

Deep learning

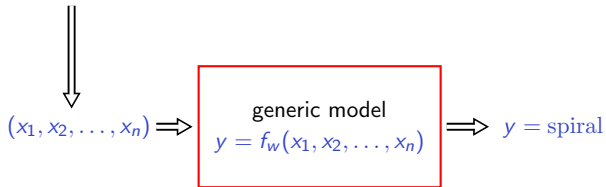
- Representation learning
- Complex pipeline learning

Traditional machine learning: Feature computation

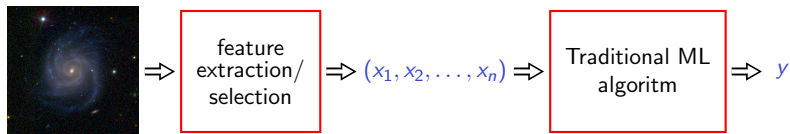


Features that describe the object

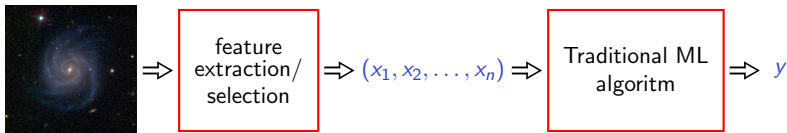
- concentration
- asymmetry
- smoothness
- entropy
- spirality, etc



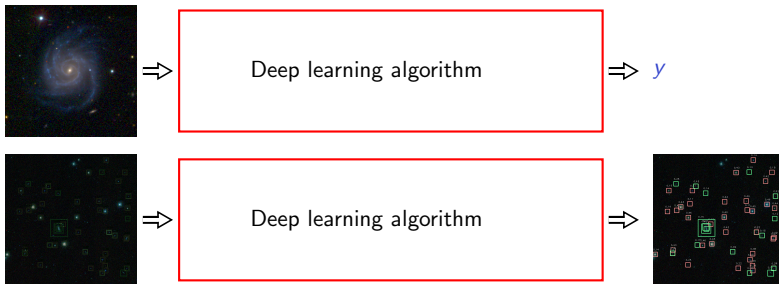
Traditional ML



Traditional ML



Deep learning: end-to-end processing





When it started: (This information is not precise)

[From Wikipedia] In 2006, publications by Geoff Hinton, Ruslan Salakhutdinov, Osindero and Teh[59] [60][61] showed how a many-layered feedforward neural network could be effectively pre-trained one layer at a time, treating each layer in turn as an unsupervised restricted Boltzmann machine, then fine-tuning it using supervised backpropagation.[62] The papers referred to learning for deep belief nets.

[59] Hinton, Geoffrey E. (2007). "Learning multiple layers of representation". Trends in Cognitive Sciences. 11(10):428–434

[60] Hinton, G. E.; Osindero, S.; Teh, Y. W. (2006). "A Fast Learning Algorithm for Deep Belief Nets". Neural Computation. 18(7):1527–1554

[61] Bengio, Yoshua (2012). "Practical recommendations for gradient-based training of deep architectures"

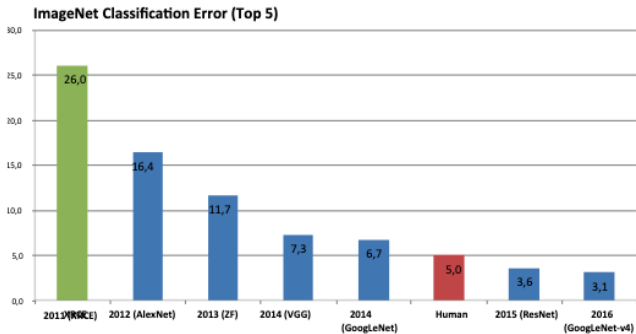
[62] G. E. Hinton., "Learning multiple layers of representation," Trends in Cognitive Sciences, 11, 428–434, 2007

Deep learning in Computer Vision

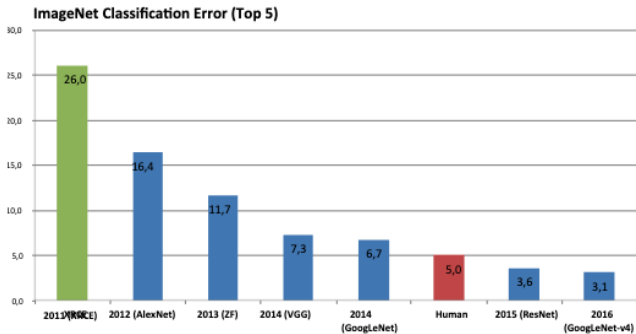
[From Wikipedia] In October 2012, ... a system by Krizhevsky et al.[5] won the large-scale ImageNet competition by a significant margin over shallow machine learning methods.

[5]. Krizhevsky, Alex; Sutskever, Ilya; Hinton, Geoffrey (2012). "ImageNet Classification with Deep Convolutional Neural Networks" (PDF). NIPS 2012: Neural Information Processing Systems, Lake Tahoe, Nevada.

Large Scale Visual Recognition Challenge (ILSVRC) - ImageNet



Large Scale Visual Recognition Challenge (ILSVRC) - ImageNet





www.image-net.org

22K categories and **14M** images

- Animals
 - Bird
 - Fish
 - Mammal
 - Invertebrate
- Plants
 - Tree
 - Flower
- Food
- Materials
- Structures
 - Artifact
 - Tools
 - Appliances
 - Structures
- Person
 - Scenes
 - Indoor
 - Geological Formations
 - Sport Activities

Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009

IMAGENET Large Scale Visual Recognition Challenge

Steel drum

The Image Classification Challenge:

1,000 object classes

1,431,167 images



Output:

Scale

T-shirt

Steel drum

Drumstick

Mud turtle



Output:

Scale

T-shirt

Giant panda

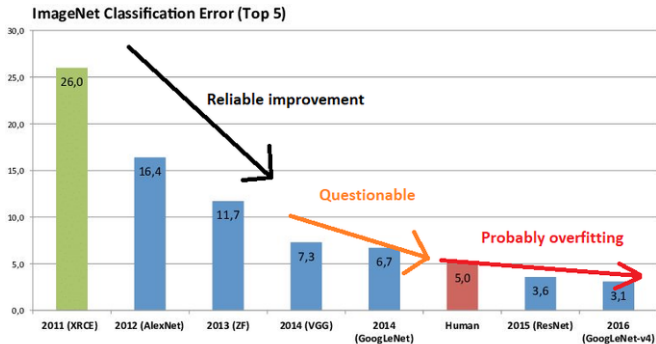
Drumstick

Mud turtle



Russakovsky et al. arXiv, 2014

Large Scale Visual Recognition Challenge (ILSVRC) - ImageNet

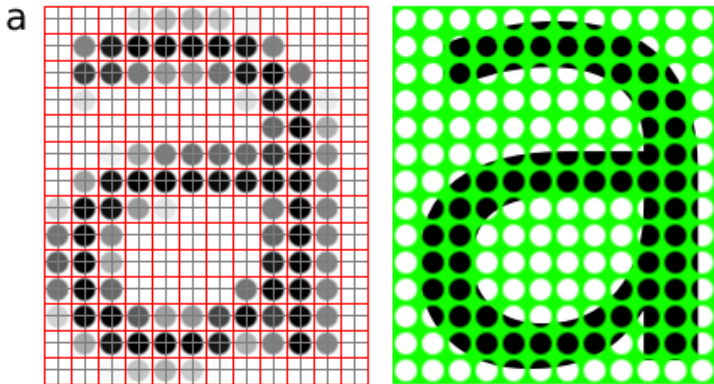


Next

1. Brief introduction to digital images
2. Examples of image processing/analysis
3. Why image processing/analysis is hard ?
4. Some traditional approaches
5. Deep learning based approaches

Brief introduction to digital images

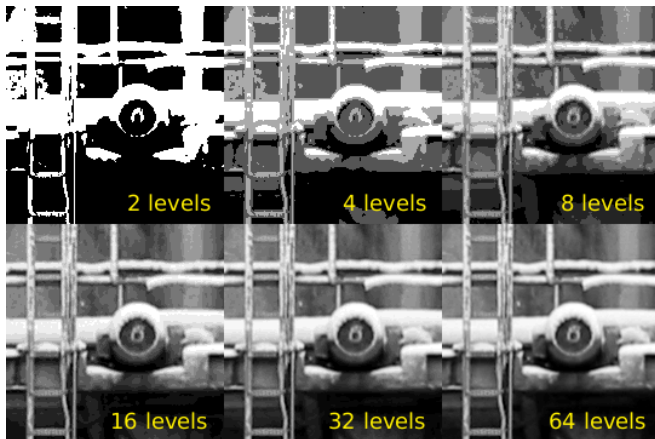
Digital images – spatial sampling



(http://pippin.gimp.org/image_processing)

Resolution: pixels per inch (ppi)

Digital images – quantization



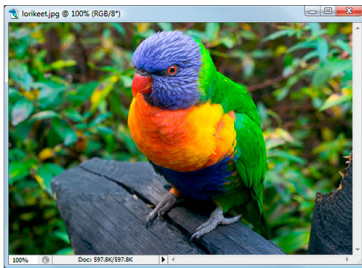
(http://pippin.gimp.org/image_processing)

Number of *bits* per pixel \rightarrow gray levels

Standard: 8 *bits* \rightarrow 256 gray levels

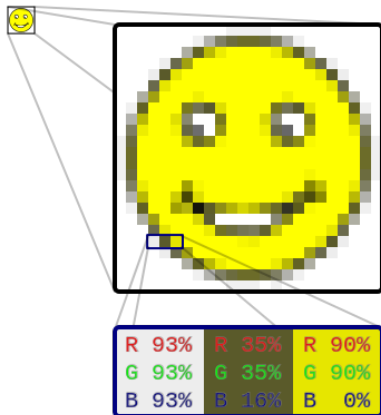
Black=0 e white=255

Digital images – Color images



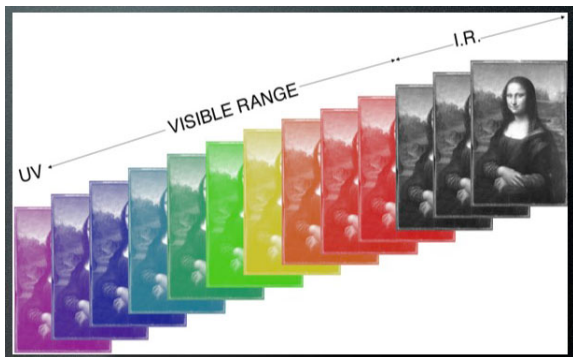
(<https://www.photoshopessentials.com/essentials/rgb/>)

Digital images – Color images



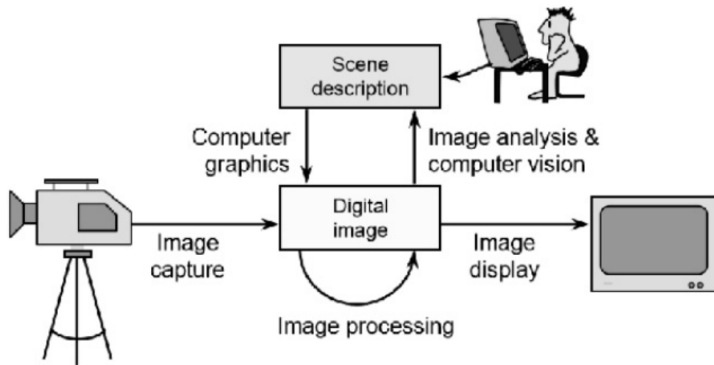
(en.wikipedia.org/wiki/Raster_graphics)

Digital images – Multispectral images



(<http://www.lumiere-technology.com/Pages/Services/services3.htm>)

Research fields related to images



(Source: slides of prof. Nidhal El Abbadi, University of Kufa)

Examples of image processing/analysis



The image shows a digital reading application interface. On the left, there is a vertical toolbar with icons for image, text, and other functions. The main window is split into two panes. The left pane shows a page from a book titled "JUVENILE READER." with the page number "37". The text on the page is highlighted in green, and there are blue boxes around certain words and phrases. The right pane shows the same text with a different highlighting scheme, including blue boxes around words and phrases. The text in both panes is as follows:

lying in bed, suffering great pain, day after day, and week after week, till he was worn quite thin, he began to recover so as to be able to hobble upon crutches.

23. After this sad misfortune, James could neither run nor play with other children, but used to sit, all day long, wishing that he had been more attentive to his mamma's request.

24. It was cause of great grief to Mrs. Cooley to see her little son in this unhappy condition; and she often regretted that she had permitted him to grow up so stubborn and obstinate.

ABSURDITY OF PRIDE.

1. EVERY man, let his state and condition in life be what they may, depends on those around him for assistance and support.

2. Men in a very low estate, may do us a great deal of good, and we often want their help. Many animals save us much labour and trouble, and supply us with many comforts.



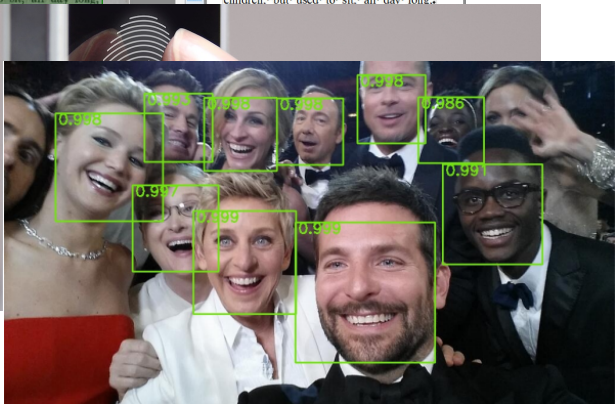
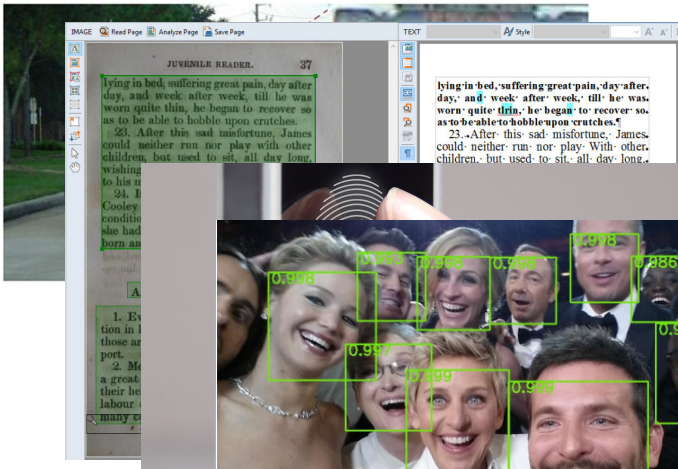
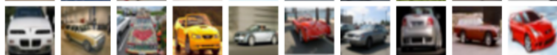


Image classification

airplane



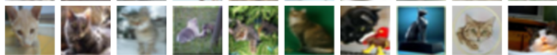
automobile



bird



cat



deer



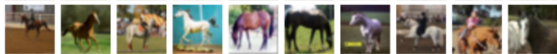
dog



frog



horse



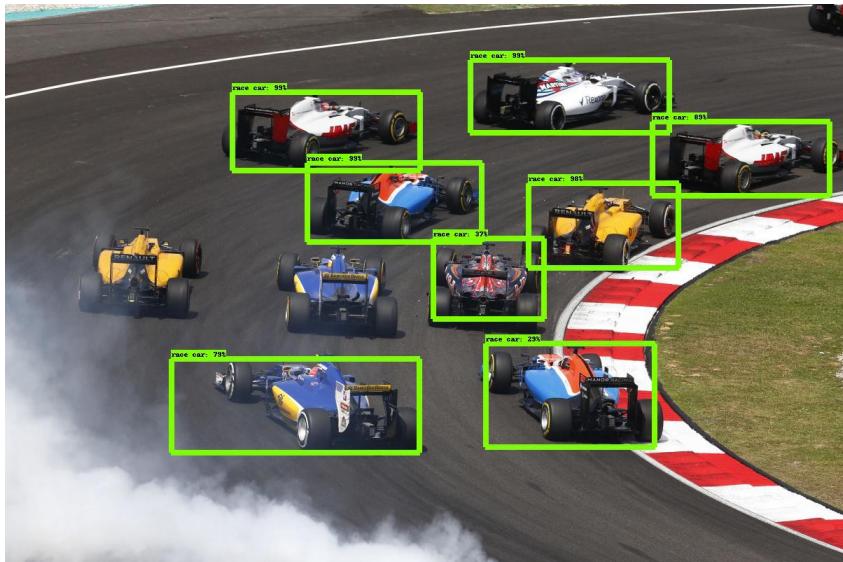
ship



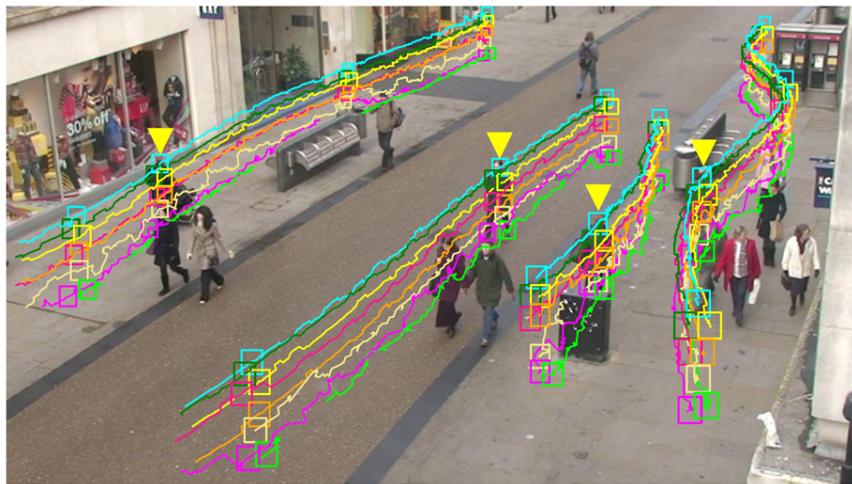
truck



Object detection



Object tracking in video



Semantic segmentation



Scene description



A female tennis player in action on the court.



A group of young men playing a game of soccer



A man riding a wave on top of a surfboard.



A baseball game in progress with the batter up to plate.



A brown bear standing on top of a lush green field.



A person holding a cell phone in their hand.

Why image processing/analysis is hard ?

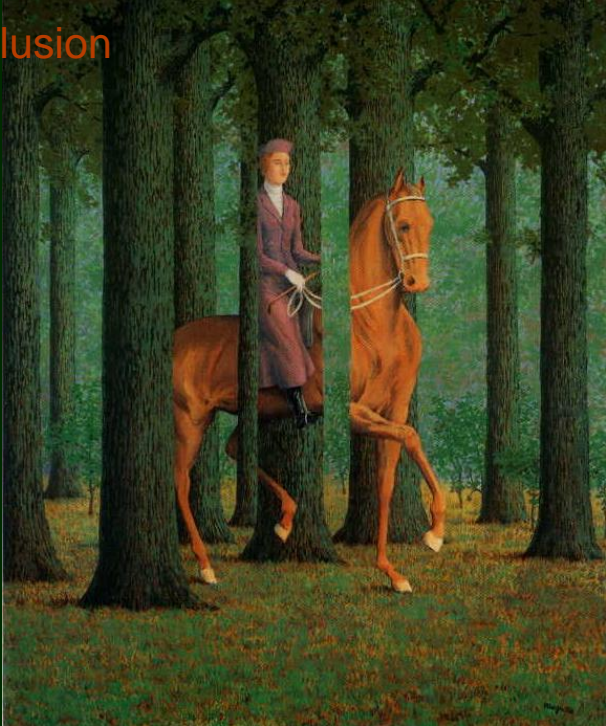
Challenges 1: view point variation



Challenges 2: illumination



Challenges 3: occlusion

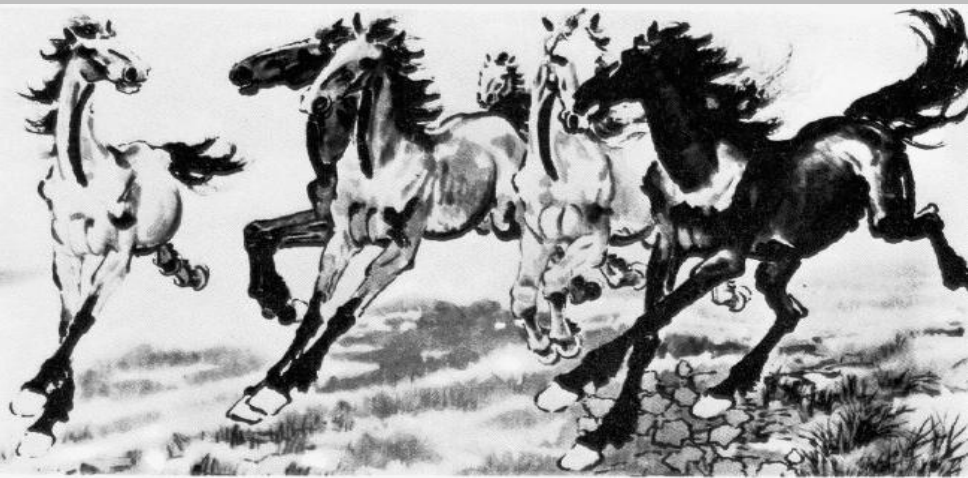


Magritte, 1957

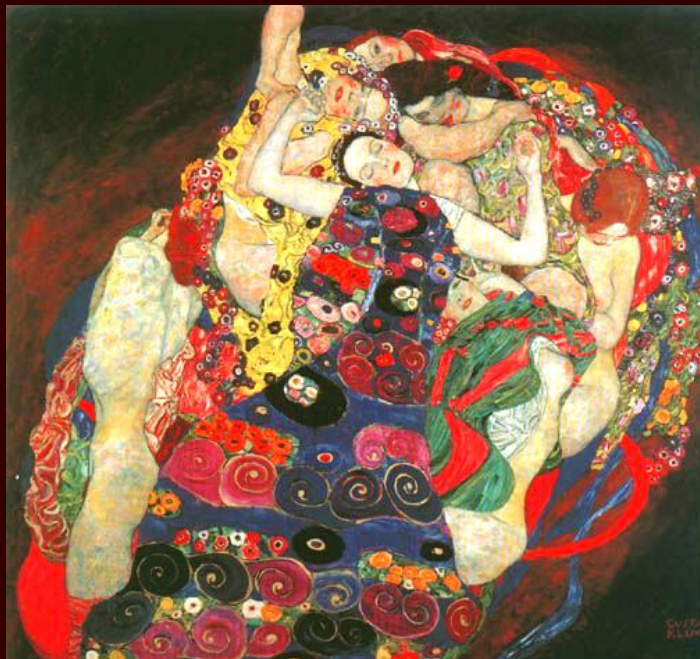
Challenges 4: scale



Challenges 5: deformation



Challenges 6: background clutter



Klimt, 1913

Within-class variations



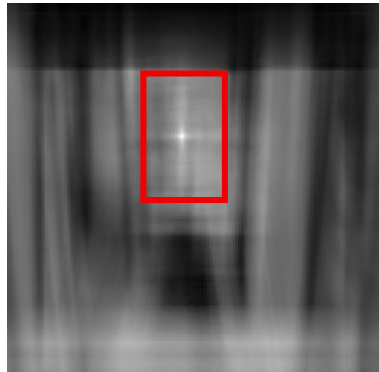
Object recognition

Is it really so hard?

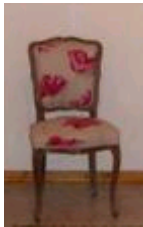
Find the chair in this image



Output of normalized correlation



This is a chair

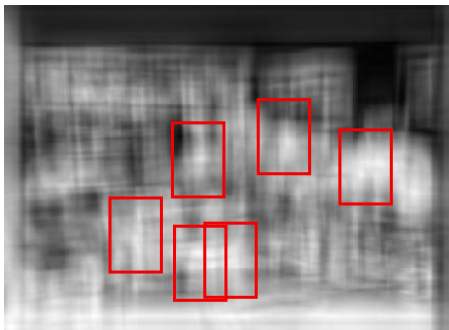
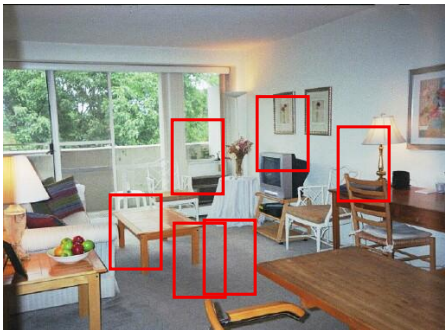




Object recognition

Is it really so hard?

Find the chair in this image



Pretty much garbage

Simple template matching is not going to make it

A “popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts.” Nivatia & Binford, 1977.

Slide: A. Torralba

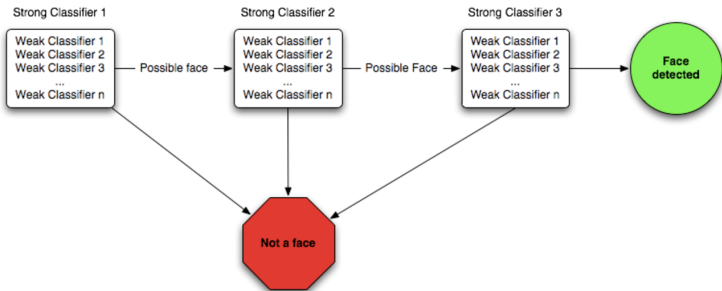
Examples of traditional approaches

- Cascade of classifiers
- SIFT
- Bag-of-visual-words

Cascade of classifiers

Rapid object detection using a boosted cascade of simple features,
P. Viola and M. Jones, CVPR 2001)

Adaboost and cascade of classifiers



(<https://www.quora.com/How-can-I-understand-Haar-like-feature-for-face-detection>)

Adaboost algorithm

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.

For $t = 1, \dots, T$:

1. Normalize the weights, $w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$
2. Select the best weak classifier with respect to the weighted error

$$\epsilon_t = \min_{f,p,\theta} \sum_i w_i |h(x_i, f, p, \theta) - y_i|.$$

See Section 3.1 for a discussion of an efficient implementation.

3. Define $h_t(x) = h(x, f_t, p_t, \theta_t)$ where $f_t, p_t,$ and θ_t are the minimizers of ϵ_t .
4. Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

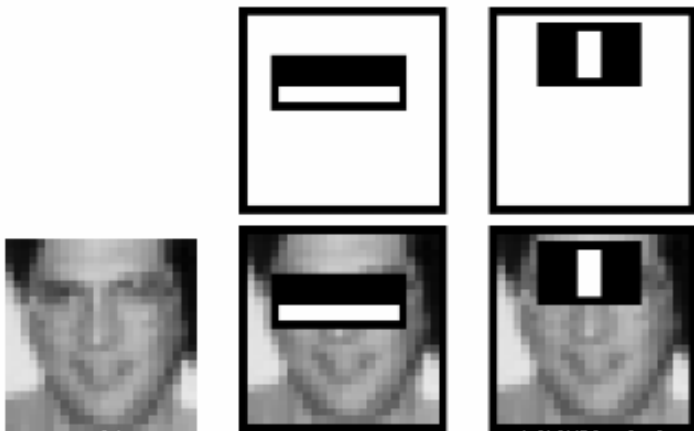
where $e_i = 0$ if example x_i is classified correctly, $e_i = 1$ otherwise, and $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$.

- The final strong classifier is:

$$C(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

where $\alpha_t = \log \frac{1}{\beta_t}$

Main Haar like features



(https://docs.opencv.org/3.4.1/d7/d8b/tutorial_py_face_detection.html)

Results



SIFT (Scale-invariant feature transform)

Registro de imagens, detecção de objetos via casamento de keypoints

Lowe, David G. (1999). *Object recognition from local scale-invariant features*. Proceedings of the International Conference on Computer Vision 2, pp. 1150–1157

Lowe, D. G. (2004). *Distinctive Image Features from Scale-Invariant Keypoints*, International Journal of Computer Vision, 60, 2, pp. 91-110

Image alignment: Applications



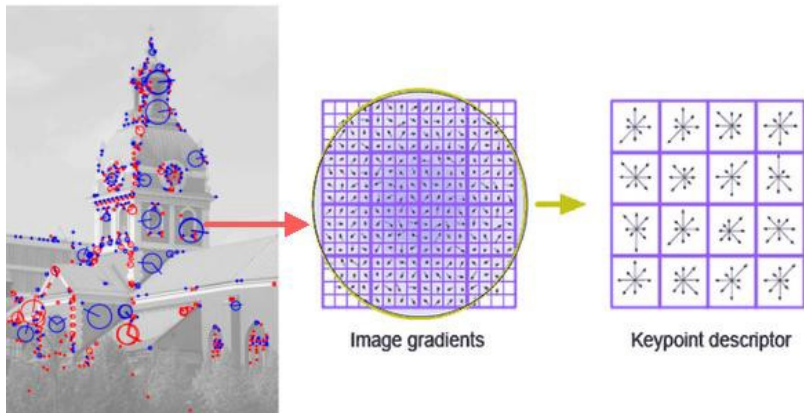
Panorama stitching



Recognition of object instances

1. **Deteção de key-points** (cantos, pontos extremos)
2. **Descrição dos key-points** (conjunto de features locais, com invariância à rotação, iluminação, etc)
3. **Correspondência entre os key-points** da query-image e da target image
Levando em conta similaridade de features
4. **Filtragem** de correspondências falsas
5. Estimar dos parâmetros de uma **transformação afim** ou **projetiva** que melhor se ajuste à correspondência

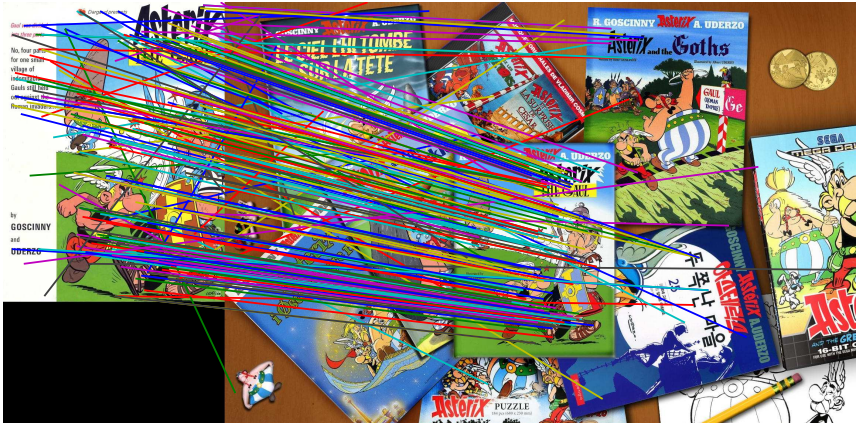
Feature description



(medium.com/machine-learning-world/

[feature-extraction-and-similar-image-search-with-opencv-for-newbies](#))

197 tentative matches



99 (50.25%) inliner matches out of 197



Mosaic

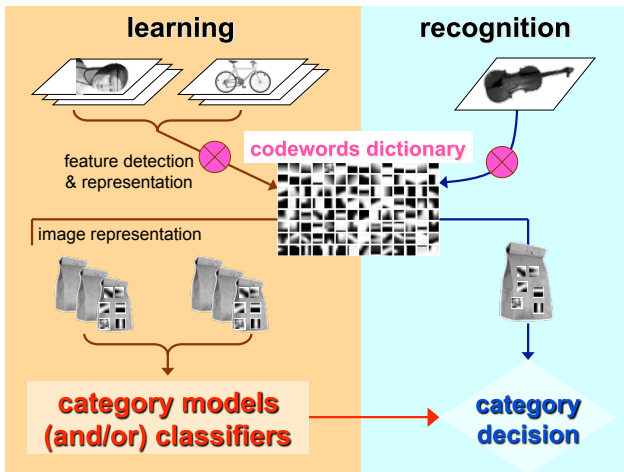


Modelo Bag of visual words (BOW)

Classificação de imagens

Image retrieval

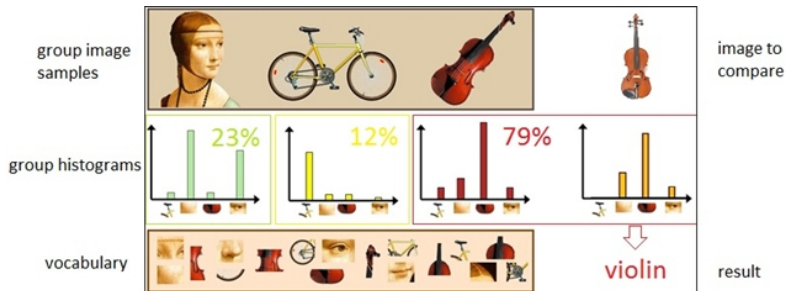
Bag-of-Words Model: Overview



Método para classificação de imagens
(qual objeto está presente na imagem?)

1. **Criar um dicionário** (padrões visuais) a partir do conjunto de imagens de treinamento
2. **Treinamento:** calcular a representação *bag-of-words* das imagens e treinar classificadores a partir dessas bags (*bag-of-words* é o conjunto de *features*)
3. **Classificação:** calcular a representação *bag-of-words* e aplicar o classificador

Modelo BOW – classificação



(<http://vgg.fiit.stuba.sk/wp-uploads/2015/02/BOW.jpg>)

Finally

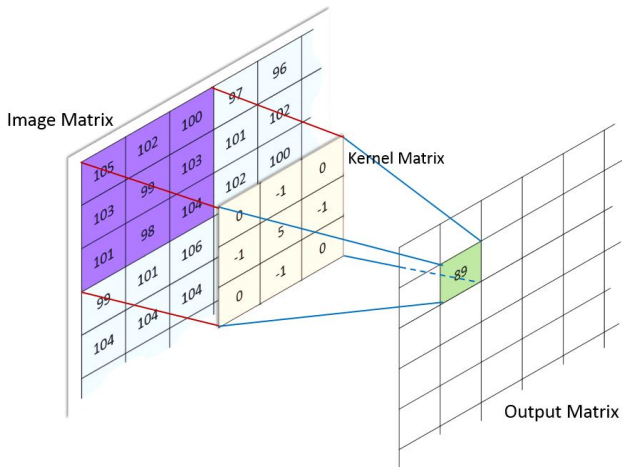
Deep convolutional networks

Convolution (or, more precisely, cross-correlation)

A basic operation used in image transformations

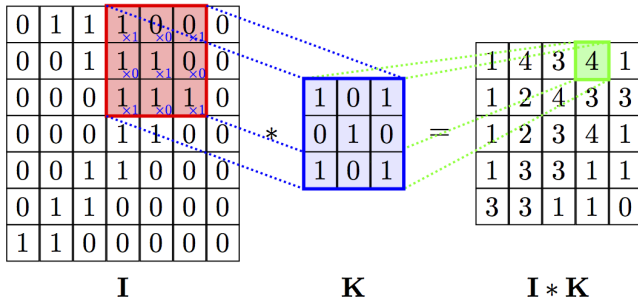
Convolution

Many image processing transformations are based on the convolution operator:



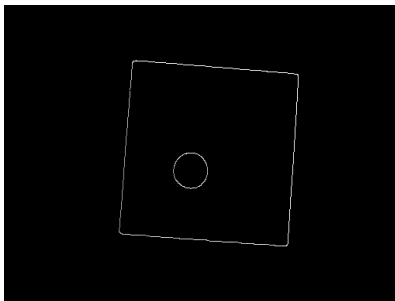
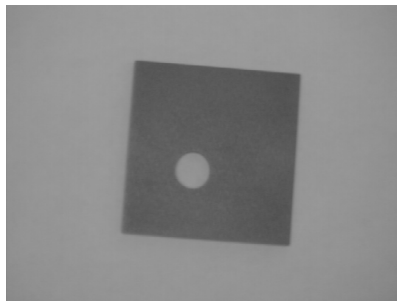
(http://machinelearninguru.com/computer_vision/basics/convolution/image_convolution_1.html)

Convolution – a second diagram



(<https://github.com/PetarV-/TikZ/tree/master/2D%20Convolution>)

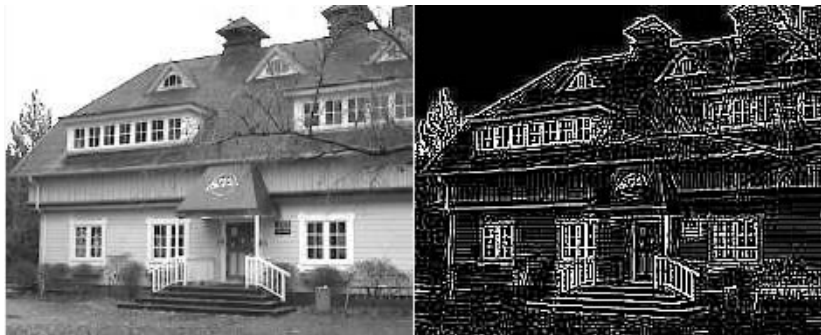
Examples of convolution kernels – edge detection



-1	-1	-1
-1	8	-1
-1	-1	-1

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

Examples of convolution kernels – edge detection



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

Smoothing



$1/9$	$1/9$	$1/9$
$1/9$	$1/9$	$1/9$
$1/9$	$1/9$	$1/9$

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

CNN (convolutional neural networks): a neural network such that

- convolution operation is part of the network
- kernels (masks) of the convolution are learned from data

Use of convolution in neural nets is not so recent:

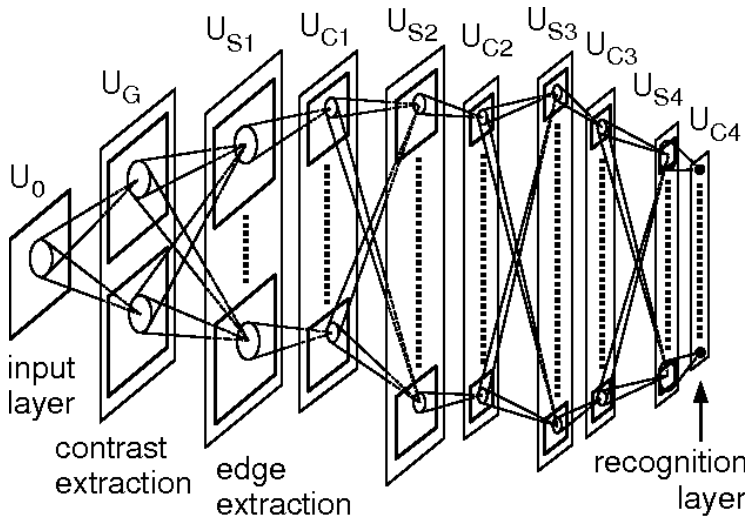
- **Neocognitron** – Fukushima, 1990

Detailed explanation: <https://www.kiv.zcu.cz/studies/predmety/uir/NS/Neocognitron/en/index.html>

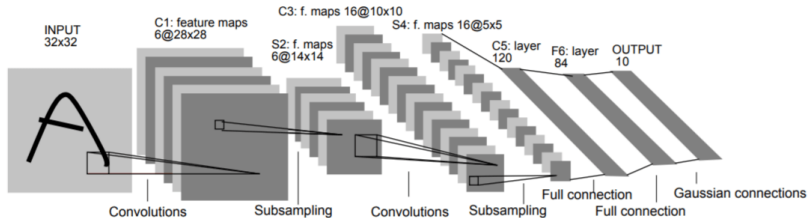
- **LeNet** - Yan Lecun, 1998, used with MNIST

<http://yann.lecun.com/exdb/lenet/index.html>

Neocognitron – Fukushima, 1990



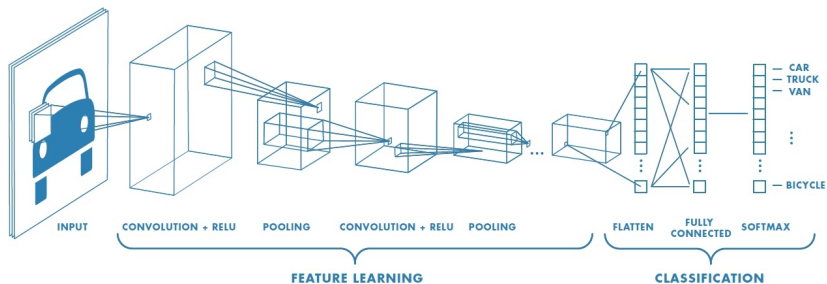
LeNet5 – Yan Lecun, 1998



This was applied on MNIST

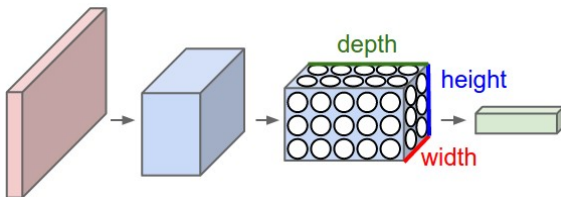
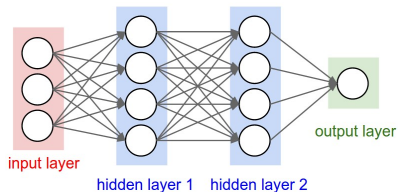
Basics of CNN

Typical architecture of Convnets



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

Graphical representation of ConvNet structure



INPUT -> [[CONV -> RELU]*N -> POOL?]*M -> [FC -> RELU]*K -> FC

Convolution on multiple channels

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+

+ 1 = -25

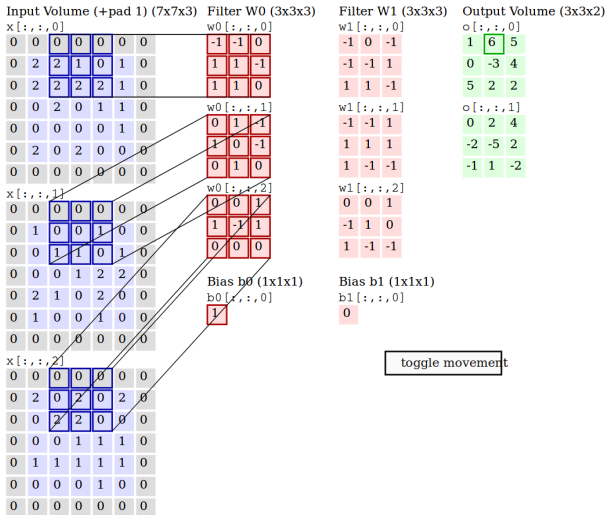
↑
Bias = 1

Output

-25			...
			...
			...
			...
...

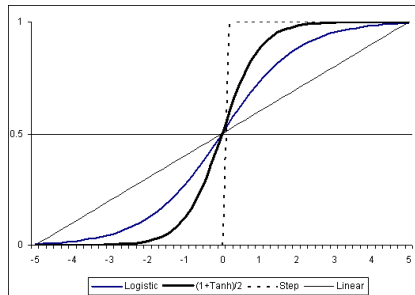
(http://machinelearningguru.com/computer_vision/basics/convolution/convolution_layer.html)

Multiple filters

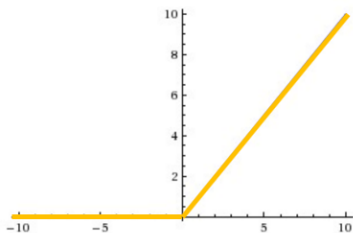


Concepts: padding, stride, kernel size, number of filters/maps

Activation functions



ReLU

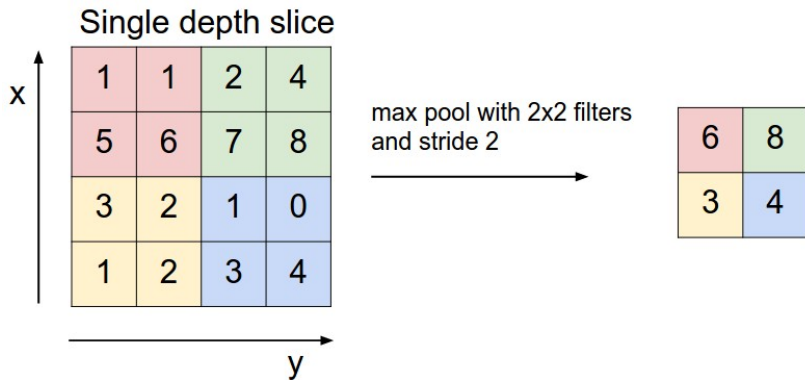


$$F(x) = \max(0, x)$$

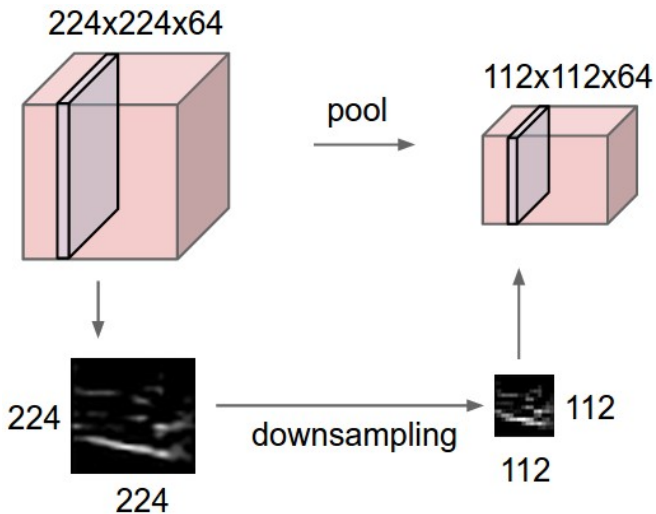
Subsampling

- via strides – skip pixels by steps of size d
- via pooling – aggregate $d \times d$ pixels

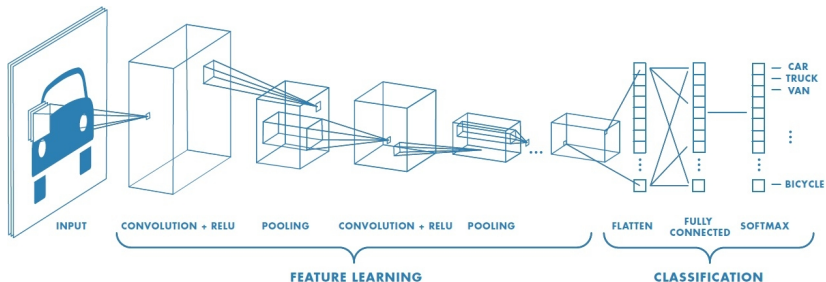
Max pooling



Pooling



Typical architecture of Convnets



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

- Convolutional layers: feature extraction
- Fully connected layers: classification

Main CNN architectures used in ILSVRC

- **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)

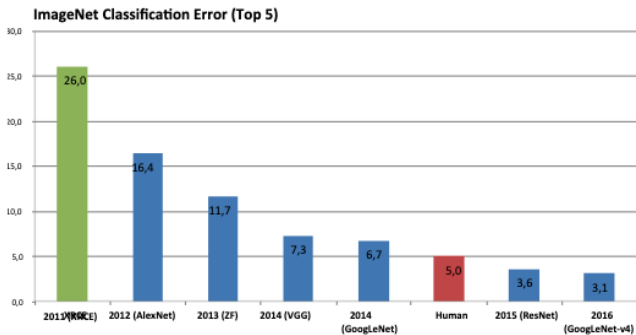
- **GoogleLeNet** (Inception), 2014: winner of ILSVRC 2014, inception layers, 7M network parameters

(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)

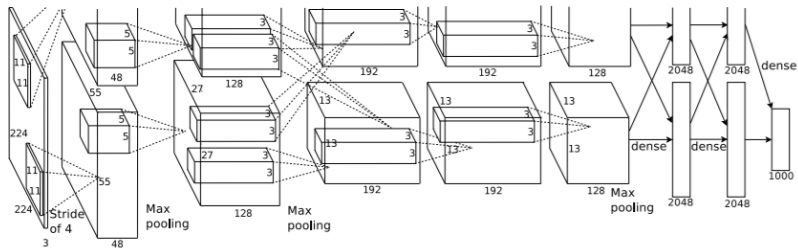
- **Residual Network** (ResNet), 2015: winner of ILSVRC 2015, 25.5M network parameters, residual block, vanishing gradient

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

Large Scale Visual Recognition Challenge (ILSVRC) - ImageNet



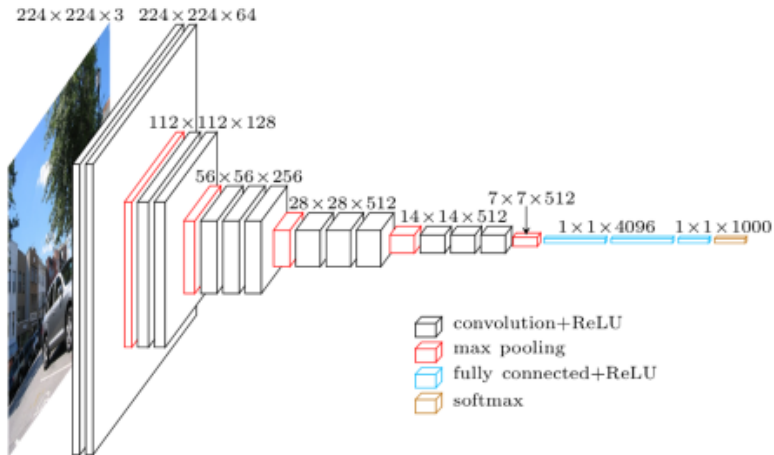
AlexNet (winner ILSVRC 2012)



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

VGGNet (2nd place ILSVRC 2014)



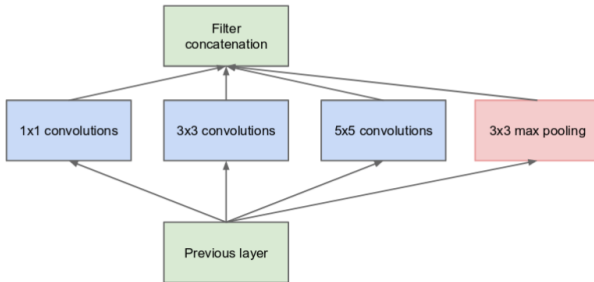
VGGNet (2nd place ILSVRC 2014)

- 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 through the layers
- 144 millions of parameters

Inception - GoogleLeNet (winner ILSVRC 2014)

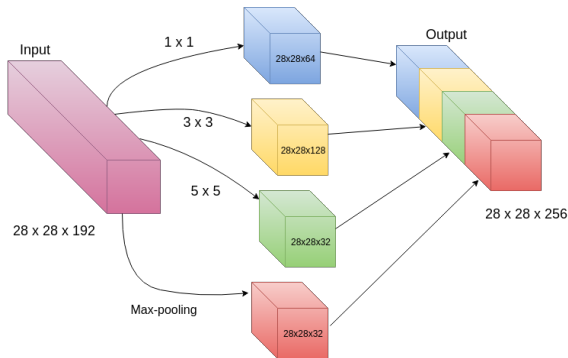


Inception – naïve version



(a) Inception module, naïve version

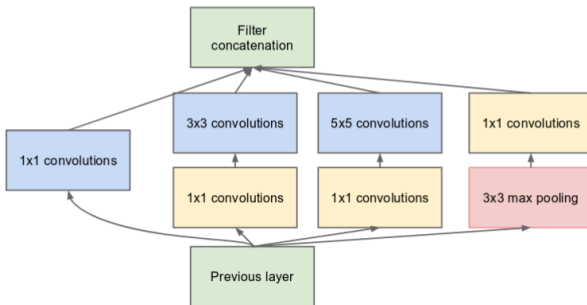
Inception – naïve version example



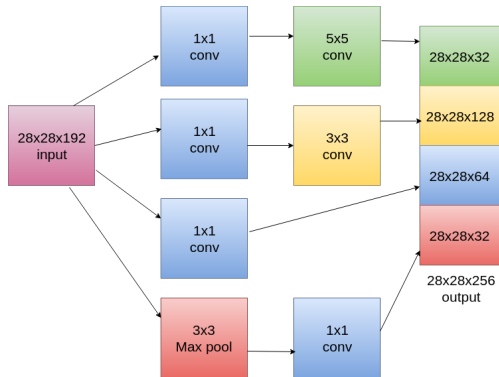
Camada de entrada = $28 \times 28 \times 192$

Convolução 5×5 : 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.

Inception – effective version



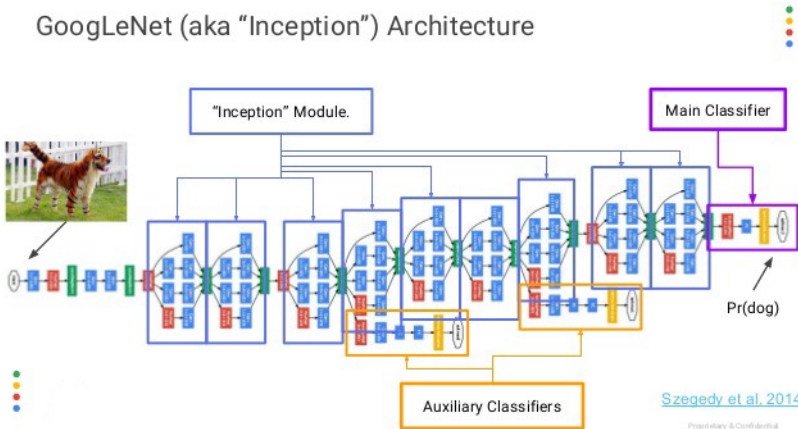
(b) Inception module with dimension reductions



Entrada $28 \times 28 \times 192 \implies 16$ filtros $1 \times 1 \implies$ camada intermediária $28 \times 28 \times 16 \implies 32$ convoluções $5 \times 5 \implies$ saída $28 \times 28 \times 32$. Total de multiplicações: cada ponto na camada intermediária requer $1 \times 1 \times 192$ multiplicações, e portanto $1 \times 1 \times 192 \times 28 \times 28 \times 16 = 2.4$ milhões para cálculo da camada intermediária; cada ponto na camada de saída requer $5 \times 5 \times 16$ multiplicações e portanto $5 \times 5 \times 16 \times 28 \times 28 \times 32 = 10$ milhões para o cálculo da camada de saída. Total = 12.4 milhões.

Inception - GoogleLeNet (winner ILSVRC 2014)

GoogLeNet (aka "Inception") Architecture

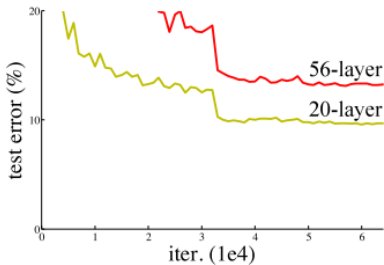
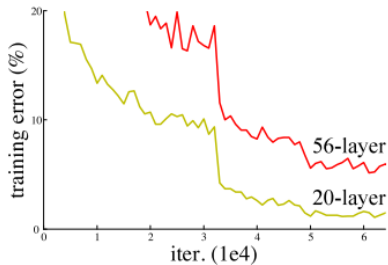


Inception - GoogleLeNet (winner ILSVRC 2014)

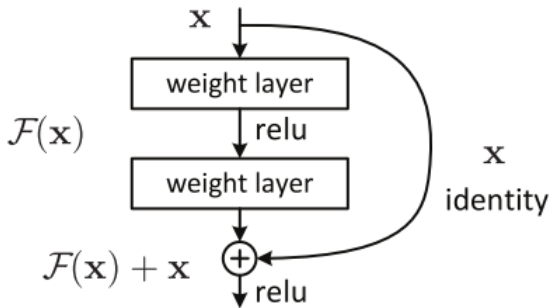
- concept of “network in network”
- inception module (“choices” made during training)
- use of 1x1 convolution
- auxiliary output to reinforce activation of intermediary layers
- more aggressive data augmentation than AlexNet
- 4 million parameters (?)
- variants (e.g., change 5x5 with two 3x3, or change 3x3 with a 1x3 and a 3x1)

ResNet (winner ILSVRC 2015)

Apesar de ser aceito que profundidade da rede é importante, observou-se que camadas adicionais degradavam a acurácia no conjunto de treinamento



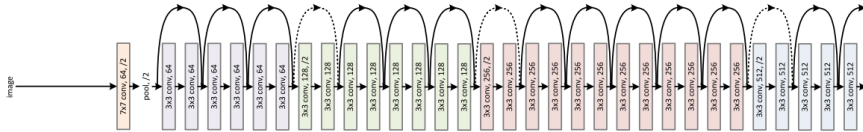
ResNet (winner ILSVRC 2015)



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 modo que $H(x) = \mathcal{F}(x) + x$ e assim a identidade corresponde a resíduo zero

ResNet (winner ILSVRC 2015)

34-layer residual



34-layer plain



VGG-19



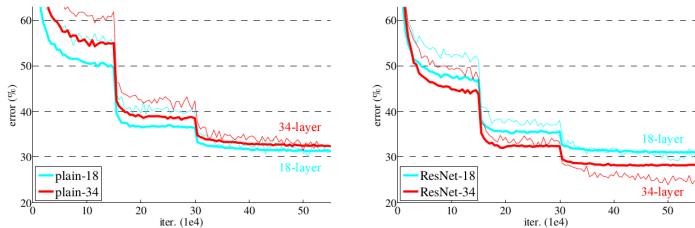


Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

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 (?)

DenseNet

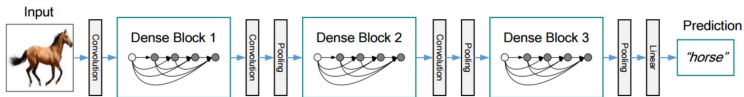
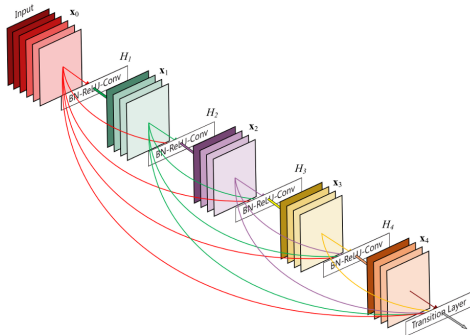
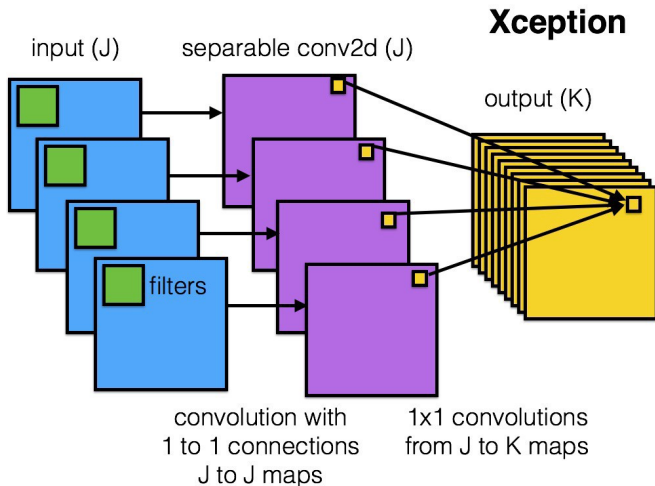


Figure 2. A deep DenseNet with three dense blocks. The layers between two adjacent blocks are referred to as transition layers and change feature map sizes via convolution and pooling.





Libraries, references, code sample

TensorFlow <https://www.tensorflow.org/>

PyTorch <https://pytorch.org/>

Keras <https://keras.io/>

Neural nets: <http://neuralnetworksanddeeplearning.com/>

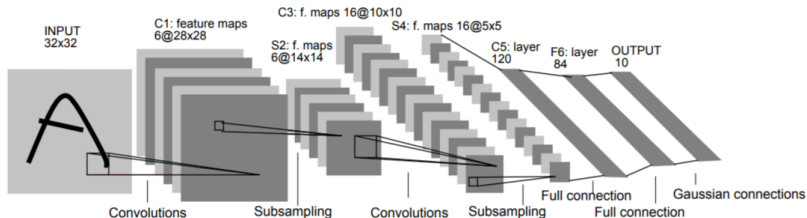
NN and CNN: Stanford CS class <http://cs231n.github.io/convolutional-networks/>

Github (F. Chollet)

<https://github.com/fchollet/deep-learning-with-python-notebooks>

LeNet-5 com Keras

LeNet-5 é a rede usada pelo LeCun no MNIST



LeNet-5 é a rede usada pelo LeCun no MNIST

Convolution: Input = $32 \times 32 \times 1$. Output = $28 \times 28 \times 6$

SubSampling: Input = $28 \times 28 \times 6$. Output = $14 \times 14 \times 6$

Convolution: Input = $14 \times 14 \times 6$. Output = $10 \times 10 \times 16$

SubSampling: Input = $10 \times 10 \times 16$. Output = $5 \times 5 \times 16$

Fully Connected: Input = $5 \times 5 \times 16$. Output = 120

Fully Connected: Input = 120. Output = 84

Output 10

```
-----  
model = keras.Sequential()  
  
model.add(layers.Conv2D(filters=6, kernel_size=(3, 3),  
                        activation='relu', input_shape=(32,32,1)))  
model.add(layers.AveragePooling2D())  
  
model.add(layers.Conv2D(filters=16, kernel_size=(3, 3),  
                        activation='relu'))  
model.add(layers.AveragePooling2D())  
  
model.add(layers.Flatten())  
  
model.add(layers.Dense(units=120, activation='relu'))  
model.add(layers.Dense(units=84, activation='relu'))  
model.add(layers.Dense(units=10, activation = 'softmax'))  
-----
```


Transfer learning

Transfer learning refers to using whatever has been learned in a certain domain to tackle similar problems in another domain

Why? Because training data is a scarce resource

Common practice: train a network with images of ImageNet, and then fine-tune the weights with data of the target domain

Comparison of traditional features vs features extracted with a network pretrained on ImageNet

Classificação de imagens de galáxias

(com Fabrício Ferrari e Mateus Espadoto)

Dataset used: EFIGI (<https://www.astromatic.net/projects/efigi>)

Procedimento executado:

1. Selecionar apenas as espirais e as elípticas
2. Extrair *features* para cada imagem
 - Usando Morfometryka (dados fornecidos por Fabrício Ferrari, FURG)
 - Usando uma rede convolucional pré-treinada com o ImageNet
3. Separar dados em treinamento e validação
4. Treinar classificador logístico usando dados de treinamento
5. Avaliar desempenho (de predição) no conjunto de validação

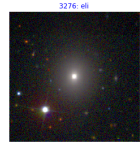
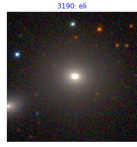
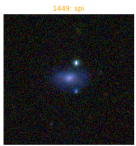
Para inspecionar visualmente a qualidade das *features* extraídas, nos slides seguintes:

- rótulo sobre a imagem é a classificação correta (azul="elíptica", laranja="espiral")
- primeira coluna: exemplos de imagens classificadas erroneamente
- colunas 2 a 5: imagens mais similares ao da primeira coluna de acordo com a similaridade de *features*

Features convolucionais

Classificação errada

4 vizinhos mais próximos



Features convolucionais

Classificação errada

4 vizinhos mais próximos



Features morfológicos

Classificação errada

4 vizinhos mais próximos



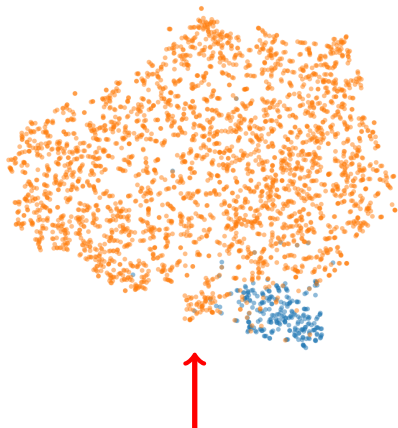
Features morfológicos

Classificação errada

4 vizinhos mais próximos

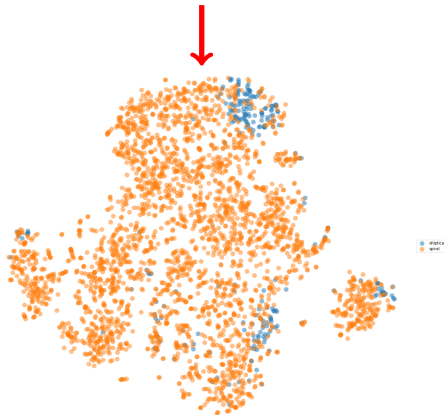


Projeções das features



Features convolucionais

Features morfométricos



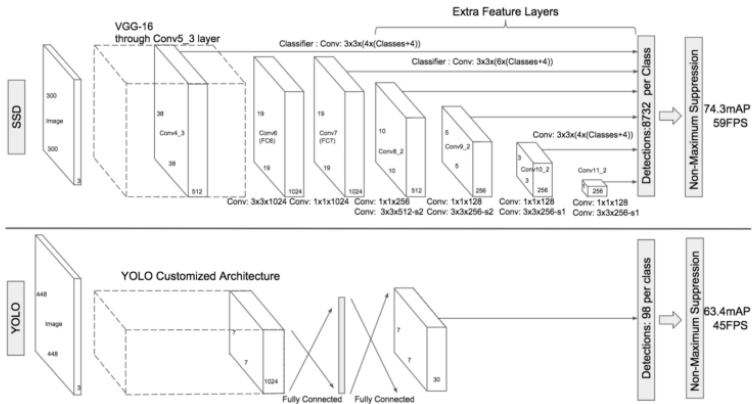
Visualization

- of the network itself
- parts of the image that affect output
- parts of the image that affect a specific neuron

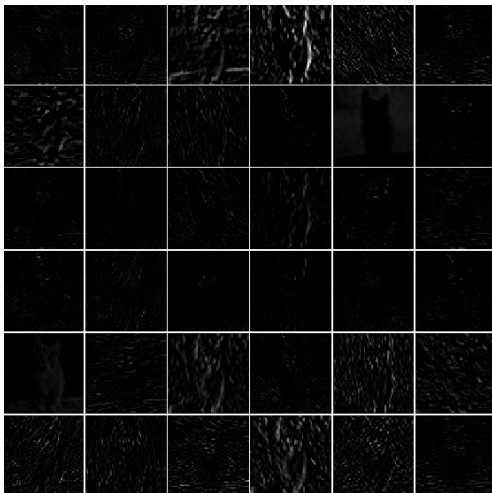
Visualizar a rede – Arquitetura

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	(None, 1000)	0
embedding_1 (Embedding)	(None, 1000, 100)	17410600
conv1d_1 (Conv1D)	(None, 996, 128)	64128
max_pooling1d_1 (MaxPooling1D)	(None, 199, 128)	0
conv1d_2 (Conv1D)	(None, 195, 128)	82048
max_pooling1d_2 (MaxPooling1D)	(None, 39, 128)	0
conv1d_3 (Conv1D)	(None, 35, 128)	82048
max_pooling1d_3 (MaxPooling1D)	(None, 1, 128)	0
global_max_pooling1d_1 (GlobalMaxPooling1D)	(None, 128)	0
dense_1 (Dense)	(None, 128)	16512
dense_2 (Dense)	(None, 20)	2580
=====		
Total params: 17,657,916		
Trainable params: 247,316		
Non-trainable params: 17,410,600		
=====		
None		

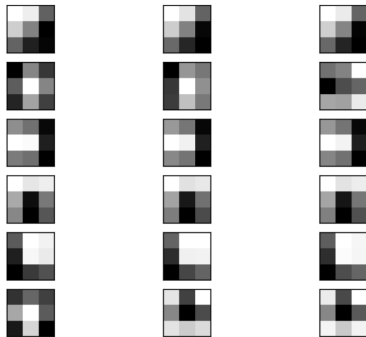
Visualizar a rede – Arquitetura



Visualizar a rede – mapa de ativação/feature map



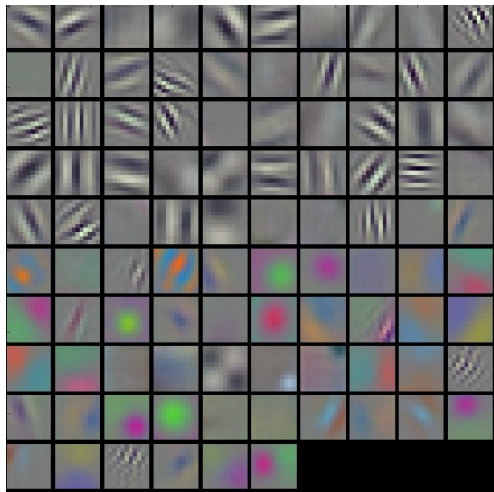
Visualizar a rede – FilTROS



<https://machinelearningmastery.com/how-to-visualize-filters-and-feature-maps-in-convolutional-neural-networks/>

6 filters (each row), from the first layer of VGG

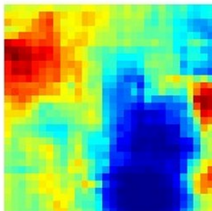
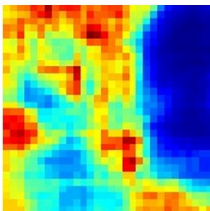
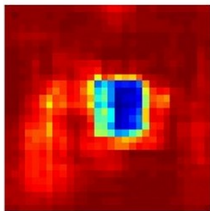
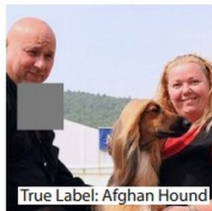
Visualizar a rede – Filtros



Images that maximize the response of the filters (AlexNet)

Qual parte da imagem afeta a classificação?

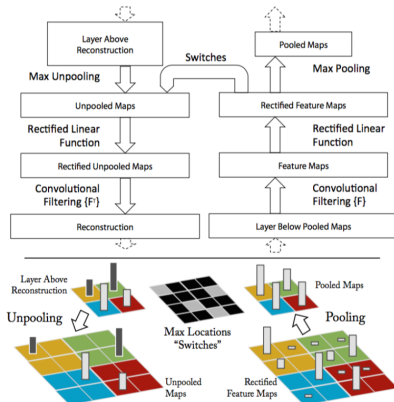
(Zeiler et al) Zerar uma região da imagem e ver como isso afeta o score de classificação.



Azul: região que mais afeta a classificação

Quais elementos da imagem ativam um neurônio ?

(Zeiler et al.) Passa-se uma imagem pela rede;
Zera-se todos os neurônios na camada, exceto aquele de interesse;
Retropropaga-se o valor do neurônio até a entrada (reconstrói-se a imagem)



Quais elementos da imagem ativam um neurônio ?

Feature map (primeira camada convolucional)



Imagens reconstruídas dos neurônios 1 a 32



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

Quais elementos da imagem ativam um neurônio ?

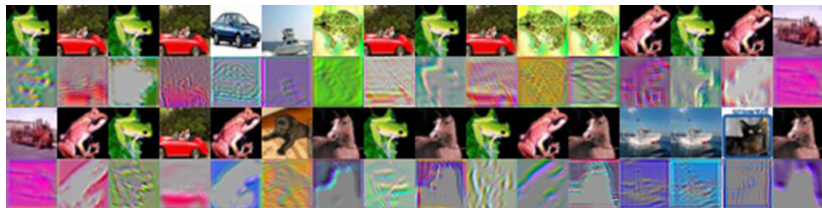
Imagens reconstruídas do neurônio 7, para diferentes imagens



Cor verde ?

Quais elementos da imagem ativam um neurônio ?

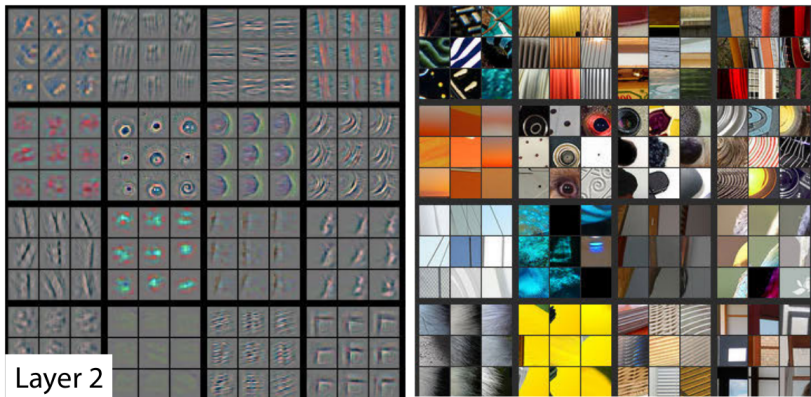
Para cada neurônio, reconstrução da imagem que provocou ativação máxima



Alguns são ativados pela presença de cor apenas.

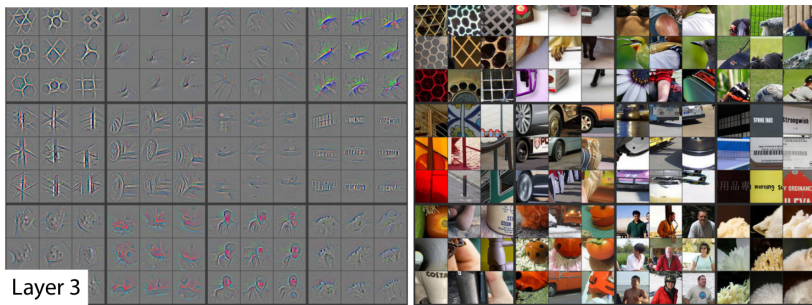
Quais elementos da imagem ativam um neurônio ?

(Zeiler et al.) 16 neurônios quaisquer na camada 2; para cada neurônio, reconstrução das 9 imagens que provocaram as maiores ativações



Quais elementos da imagem ativam um neurônio ?

12 neurônios quaisquer na camada 3; para cada neurônio, reconstrução das 9 imagens que provocaram as maiores ativações



Beyond classification

CNNs desenvolvidas para classificação de imagens passaram a servir como base para outras tarefas

CNN para object detection/recognition

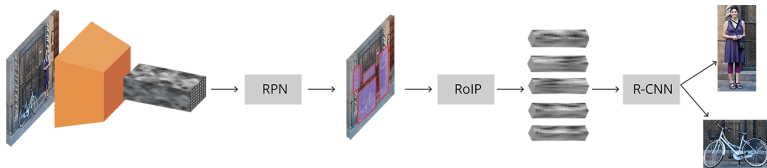
<https://towardsdatascience.com/deep-learning-for-object-detection-a-comprehensive-review-73930816d8d9>

O conjunto de dados de treinamento consiste de imagens mais uma lista de *bounding box* dos objetos com os respectivos rótulos de classe.

Função objetivo é uma combinação de classificação (acertar rótulos) e regressão (acertar tamanho e localização dos *boxes*).

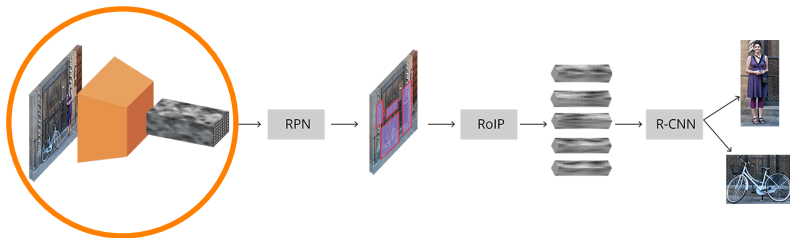
Existem os modelos two-step (ex., Faster R-CNN) e *single shot* (SSD ou YOLO)

Faster R-CNN architecture



<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

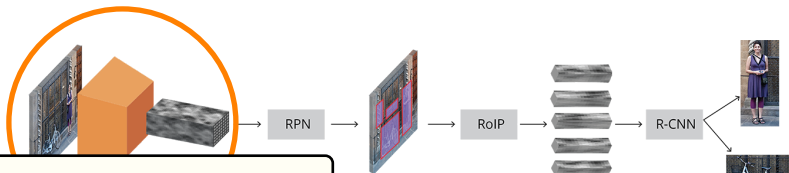
Faster R-CNN architecture



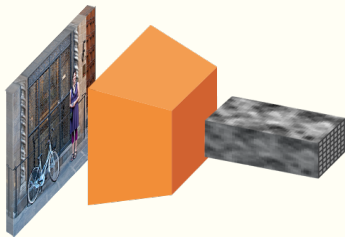
Feature extraction

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture



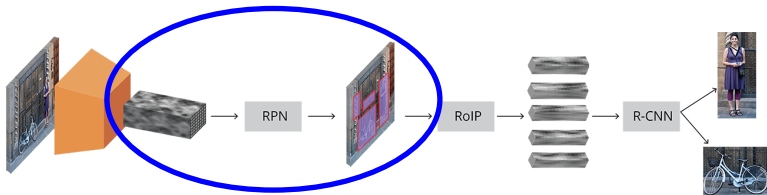
Feature extraction



- Receives an image as input
- Feature extraction part (convolutional and pooling layers) of existing CNNs (VGG, Resnet, ...)

cnn-combines-fast-rpn-for-modern-object-detection

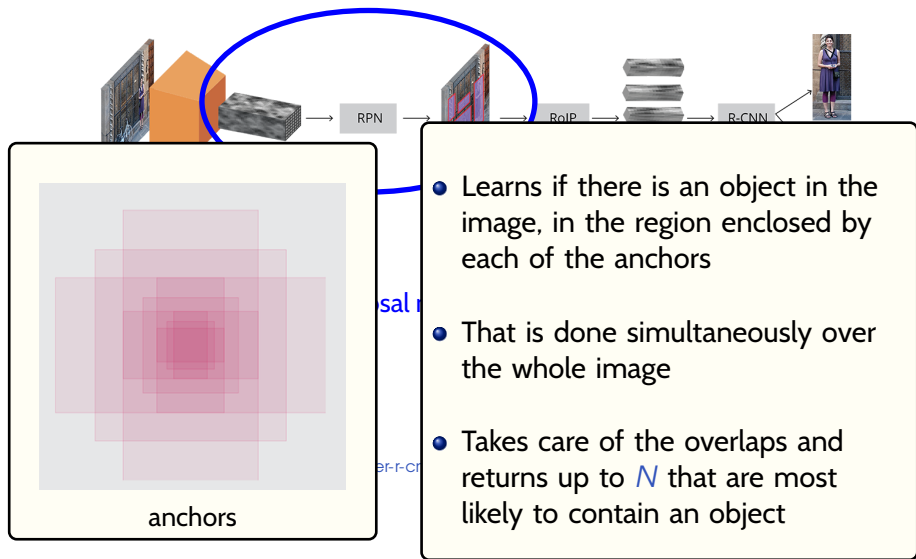
Faster R-CNN architecture



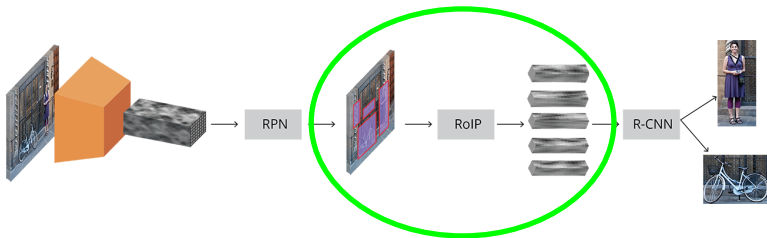
Region proposal network

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture



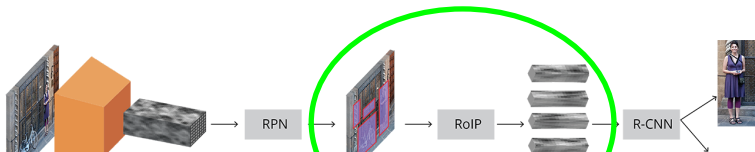
Faster R-CNN architecture



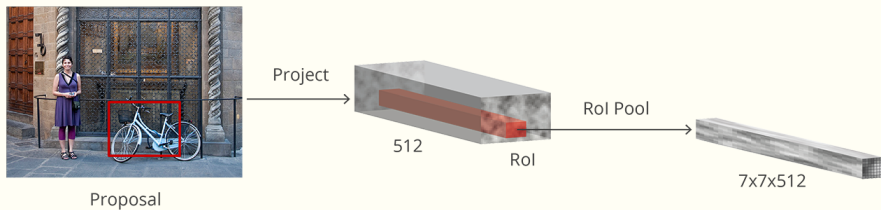
Region of interest pooling

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

Faster R-CNN architecture

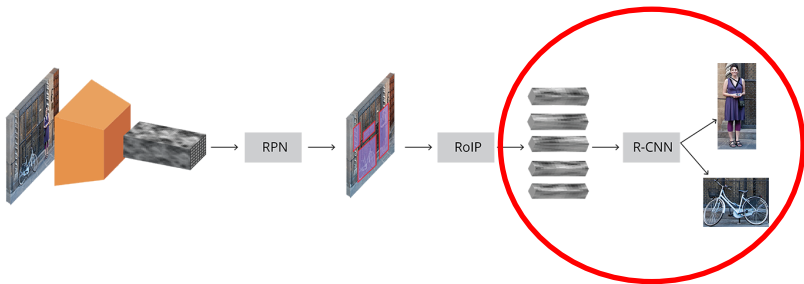


Region of interest pooling



- For each region proposal, crops the corresponding cells from the feature map
 - pools it to a fixed size

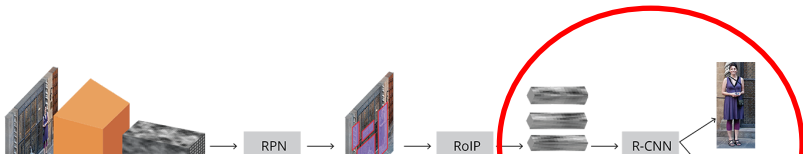
Faster R-CNN architecture



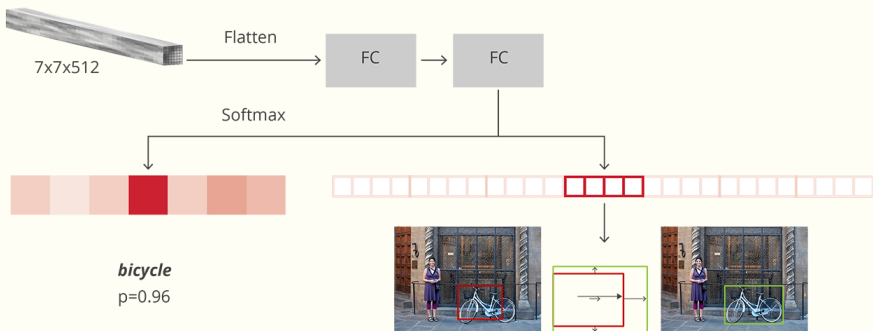
Region classification

<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

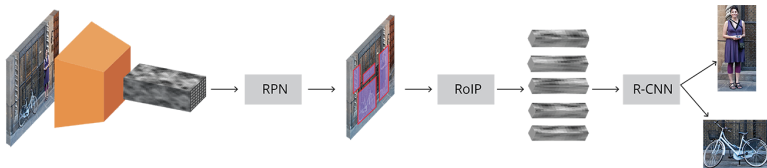
Faster R-CNN architecture



Region classification



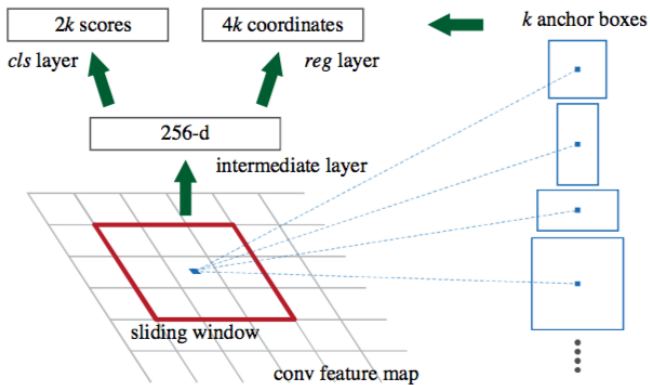
Faster R-CNN architecture

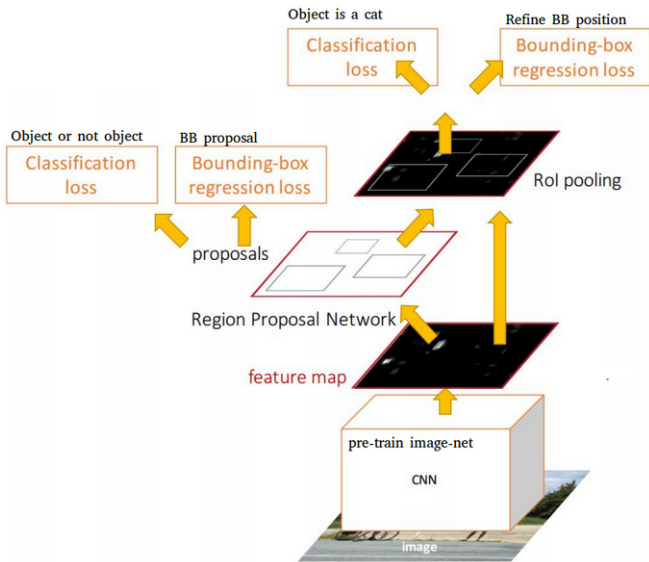


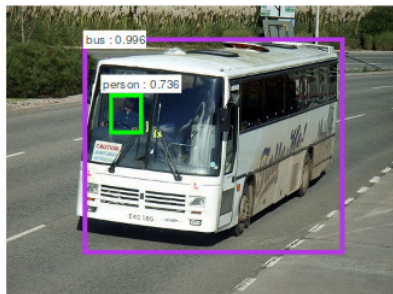
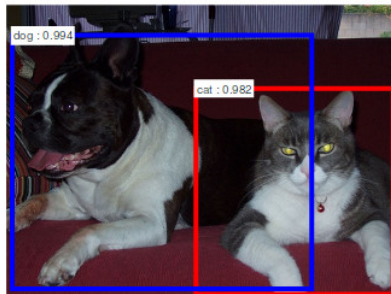
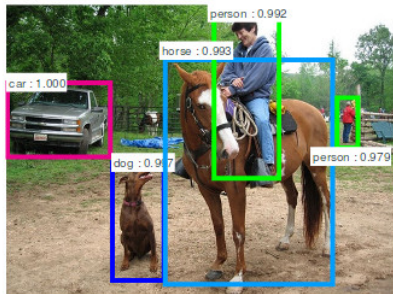
Input is an image

Output is a list of detected objects (bounding box and class probabilities)

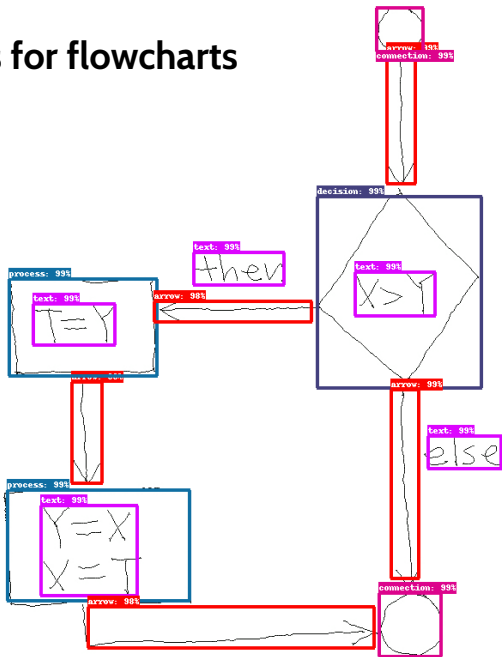
<https://tryolabs.com/blog/2018/01/18/faster-r-cnn-down-the-rabbit-hole-of-modern-object-detection/>

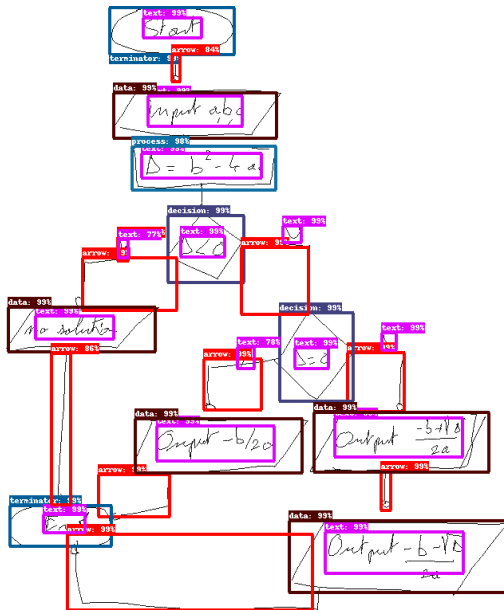


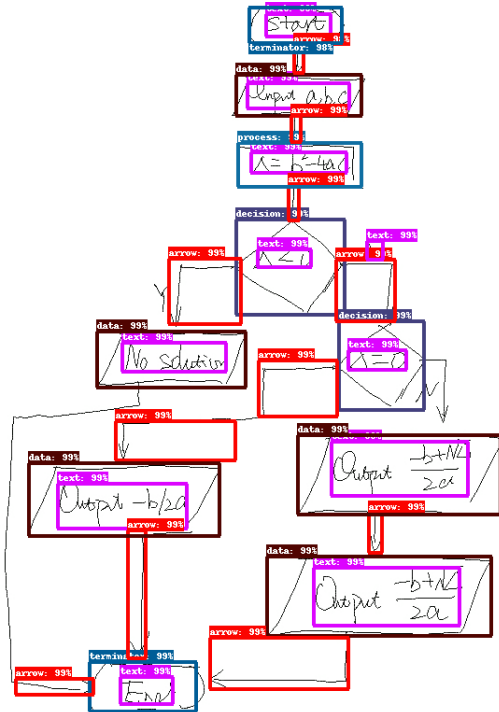


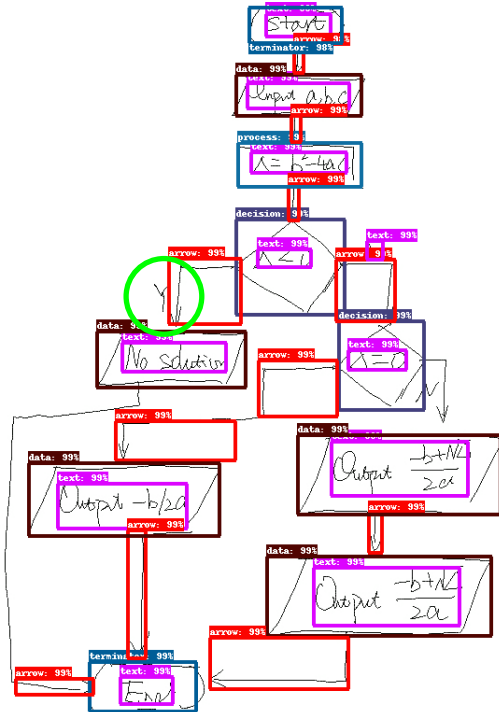


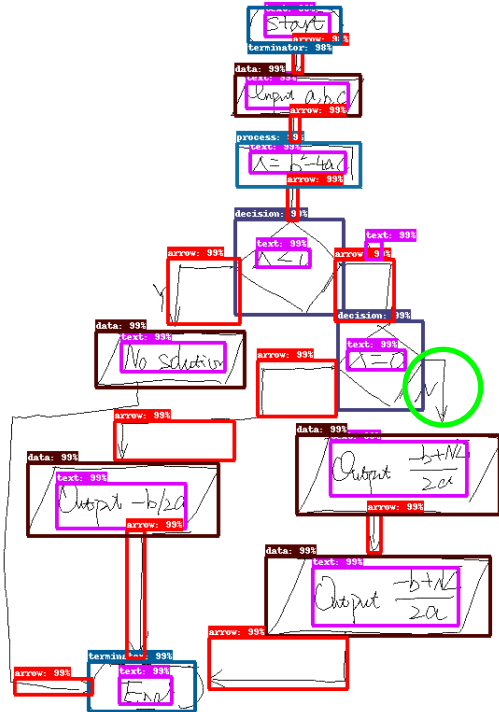
Results for flowcharts

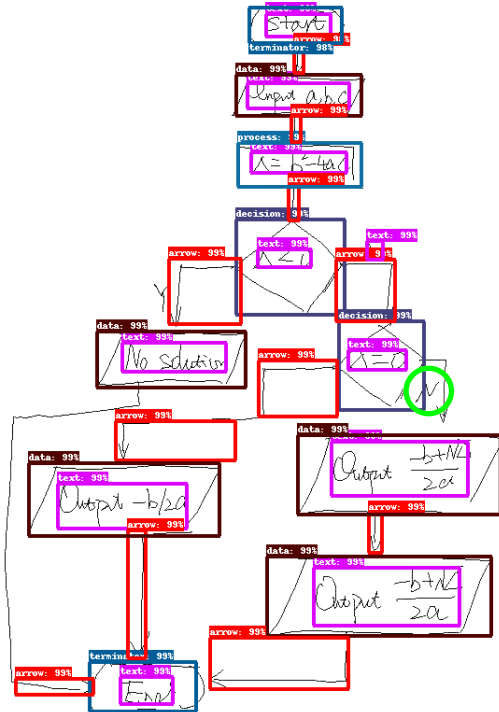


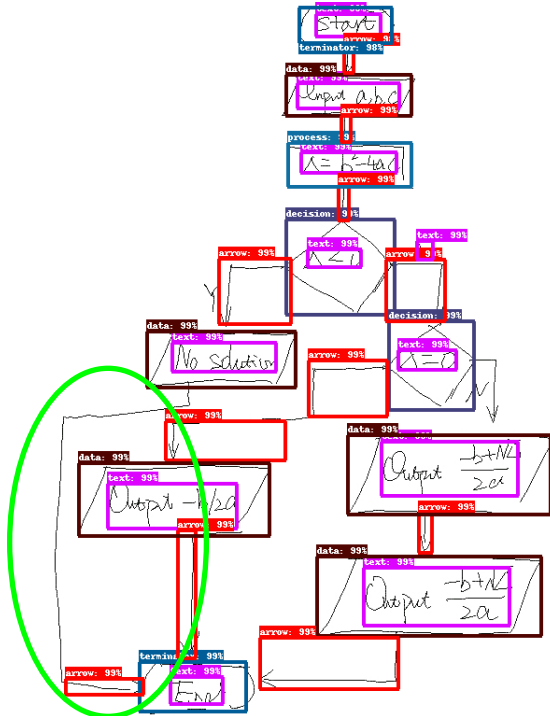




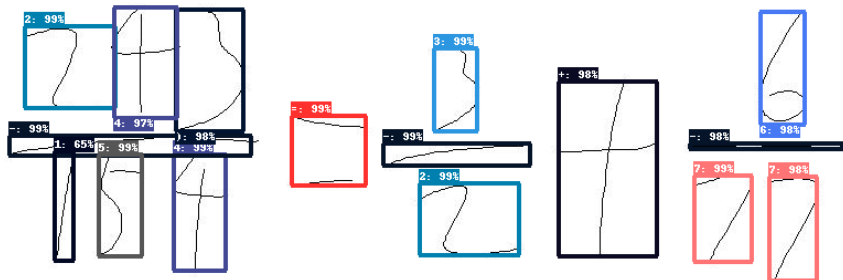








Results for math expressions



\sin: 99%

SIN

b:

98%

+: 99%

+

\sin: 99%

SIN

C:

96%

=: 99%

=

\cos: 99%

COS

p: 72%

p

+: 99%

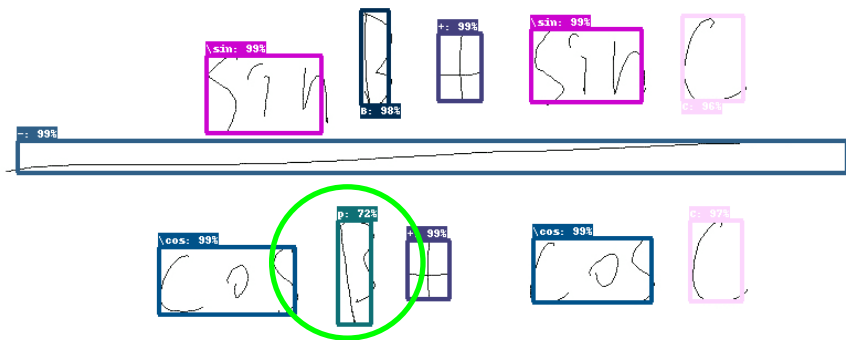
+

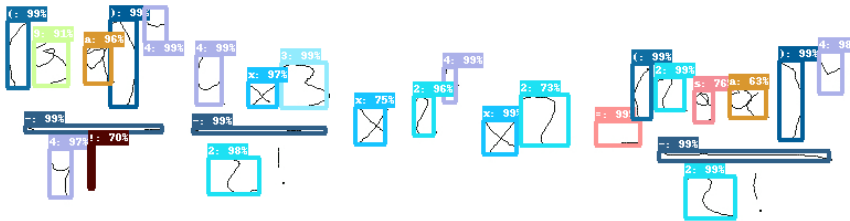
\cos: 99%

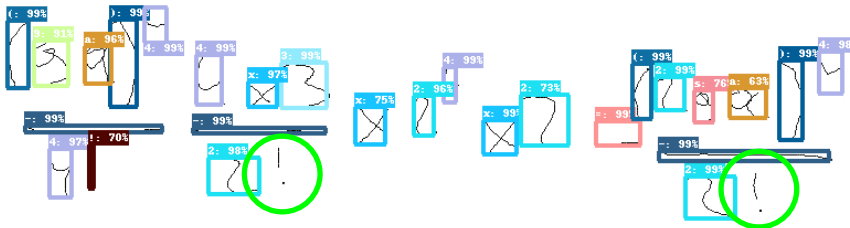
COS

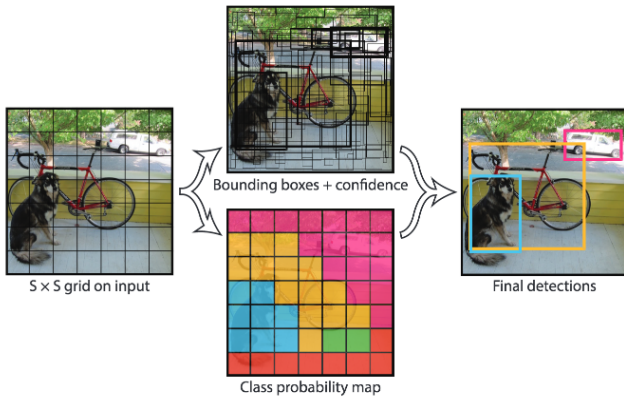
C: 97%

C









CNN para image segmentation

Mais: <http://blog.qure.ai/notes/semantic-segmentation-deep-learning-review>



Objetivo: Delimitar o contorno dos objetos/componentes de interesse presentes na imagem (ou atribuir um rótulo único para todos os *pixels* de um objeto)

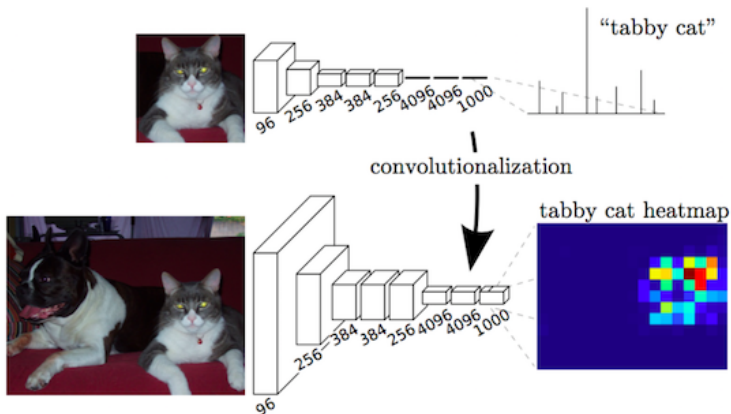
A arquitetura das redes segue a forma *encoder-decoder*

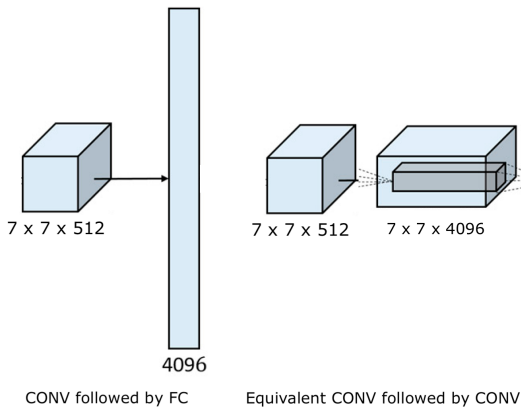
A parte *encoder* em geral é a parte convolucional de uma CNN.

A parte *decoder* busca restaurar a resolução inicial.

FCN – Fully convolutional network

First *end-to-end* network for semantic segmentation





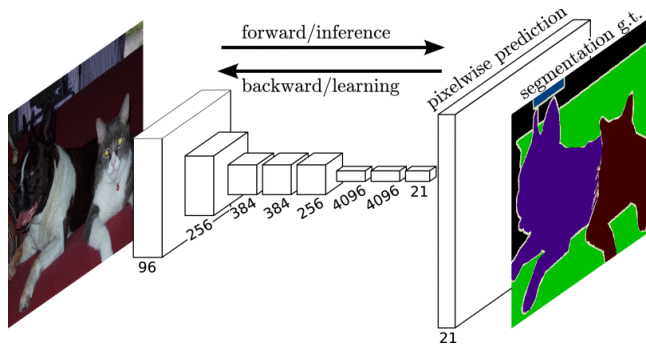
<http://cs231n.github.io/convolutional-networks/#convert>

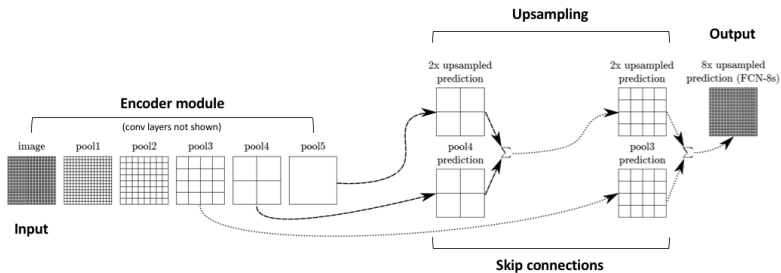
Example: Last feature map $7 \times 7 \times 512$ + two fully connected layers with 1024 nodes + output layer with 1000 nodes

Replace the first FC layer that looks at $[7 \times 7 \times 512]$ volume with a CONV layer that uses filter size $F=7$, giving output volume $[1 \times 1 \times 4096]$.

Replace the second FC layer with a CONV layer that uses filter size $F=1$, giving output volume $[1 \times 1 \times 4096]$.

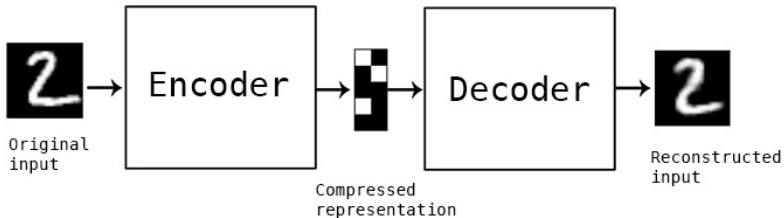
Replace the last FC layer similarly, with $F=1$, giving final output $[1 \times 1 \times 1000]$





Encoder-decoder schema

A maioria das redes imagem-a-imagem segue a estrutura encoder-decoder, que é um caso geral do autoencoder



U-Net

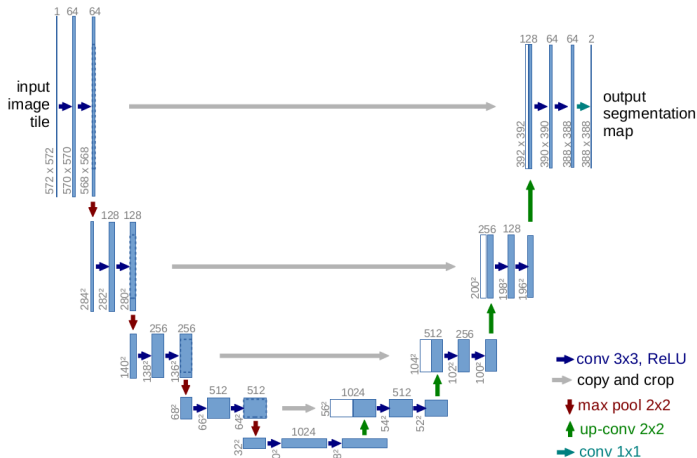
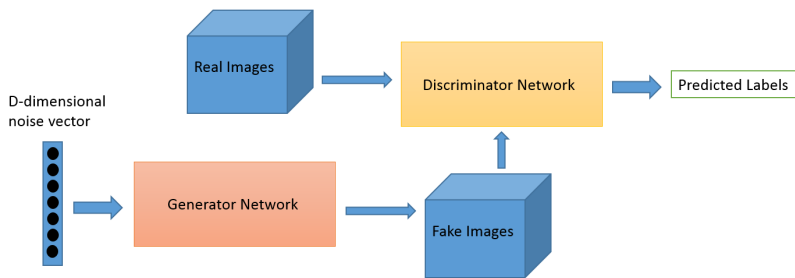


Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

- **GAN** – gerar imagens
- **Siamese network, triplet network** — image retrieval
- **RNN** – dados temporais/sequenciais (recorrência)
- **LSTM** – recorrência
- **Geometric Deep Learning** deals with the extension of Deep Learning techniques to graph/manifold structured data.
<http://geometricdeeplearning.com/>

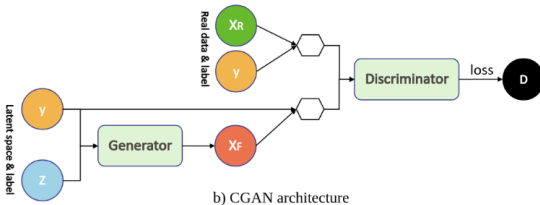
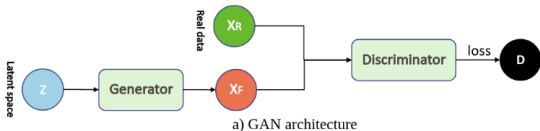
GAN schema



(fonte: O'Reilly)

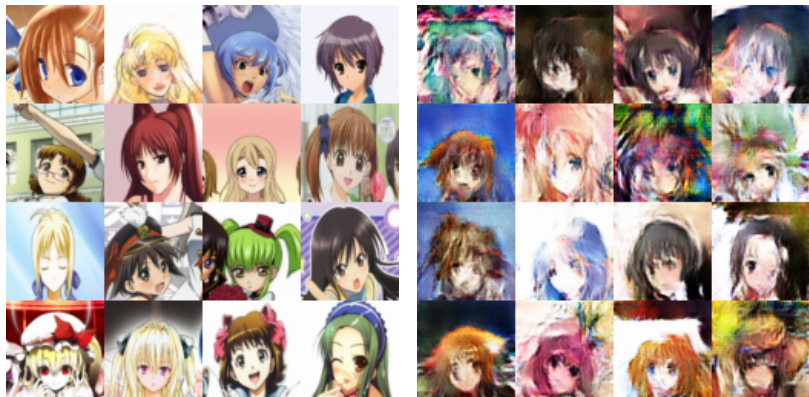
$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + E_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + E_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$



<https://mc.ai/a-tutorial-on-conditional-generative-adversarial-nets-keras-implementation/>

GAN



(<https://github.com/nashory/gans-awesome-applications>)

Você já viu esses rostos ?



es rostos ?



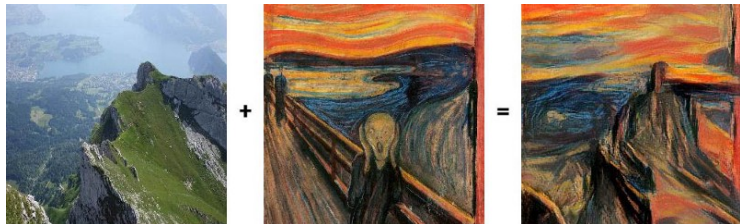






Esses rostos não existem. Foram gerados por computador!
<https://thispersondoesnotexist.com/>

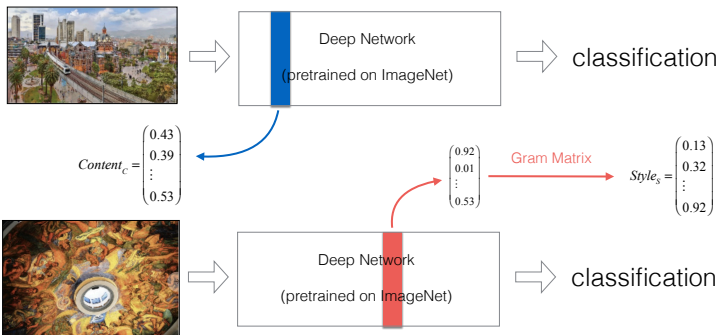
Artistic Style transfer



Artistic Style Transfer with Convolutional Neural Network

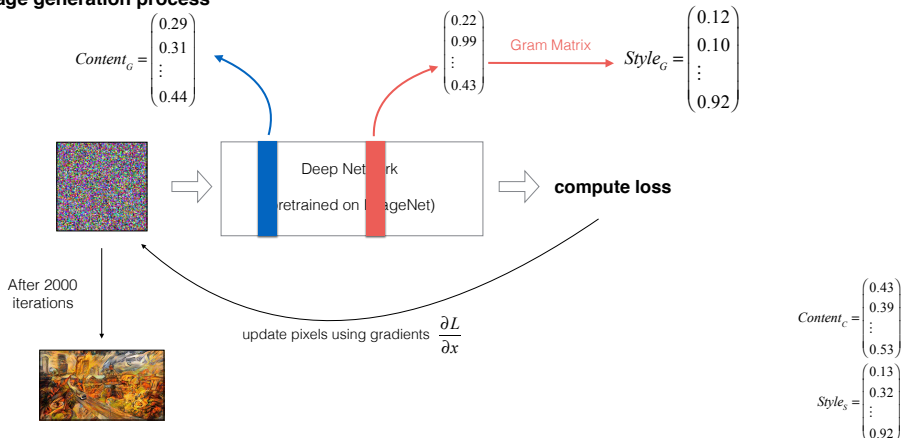
Case study 3: Art Generation

4. Architecture? We use a **pre-trained model** because it **extracts important information** from images.

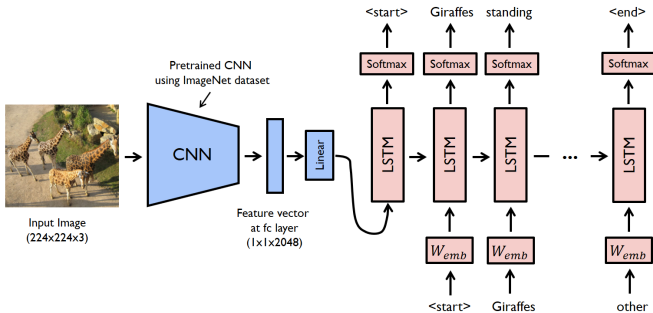


Case study 3: Art Generation

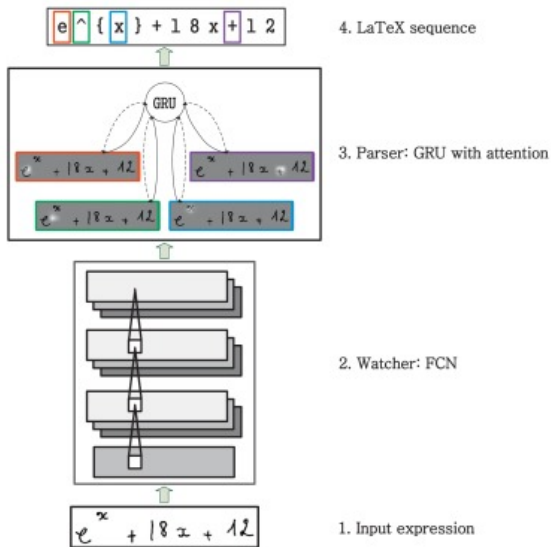
Image generation process

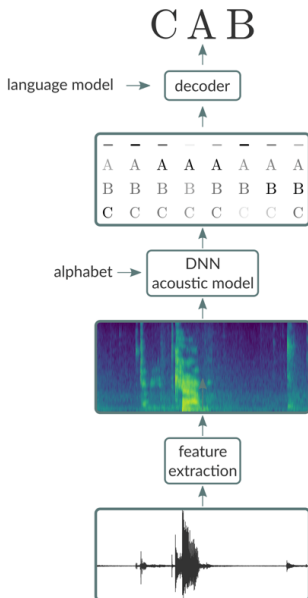


Automatic captioning



Automatic 2D to 1D





Transforming Paintings and Photos Into Animations With AI

<https://news.developer.nvidia.com/transforming-paintings-and-photos-into-animations-with-ai/>

