partial O}$$

Convolution

F

$$\frac{\partial L}{\partial F} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial F}$$

**Update weights of F**

https://medium.com/@pavisj/convolutions-and-backpropagations-46026a8f5d2c

15

# CNNs Back Propagation

**Pooled Activations of Layer N (Output)**

Output of Layer N-1

**Input of Layer N**

**Filter Output of Layer N**

**Activations of N**

Input to Layer N+1

$$\frac{\partial O^{(N)}}{\partial X^{(N)}}$$

$$X^{(N)}_{i,j,c}$$

$$O^{(N)}_{i,j,c}$$

$$\phi(O^{(N)}_{i,j,c})$$

OR

$$L^{(N)}_{i,j,c}$$

$$X^{(N+1)}_{i,j,c}$$

Filtering

Activation

Pooling

$$F^{(N)}_{a,b,c}$$

one filter in Layer N

$$\frac{\partial O^{(N)}}{\partial F^{(N)}}$$

$$\frac{\partial L^{(N)}}{\partial O^{(N)}}$$

$$V^{(N+1)}_{pool}$$

sensitivity through pooling

$$V^{(N+1)}$$

sensitivity

$$\frac{\partial O^{(N+1)}}{\partial X^{(N+1)}}$$

Now we can calc partial derivative

$$\frac{\partial L^{(N)}}{\partial F^{(N)}} = \frac{\partial O^{(N)}}{\partial F^{(N)}} \cdot \frac{\partial L^{(N)}}{\partial O^{(N)}}$$

Just incorporate sensitivity, to get weight update

$$\frac{\partial J_{obj}}{\partial F^{(N)}} = \frac{\partial O^{(N)}}{\partial F^{(N)}} \cdot \frac{\partial L^{(N)}}{\partial O^{(N)}} \cdot V^{(N+1)}_{pool}$$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

$$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22}$$

$$O_{21} = X_{21}F_{11} + X_{22}F_{12} + X_{31}F_{21} + X_{32}F_{22}$$

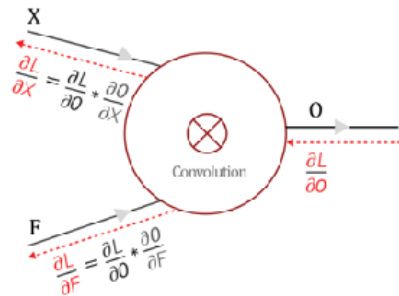$$O_{22} = X_{22}F_{11} + X_{23}F_{12} + X_{32}F_{21} + X_{33}F_{22}$$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

https://medium.com/@pavisj/convolutions-and-backpropagations-46026a8f5d2c

# Gradient of Convolution

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial O}\frac{\partial O}{\partial X}$$

for back propagation

$$\frac{\partial L}{\partial F} = \frac{\partial L}{\partial O}\frac{\partial O}{\partial F}$$

for weight updates



$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

*Finding derivatives with respect to* $F_{11}$, $F_{12}$, $F_{21}$ *and* $F_{22}$

$$\frac{\partial O_{11}}{\partial F_{11}} = X_{11} \qquad \frac{\partial O_{11}}{\partial F_{12}} = X_{12} \qquad \frac{\partial O_{11}}{\partial F_{21}} = X_{21} \qquad \frac{\partial O_{11}}{\partial F_{22}} = X_{22}$$

derivative of every $O_{ij}$ w.r.t. $F_{11}$

$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}}*\frac{\partial O_{11}}{\partial F_{11}} + \frac{\partial L}{\partial O_{12}}*\frac{\partial O_{12}}{\partial F_{11}} + \frac{\partial L}{\partial O_{21}}*\frac{\partial O_{21}}{\partial F_{11}} + \frac{\partial L}{\partial O_{22}}*\frac{\partial O_{22}}{\partial F_{11}}$$

$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}}*\frac{\partial O_{11}}{\partial F_{12}} + \frac{\partial L}{\partial O_{12}}*\frac{\partial O_{12}}{\partial F_{12}} + \frac{\partial L}{\partial O_{21}}*\frac{\partial O_{21}}{\partial F_{12}} + \frac{\partial L}{\partial O_{22}}*\frac{\partial O_{22}}{\partial F_{12}}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}}*\frac{\partial O_{11}}{\partial F_{21}} + \frac{\partial L}{\partial O_{12}}*\frac{\partial O_{12}}{\partial F_{21}} + \frac{\partial L}{\partial O_{21}}*\frac{\partial O_{21}}{\partial F_{21}} + \frac{\partial L}{\partial O_{22}}*\frac{\partial O_{22}}{\partial F_{21}}$$

$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}}*\frac{\partial O_{11}}{\partial F_{22}} + \frac{\partial L}{\partial O_{12}}*\frac{\partial O_{12}}{\partial F_{22}} + \frac{\partial L}{\partial O_{21}}*\frac{\partial O_{21}}{\partial F_{22}} + \frac{\partial L}{\partial O_{22}}*\frac{\partial O_{22}}{\partial F_{22}}$$

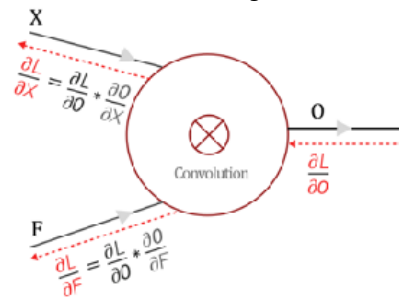$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}}*X_{11} + \frac{\partial L}{\partial O_{12}}*X_{12} + \frac{\partial L}{\partial O_{21}}*X_{21} + \frac{\partial L}{\partial O_{22}}*X_{22}$$

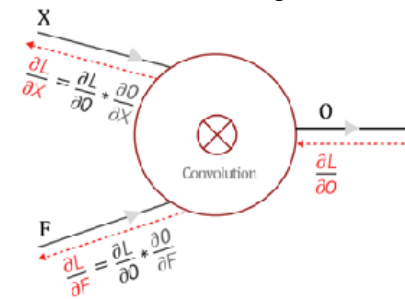$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}}*X_{12} + \frac{\partial L}{\partial O_{12}}*X_{13} + \frac{\partial L}{\partial O_{21}}*X_{22} + \frac{\partial L}{\partial O_{22}}*X_{23}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}}*X_{21} + \frac{\partial L}{\partial O_{12}}*X_{22} + \frac{\partial L}{\partial O_{21}}*X_{31} + \frac{\partial L}{\partial O_{22}}*X_{32}$$

$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}}*X_{22} + \frac{\partial L}{\partial O_{12}}*X_{23} + \frac{\partial L}{\partial O_{21}}*X_{32} + \frac{\partial L}{\partial O_{22}}*X_{33}$$



Filter updates = Convolution ( Input , Derivative From activation! )

# Gradient of Convolution

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial O}\frac{\partial O}{\partial X}$$

for back propagation

$$\frac{\partial L}{\partial F} = \frac{\partial L}{\partial O}\frac{\partial O}{\partial F}$$

for weight updates



$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

Differentiating with respect to $X_{11}, X_{12}, X_{21}$ and $X_{22}$

$$\frac{\partial O_{11}}{\partial X_{11}} = F_{11} \quad \frac{\partial O_{11}}{\partial X_{12}} = F_{12} \quad \frac{\partial O_{11}}{\partial X_{21}} = F_{21} \quad \frac{\partial O_{11}}{\partial X_{22}} = F_{22}$$

Similarly, we can find local gradients for $O_{12}, O_{21}$ and $O_{22}$

New sensitivity

Rotated Filter

Derivative From activation! (zero padded)

# Summary
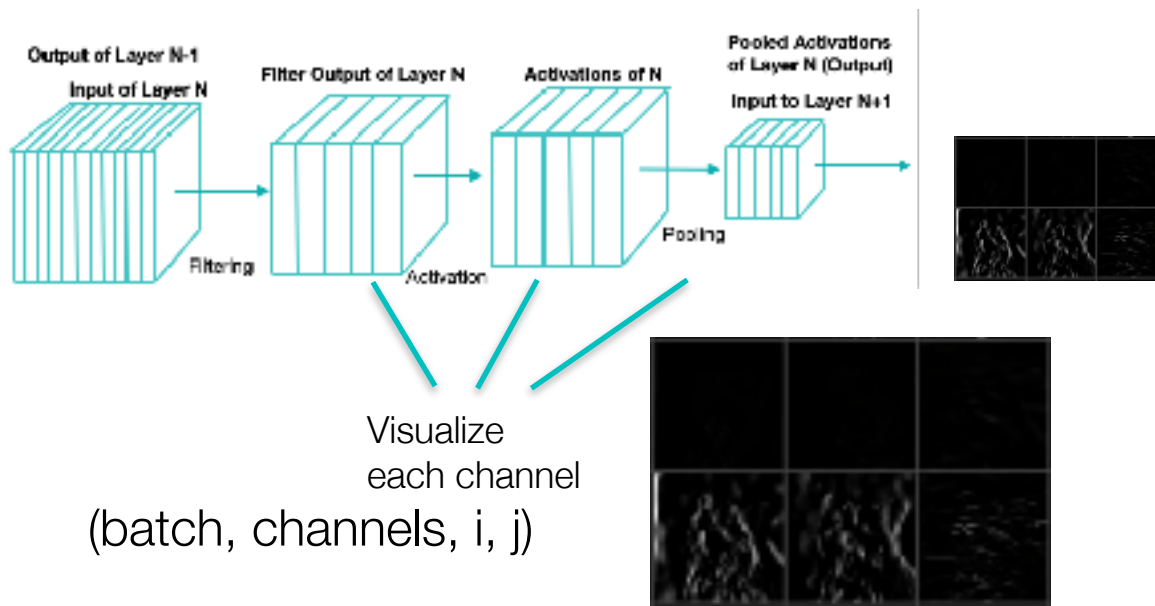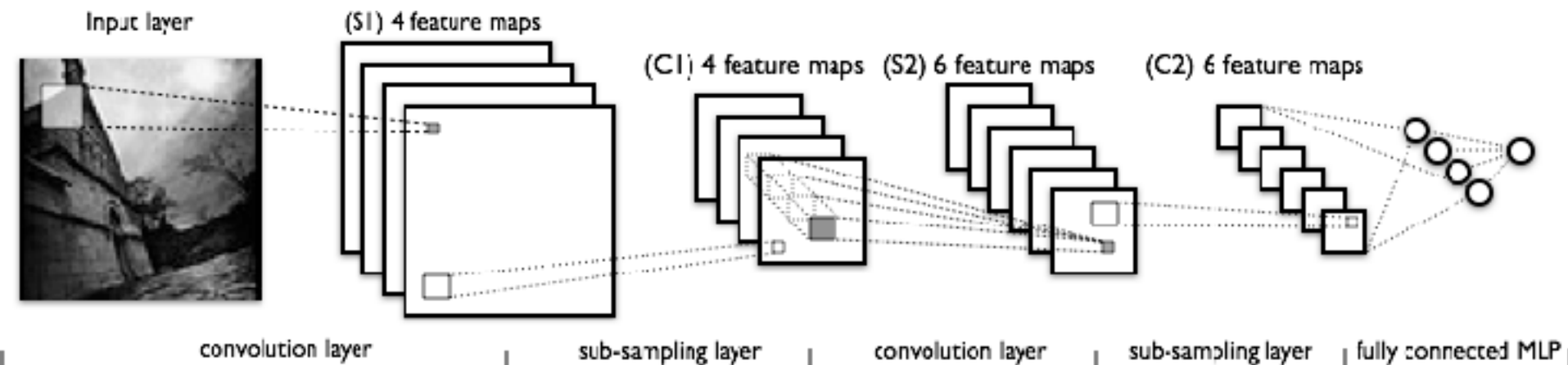
# CNN Gradient

- Takeaways:
  - Derivative of a convolutional layer is calculated through two additional convolutions
    - One for filter updates
    - One for calculating a new sensitivity
  - We need to run convolution fast in order to speed up both:
    - feedforward operations (inference and training)
    - back propagation (training)
  - Another great resource:
    - https://becominghuman.ai/back-propagation-in-convolutional-neural-networks-intuition-and-code-714ef1c38199

Visualize
each channel
(batch, channels, i, j)

Naming:
conv1 (output of conv)
p1 (output of pooling)
n1 (output of normalization)

Deep Visualization Toolbox

yosinski.com/deepvis

#deepvis

Jason Yosinski    Jeff Clune    Anh Nguyen    Thomas Fuchs    Hod Lipson

Cornell University    UNIVERSITY OF WYOMING    NASA Jet Propulsion Laboratory California Institute of Technology

https://github.com/yosinski/deep-visualization-toolbox

**Demo**

Convolutional Neural Networks
  in TensorFlow
  with Keras



`11. Convolutional Neural Networks.ipynb`

# Next Lecture

- More CNN architectures and CNN history