# MAC5921 – Deep Learning

Aula 04 – 22/08/2023

Nina S. T. Hirata

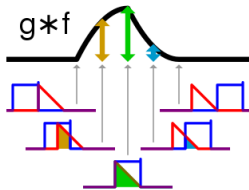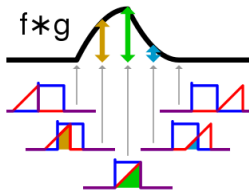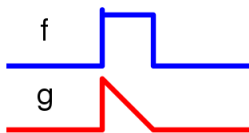## Convolution

Domínio contínuo

$$(\mathbf{f} * \mathbf{g})(\mathbf{t}) = \int_{-\infty}^{+\infty} \mathbf{f(a)g(t-a)da} = \int_{-\infty}^{+\infty} \mathbf{f(t-a)g(a)da}$$

Domínio discreto

$$(\mathbf{f} * \mathbf{g})(\mathbf{t}) = \sum_{-\infty}^{+\infty} \mathbf{f(a)g(t-a)} = \sum_{-\infty}^{+\infty} \mathbf{f(t-a)g(a)}$$

Convolução domínio discreto 2D

$$S(i,j) = \sum_m \sum_n I(m,n)K(i-m,j-n) = \sum_m \sum_n I(i-m,j-n)K(m,n)$$

Cross-correlation domínio discreto 2D

$$S(i,j) = \sum_m \sum_n I(i+m,j+n)K(m,n)$$

Supondo que $K$ tem suporte finito,

$$S(i,j) = \sum_m \sum_n I(i-m, j-n)K(m,n)$$

representa a ideia de "sliding window"

# 1D Convolution

Image Matrix

Kernel Matrix

Output Matrix

(http://machinelearninguru.com/computer_vision/basics/convolution/image_convolution_1.html)

| 1/9 | 1/9 | 1/9 |
|-----|-----|-----|
| 1/9 | 1/9 | 1/9 |
| 1/9 | 1/9 | 1/9 |

(http://aishack.in/tutorials/image-convolution-examples/)

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 9 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

| -1 | -1 | -1 |
|----|----|----|
| -1 | 8  | -1 |
| -1 | -1 | -1 |

(http://homepages.inf.ed.ac.uk/rbf/HIPR2/canny.htm)

(http://aishack.in/tutorials/image-convolution-examples/)

$$\begin{bmatrix} 0 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Input Channel #1 (Red)  Input Channel #2 (Green)  Input Channel #3 (Blue)

Kernel Channel #1  Kernel Channel #2  Kernel Channel #3

$$308 \quad + \quad -498 \quad + \quad 164 \quad + 1 = -25$$

Bias = 1

Output

(http://machinelearninguru.com/computer_vision/basics/convolution/convolution_layer.html)

Concepts: padding, stride, kernel size, number of filters/maps

ReLU

$$F(x) = \max(0, x)$$

- via strides – skip pixels by steps of size $d$

- via pooling – aggregate $d \times d$ pixels

224x224x64

pool

112x112x64

224

224

downsampling

112

112

Source: https://iitmcvg.github.io/summer_school/DLSession3/

- Convolutonal layers: feature extraction

- Fully connected layers: classification

**Por que usar convolução?**

weight sharing: Redução no número de pesos a serem ajustados

propriedade de "equivariance to translation"

**Main CNN architectures used in ILSVRC**

## CNN models

- **AlexNet**, 2012: winner of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), 60M network parameters

  (Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012, pages 1097–1105)

- **VGG**-**11, 16 e 19**, 2014: 8, 13 e 16 convolutional layers, VGG-19 138M network parameters

  (Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for largescale image recognition. CoRR, abs/1409.1556, 2014)

## CNN models

- **GoogleLeNet** (Inception), 2014: winner of ILSVRC 2014, inception layers, 7M network parameters

  (C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. CVPR 2015, pages 1–9)

- **Residual Network** (ResNet), 2015: winner of ILSVRC 2015, 25.5M network parameters, residual block, vanishing gradient

  (K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. CVPR 2016, pages 770–778)

- **AlexNet**, 2012: winner of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC), 60M network parameters

  (Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012, pages 1097–1105)

## AlexNet (winner ILSVRC 2012)

- succeeded on training a large CNN
- 60 million parameters
- used 2 GPUs
- proposed ReLU
- used dropout
- used data augmentation (rotation, scale, crops, etc)
- indication that number of layers is important

Images that maximize the response of the filters (AlexNet)

- **VGG-11, 16 e 19**, 2014: 8, 13 e 16 convolutional layers, VGG-19 138M network parameters

(Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for largescale image recognition. CoRR, abs/1409.1556, 2014)
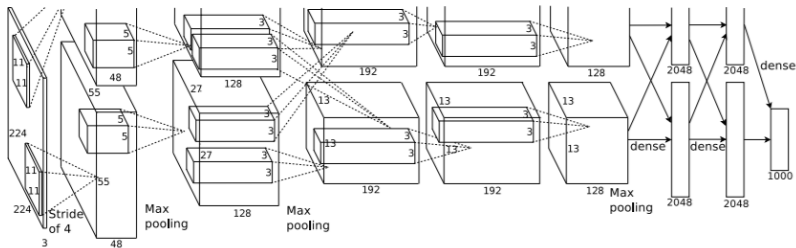
## VGGNet (2nd place ILSVRC 2014)

- multiple layers (ex., 16 e 19)
- only 3x3 convolutions
- stride 1 for the convolutions (AlexNet used stride ¿ 1)
- first layers with few filters and gradually increasing as we advance throuh the layers
- 144 millions of parameters

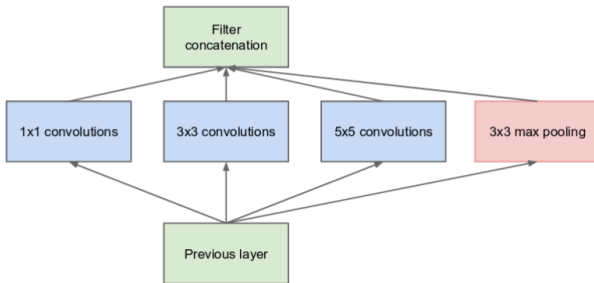- **GoogleLeNet** (Inception), 2014: winner of ILSVRC 2014, inception layers, 7M network parameters

  (C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A.

  Rabinovich. Going deeper with convolutions. CVPR 2015, pages 1–9)

(a) Inception module, naïve version

Camada de entrada $= 28{\times}28{\times}192$

Convolução 5x5: com 32 filtros, camada de saída $= 28{\times}28{\times}32$. Para cada ponto da camada de saída, o número de multiplicações necessárias é $5{\times}5{\times}192$. Como são $28{\times}28{\times}32$ pontos na camada de saída, o total de multiplicações necessárias é $5{\times}5{\times}192 \times 28{\times}28{\times}32 = 120$ milhões.

(b) Inception module with dimension reductions

Entrada 28x28x192 $\implies$ 16 filtros 1x1 $\implies$ camada intermediária 28x28x16 $\implies$ 32 convoluções 5x5 $\implies$ saída 28x28x32. Total de multiplicações: cada ponto na camada intermediária requer 1x1x192 multiplicações, e portanto 1x1x192 x 28x28x16 = 2.4 milhões para cálculo da camada intermediária; cada ponto na camada de saída requer 5x5x16 multiplicações e portanto 5x5x16 x 28x28X32 = 10 milhões para o cálculo da camada de saída. Total = 12.4 milhões.

GoogLeNet (aka "Inception") Architecture

"Inception" Module.

Main Classifier

Pr(dog)

Auxiliary Classifiers

Szegedy et al, 2014

Proprietary & Confidential
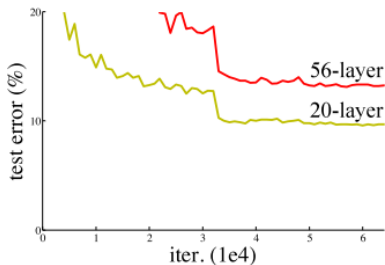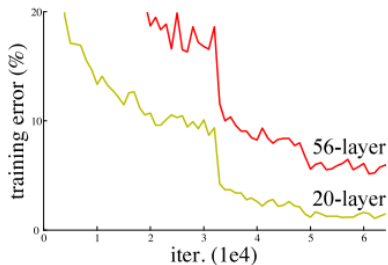
## Inception - GoogleLeNet (winner ILSVRC 2014)

- concept of "network in network"
- inception module ("choices" made during training)
- use of 1x1 convolution
- auxiliary output to reinforce activation of intermediary layers
- more aggressive data augmentation tha AlexNet
- 4 million parameters (?)
- variants (e.g., change 5x5 with two 3x3, or change 3x3 with a 1x3 and a 3x1)
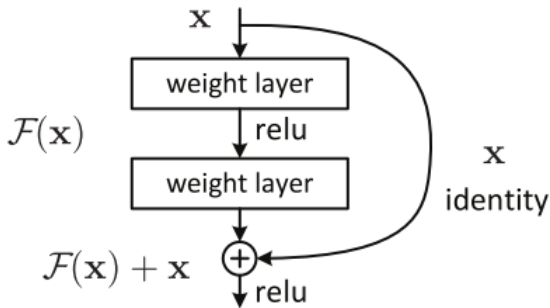
Apesar de ser aceito que profundidade da rede é importante, observou-se que camadas adicionais degradavam a acurácia no conjunto de treinamento



- **Residual Network** (ResNet), 2015: winner of ILSVRC 2015, 25.5M network parameters, residual block, vanishing gradient

  (K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. CVPR 2016, pages 770–778)
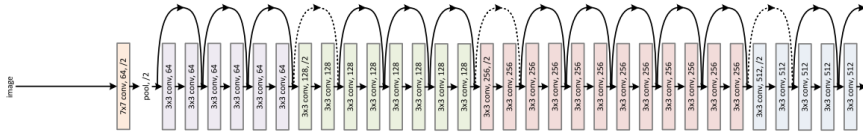
Entrada $x$. Pode ser conveniente aprender um mapeamento ($H(x)$) que seja a identidade. Mas representar identidade por meio de transformações não lineares não é simples. Assim, em vez de $H$, aprende-se o resíduo $\mathcal{F}(x)$, de mode que $H(x) = \mathcal{F}(x) + x$ e assim a identidade corresponde a resíduo zero
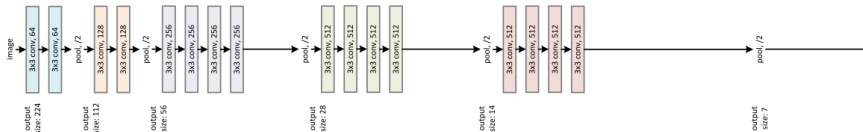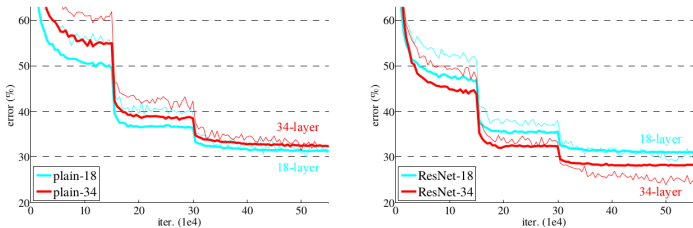
Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

## ResNet (winner ILSVRC 2015)

- identificam uma degradação associada ao aumento de camadas, que aparentemente não é *overfitting* nem *vanishing gradient*
- propõe módulo para aprender o resíduo da transformação desejada
- aplica *batch normalization* após cada convolução e antes da ativação
- não usa *dropout*
- conseguem treinar redes com mais de 100 camadas
- testado em outros datasets (CIFAR, PASCAL, MS-COCO)
- ResNet 56 layers: 0.85M parameters (?)
- ResNet 110 layers: 1.7M parameters (?)

**Training of CNN**

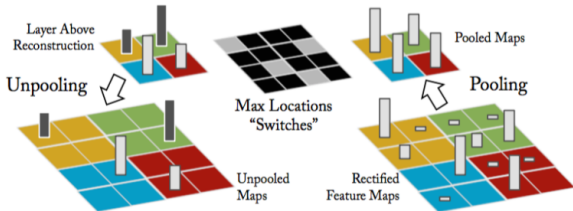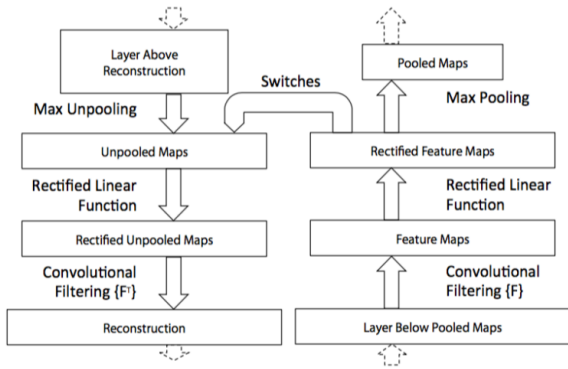It is done using the backpropagation algorithm
(gradient descent)

Kernels são ajustados nesse processo – detalhes ficam para depois

**Convolutional layers act as feature extractors**. How ?

Zeiler et al., 2013, proposed a visualizaton technique

At a given convolutional layer, take the node with largest activation

Reconstruct the input image by backpropagating from that node

(Zeiler et al., 2013)

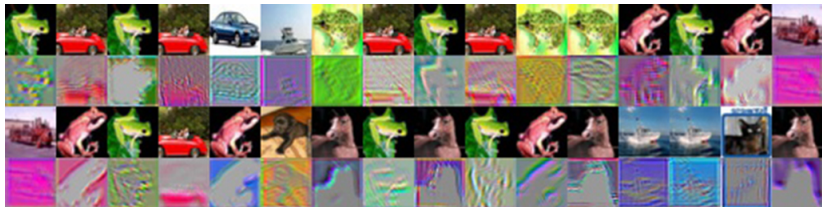Images reconstructed from neurons 1 to 32 of layer 1, for one input image



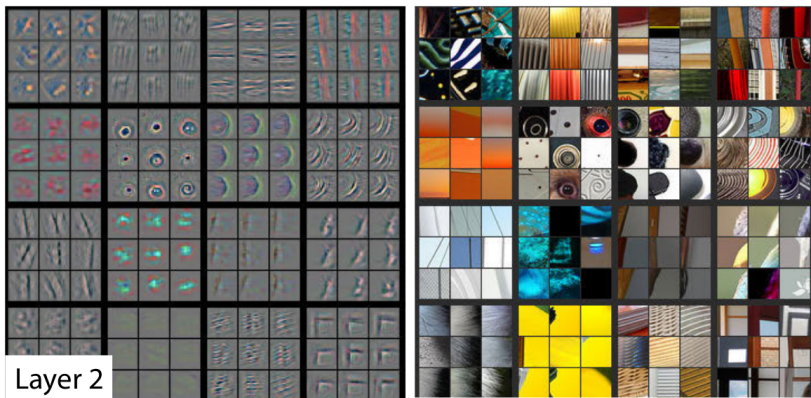(http://kvfrans.com/visualizing-features-from-a-convolutional-neural-network/)

Images reconstructed from neuron 7 of layer 1, for multiple input images

For each of the 32 neurons, reconstruction of the images that most
activated the neuron

(Zeiler et al.) 16 neurons in layer 2; for each neuron, reconstruction of the 9 images that correspond to strongest activation of the node



Layer 2

(Zeiler et al.) 12 neurons in layer 3; for each neuron, reconstruction of the 9 images that correspond to strongest activation of the node



Layer 3