# Lecture Notes for
# **Machine Learning in Python**

[ 👨‍🏫, 👨‍💻, 🐍, 👨‍🔬 ]

## Professor Eric Larson
## **Basic Convolutional Neural Networks**

# Logistics and Agenda

- Logistics
  - Wide/Deep due soon!
  - Remember: late turn in…
- Agenda
  - Wide/Deep Finish Demo and Town Hall
  - Basic CNN architectures and Demo

# Class Overview, by topic

| | |
|---|---|
| **Table Data Visualization** | Numpy, Pandas, Seaborn<br>Overviews with some in-depth discussion |
| **Dimension Reduction and Image Processing** | Scikit-learn, Scikit Image,<br>Intuition only, Some mathematics |
| **Linear and Logistic Regression** | Numpy, Recreate API for Scikit-learn<br>Detailed mathematics for simple optimization<br>intuition for advanced optimization |
| **Neural Networks and Back Prop.** | Numpy<br>Detailed mathematics for NN operations |

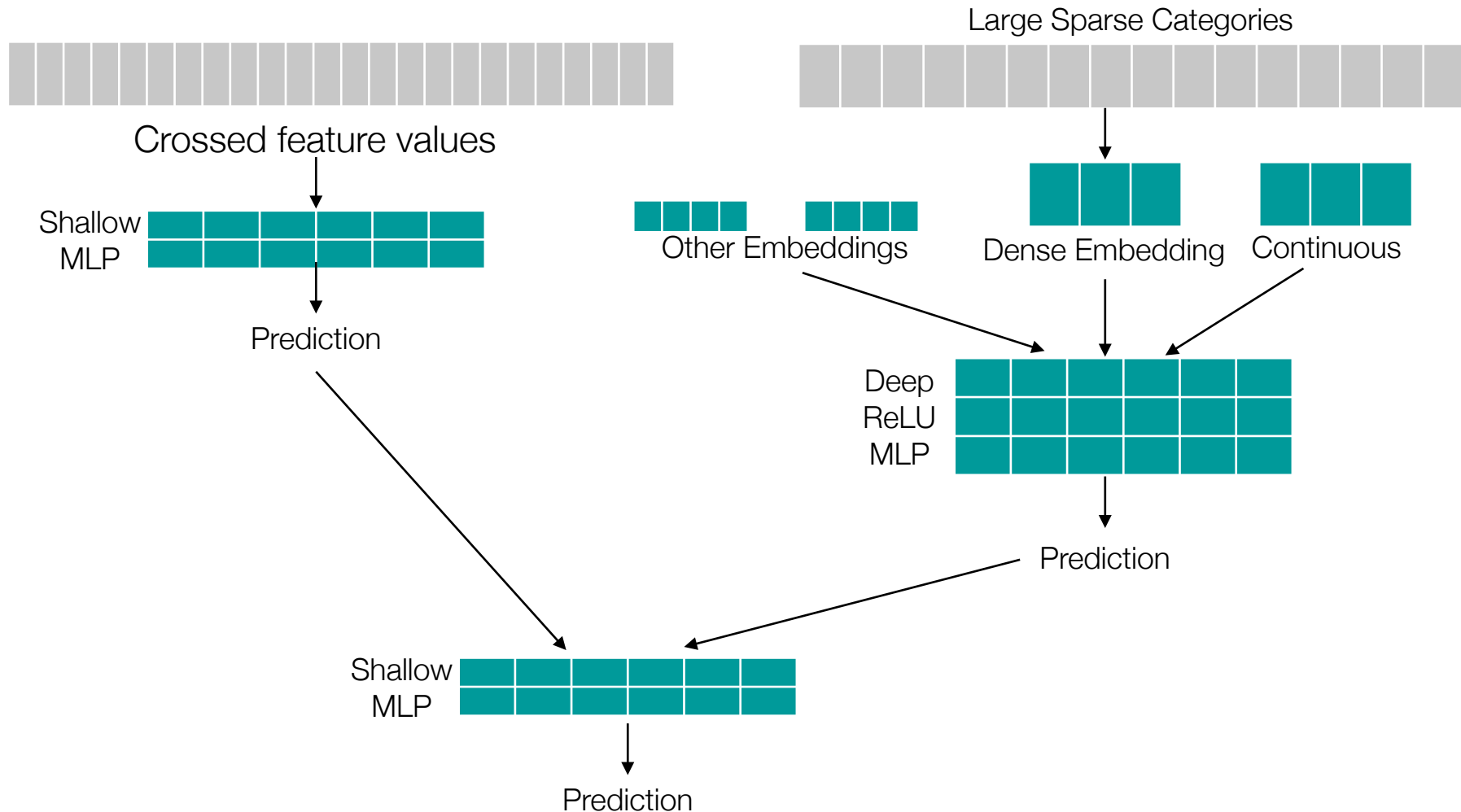**Wide and Deep Networks**   **Convolutional Networks**   **Recurrent Networks**

Keras, Tensorflow
Intuition, Detailed implement.

**Ethics in Language Models**

ConceptNet
Case studies

# Last Time:

- Deep refers to increasingly smaller hidden layers
- Embed into sparse representations via ReLU

Large Sparse Categories

Crossed feature values

Shallow MLP

Prediction

Other Embeddings

Dense Embedding

Continuous

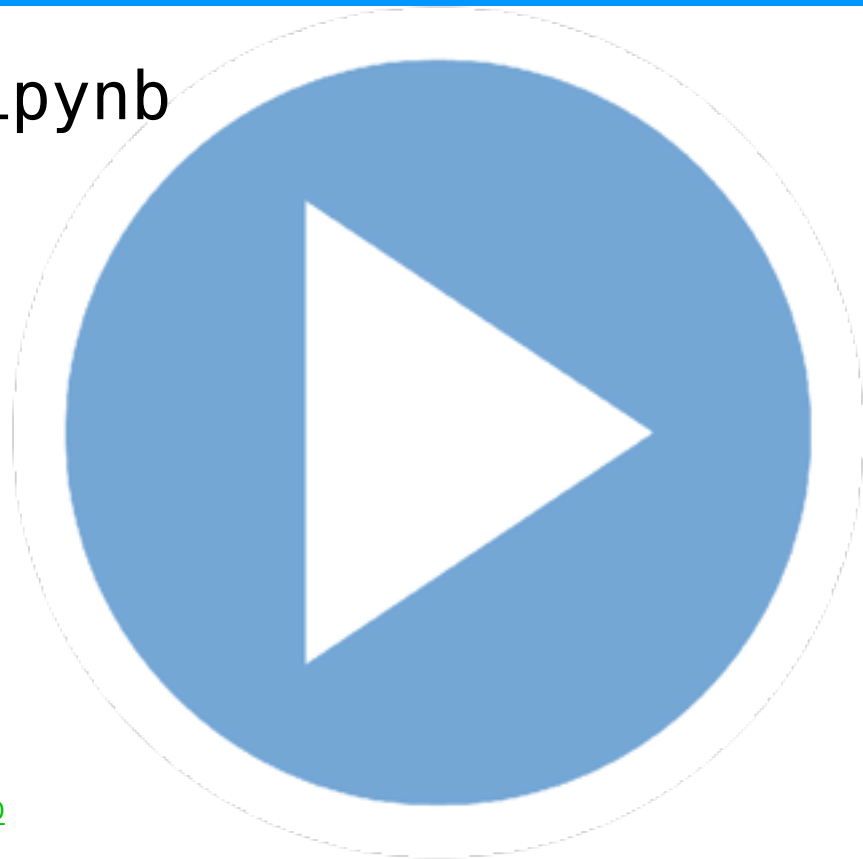Deep ReLU MLP

Prediction

Shallow MLP

Prediction

`10. Keras Wide and Deep.ipynb`

The awful dataset:
Toy Census Data Example

Other tutorials:

https://www.tensorflow.org/tutorials/wide_and_deep

# Town Hall, Wide and Deep Networks

When $p < 0.05$

# Convolutional Neural Networks

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

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, $\mathbf{k}$ 3x3

| 20 | 21 | 36 | ... | ... | ... | ... |
|----|----|----|-----|-----|-----|-----|
| ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |

output image, $\mathbf{O}$

8

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

**Observe**: Can also express the convolution as matrix multiplication with a reshaped input and filter!

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

# What we did before

- the gradient (2D derivative)



$$\|\nabla J\| = \sqrt{J_x^2 + J_y^2}$$

$$O = \tan^{-1}(Jx \,/\, Jy)$$

images: Jianbo Shi, Upenn

take normalized histogram at point $u,v$

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

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

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

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

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

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

flatten and pass through layer

pass through layer

pass through layer

128x128x3

convolve with 64 filters

down sample
convolve with 128 filters

**Blue Tensors**: Outputs of Each Layer

**Learned Params**: Weights in Each Filter and Fully Connected Layer

# Convolution in a CNN

**Input**

**Output**

**3x3 Filter (Kernel)**

**Output of Layer N-1**

**Input of Layer N**

**Filter Result of Layer N**

**Activations of N**

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

**Input to Layer N+1**

Filtering

Activation Function
(e.g., ReLU)

Pooling

one filter in Layer N

x Num Filters

another filter in Layer N

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

What are the learned parameters?
A. Activations
B. Pooling
C. Filters

max pooling

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

average pooling

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

One channel input

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

One channel output

14

- Initial layer(s):
  - convolution
  - 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

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)?

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

for 9x9 = $200 + 20 \times 9^2 =$ 1,820 parameters

# Simple Example: From Fully Connected to CNN

3x3 kernels

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

output convolutions

max pool

max pool

max pool

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 weights**

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

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

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

$\varphi$

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

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

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

$\varphi$

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

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

0
1
2
3
4
5
6
7
8
9

**kernel size = k x k**

**num kernels = $N_k$**

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

3x3 kernels

output convolutions

max pool

max pool

max pool

pooling outputs
if stride = k
(K/k) x (K/k) x $N_k$

**convolutional params**

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

num filters

filter dimension

**Input to MLP**

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

image dimension

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

$(N+1) \times 1$

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

$(N+1) \times S^1$

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

$S^1 \times 1$

$\boldsymbol{\varphi}$

$\mathbf{a}^{(2)}$

$(S^1+1) \times 1$

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

$(S^1+1) \times S^2$

$\mathbf{z}^{(2)}$

$S^2 \times 1$

$\boldsymbol{\varphi}$

$\mathbf{a}^{(3)}$

$S^2 \times 1$

# CNN gradient



**X**(L-1)

3x3 kernels

**F**a,b

output convolutions **O**

max pool

max pool

max pool

Derivative of max pool is easy:

for each input $X_i$

$$\text{pool}'(X_i) = 1 \text{ if } X_i \text{ is max}$$
$$0 \text{ else}$$

$$\mathbf{X}^{(L)} = \text{pool}(\sigma(\mathbf{O}))$$

**Back propagate to previous Layer**

**Use to update weights of F**

Derivative of convolution is more involved:

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

**Output of Layer N-1**

**Input of Layer N**

**Filter Output of Layer N**

**Activations of N**

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

**Input to Layer N+1**

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

$$X_{i,j}$$

$$O_{i,j}$$

$$\phi(O_{i,j})$$

$$L_{i,j}$$

Pooling

Filtering

Activation

$$F_{a,b}$$

one filter in Layer N

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

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

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

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

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

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

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

https://medium.com/@pavisj/convolutions-and-backpropagations-46026a8f5d2c
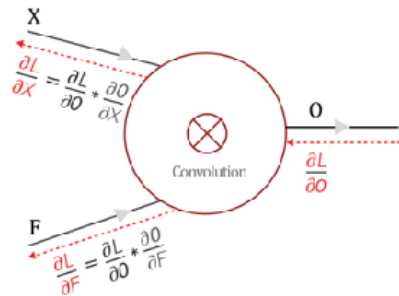
$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial O}\frac{\partial O}{\partial X}$$

for back propagation

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

for weight updates

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

Convolution

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

$$\frac{\partial L}{\partial F} = \frac{\partial L}{\partial O} * \frac{\partial O}{\partial F}$$
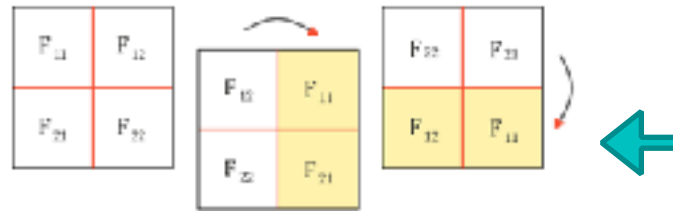
$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$
$$O_{12} = X_{12}F_{11} + X_{13}F_{12} + X_{22}F_{21} + X_{23}F_{22}$$
$$O_{21} = X_{21}F_{11} + X_{22}F_{12} + X_{31}F_{21} + X_{32}F_{22}$$
$$O_{22} = X_{22}F_{11} + X_{23}F_{12} + X_{32}F_{21} + X_{33}F_{22}$$

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

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

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

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

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

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

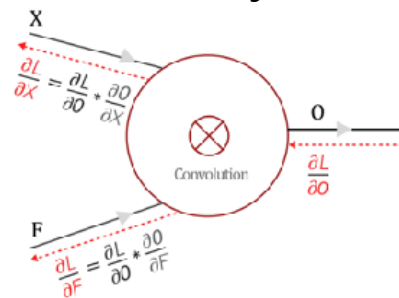$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}} * X_{22} + \frac{\partial L}{\partial O_{12}} * X_{23} + \frac{\partial L}{\partial O_{21}} * X_{32} + \frac{\partial L}{\partial O_{22}} * X_{33}$$

Filter updates

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

Input

Sensitivity from next layer

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

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

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 $=$ Full Convolution ( Rotated Filter , Sensitivity from next layer (zero padded) )

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

# Summary

**Filters at layer L-1**



**Filters at layer L**



New sensitivity · Rotated Filter · Sensitivity from next layer

Filter updates · Input · Sensitivity from next layer

New sensitivity · Rotated Filter · Sensitivity from next layer

Filter updates · Input · Sensitivity from next layer

# 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

# Next Lecture

- More CNN architectures and CNN history