

Indexación, búsqueda y análisis en repositorios multimedia: imagen-video

Juan Carlos San Miguel

Agenda: Multimedia (imagen, video)

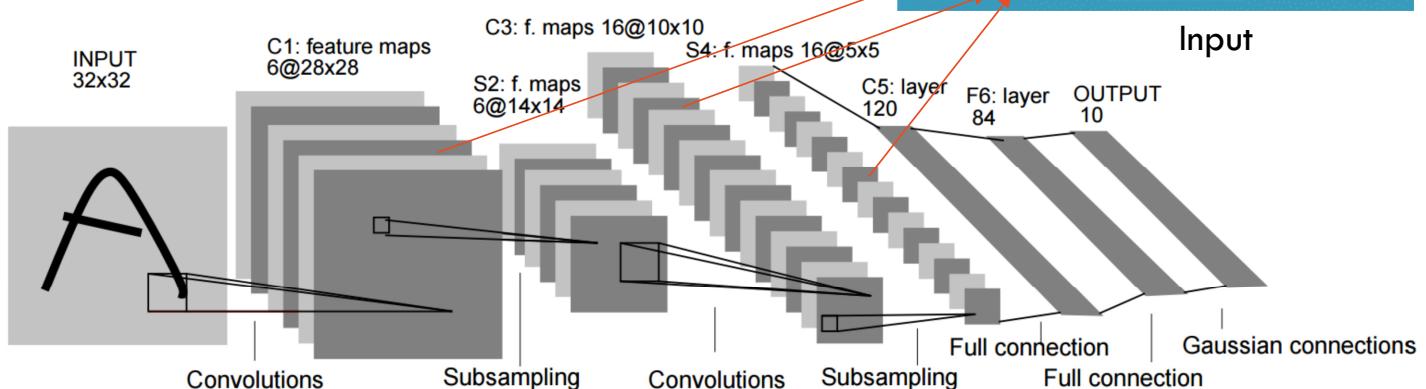
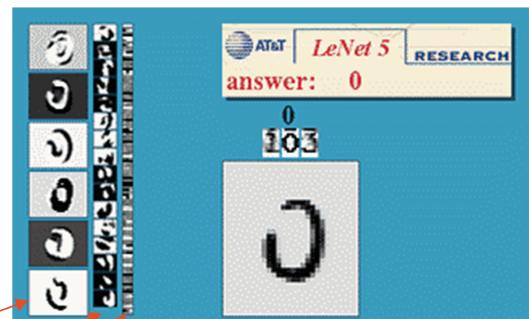
Temas avanzados

- Arquitecturas populares en clasificación
 - ZFNet
 - VGG
 - GoogleLeNet
 - ResNet
- Transfer Learning

Arquitecturas populares en clasificación

➤ LeNet: inicios

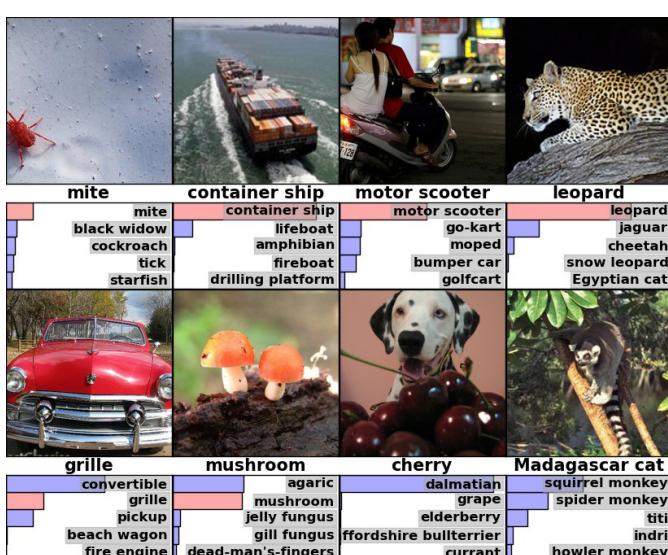
- Conv filtros son 5×5 , aplicados con stride 1
- 2×2 submuestreo (Pooling) con stride 2



Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proc. IEEE 86(11): 2278–2324, 1998.

Arquitecturas populares en clasificación

➤ ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

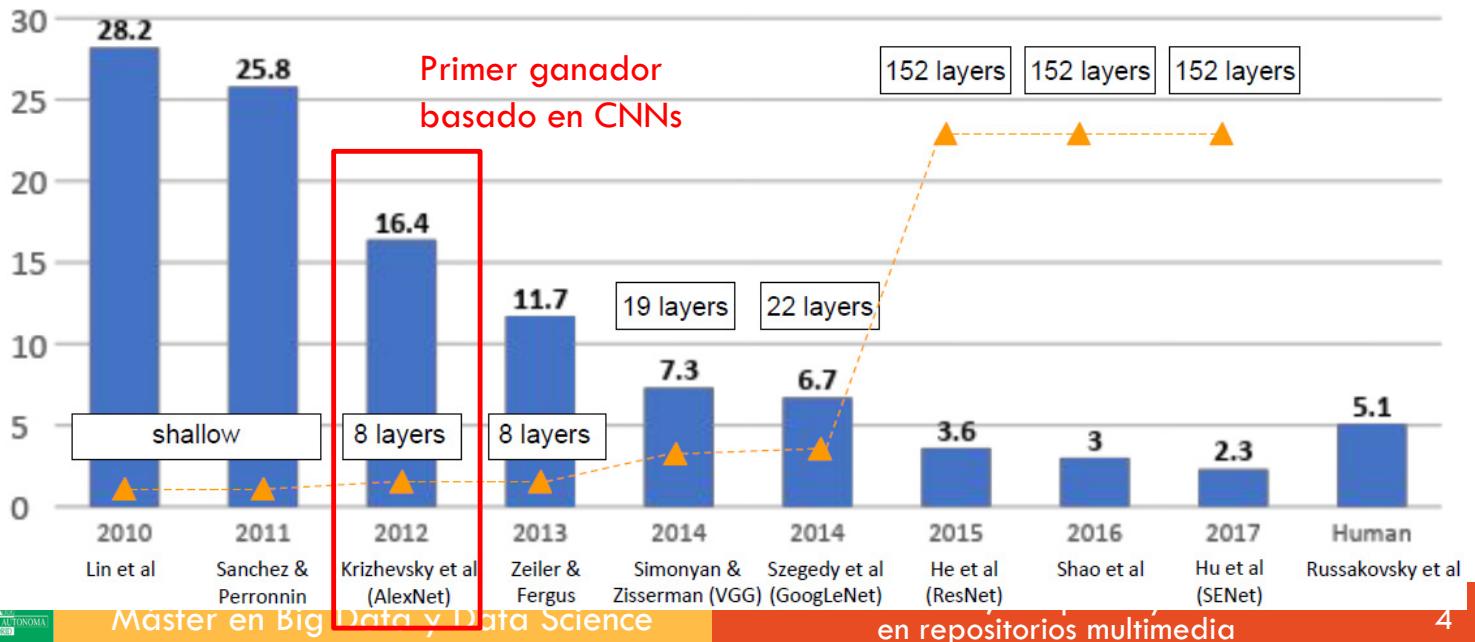


IMAGENET

- ~14M de imágenes obtenidas de internet
- Anotadas manualmente con 20k categorías
- Competición/Challenge:
 - 1k categorías
 - 1.2M imágenes train
 - 150K imágenes test

Arquitecturas populares en clasificación

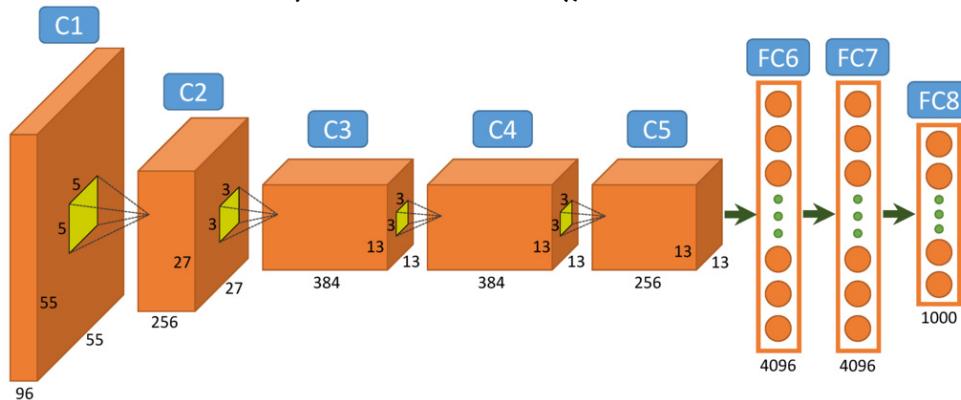
➤ ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Arquitecturas populares en clasificación

➤ AlexNet: ganador ILSVRC 2012

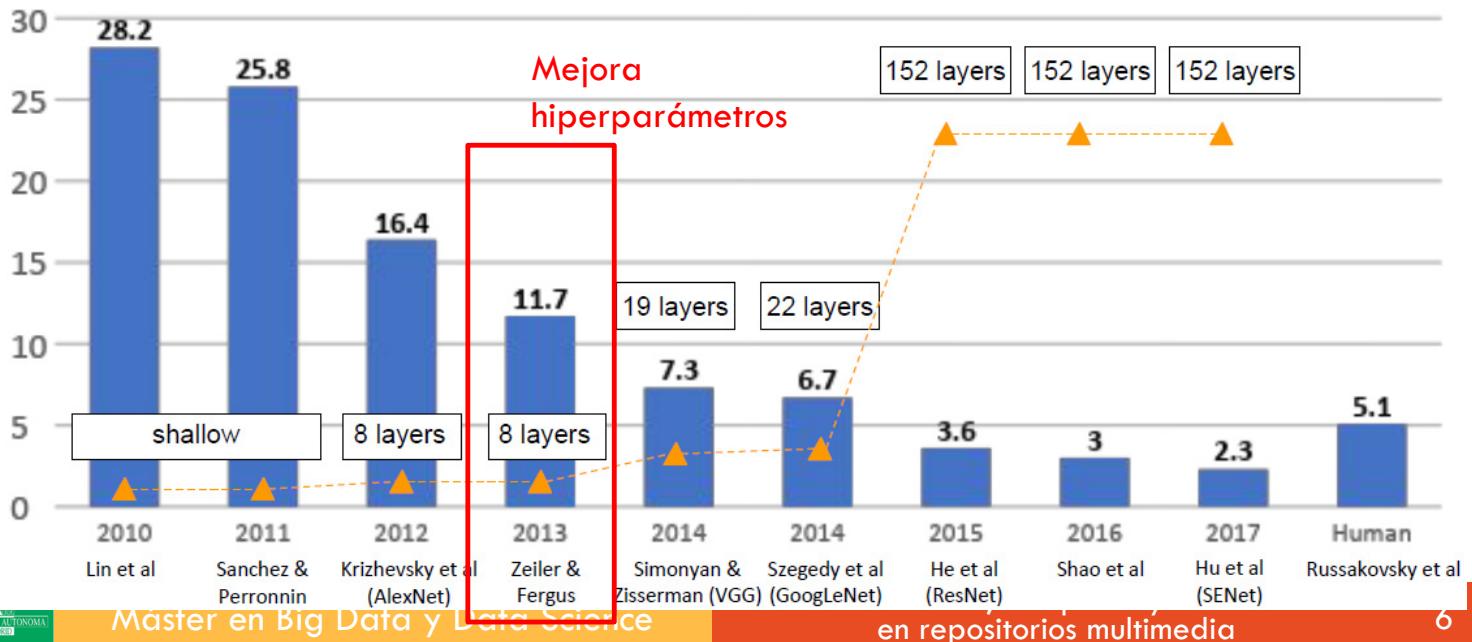
- Max pooling, ReLU nonlinearity, regularización Dropout
- Más datos y modelo más complejo (7 capas ocultas, ~650K unidades, ~60M params)
- Aceleración con GPUs (~50x sobre CPU), entrenamiento ~1 semana con 2 GPUs



A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012

Arquitecturas populares en clasificación

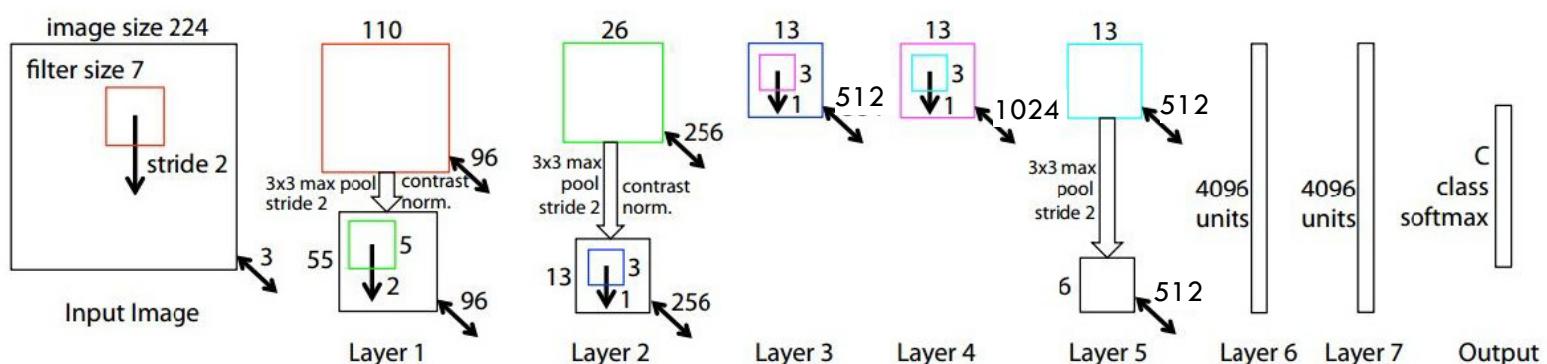
➤ ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Arquitecturas populares en clasificación: ZFNet

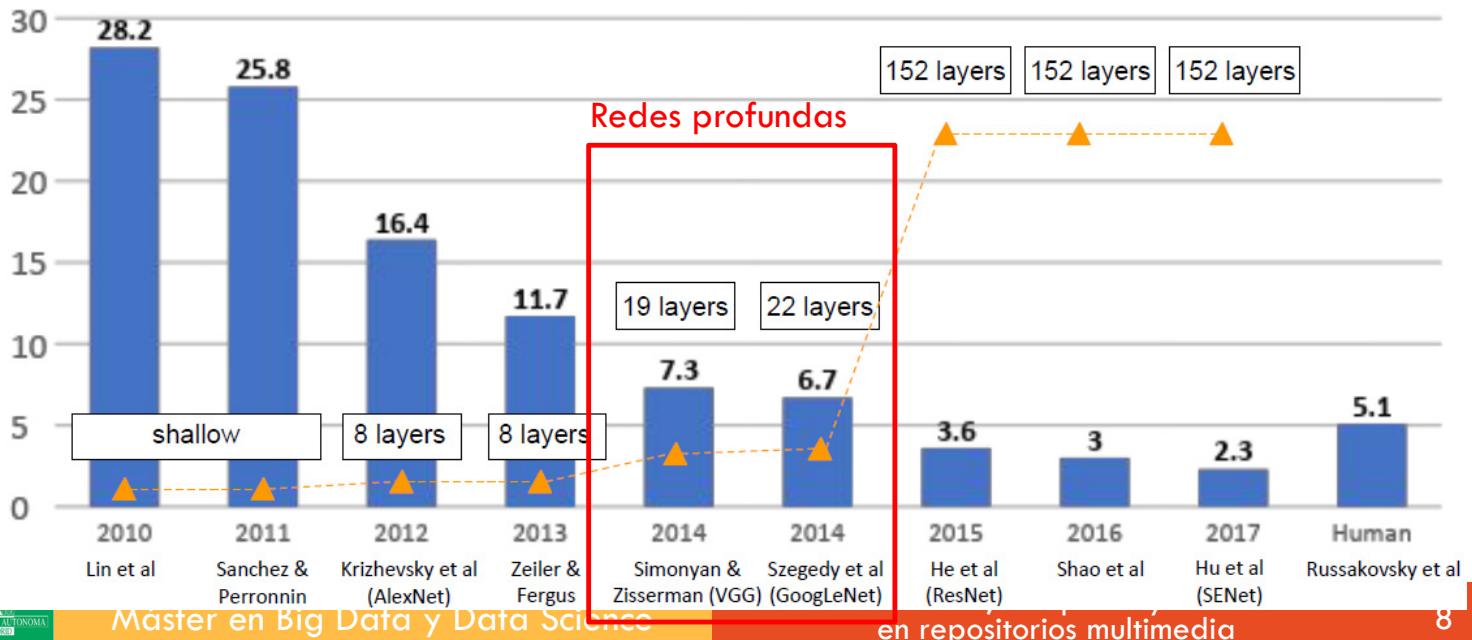
➤ Cambios con respecto a AlexNet

- CONV1: 11×11 stride 4 $\rightarrow 7 \times 7$ stride 2
- CONV3,4,5: 384/384/256 filtros \rightarrow 512/1024/512 filtros
- Mejora error Top-5 16.4% \rightarrow 11.7%



Arquitecturas populares en clasificación

➤ ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Máster en Big Data y Data Science

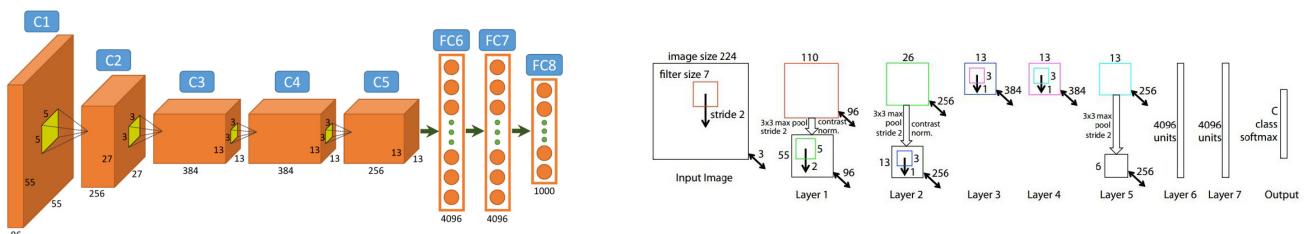
en repositorios multimedia

8



➤ Seleccione la afirmación correcta

- a) AlexNet puede procesar una imagen de cualquier tamaño
- b) ZFNet puede procesar una imagen de cualquier tamaño
- c) Ambas AlexNet y ZFNet que las imágenes tengan tamaño fijo



Máster en Big Data y Data Science

Indexación, búsqueda y análisis
en repositorios multimedia

9

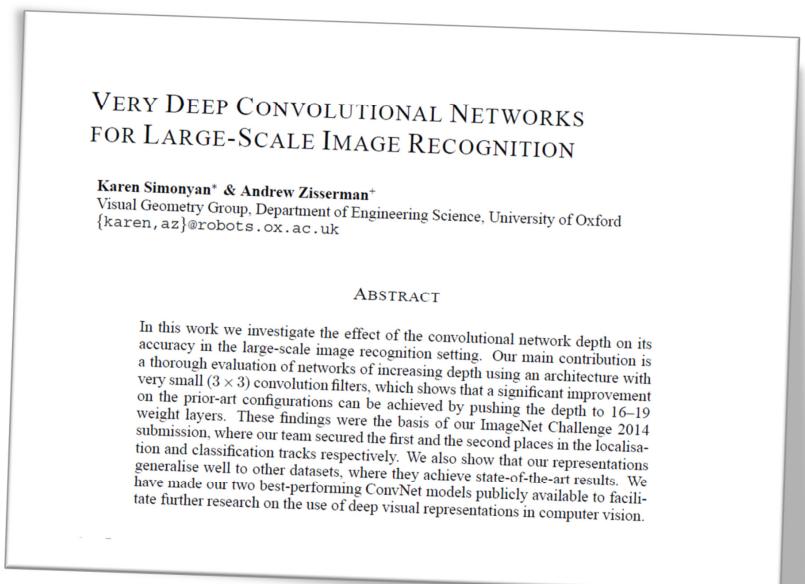
Arquitecturas populares en clasificación

➤ VGG Net

- Oxford University (UK)
- Publicado en 2015
- International Conference on Learning Representation
- +33K citas (Ene 2020)

Disponible en

http://www.robots.ox.ac.uk/~vgg/research/very_deep/



Arquitecturas populares en clasificación: VGG

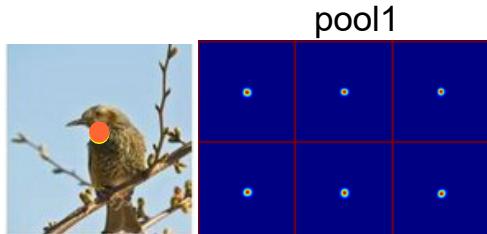
- 2^a posición en ILSVRC 2014
- Filtros pequeños, Redes más profundas
- 8 capas (AlexNet)
→ 16-19 capas en VGG
- Solo 3x3 CONV stride 1, pad 1
y 2x2 MAX POOL stride 2



Arquitecturas populares en clasificación: VGG

- ¿Por qué utilizar filtros más pequeños? (3x3 conv)
 - Apilamiento de tres capas 3x3 conv (stride 1) tiene un efecto similar (**effective receptive field**) como una capa con filtros 7x7 conv
 - Tres capas 3x3 convs para obtener C mapas de características ($3 \times 3 \times 3 \times C^2 = 27C^2$ params)
 - una capa de filtros 7x7 convs ($1 \times 7 \times 7 \times C^2 = 49C^2$ params)

Campo receptivo estimado



pool1



El tamaño real del campo receptivo es mucho menor que el tamaño teórico

conv3

pool5

Indexación, búsqueda y análisis
en repositorios multimedia

Máster en Big Data y Data Science

12

Arquitecturas populares en clasificación: VGG

- Complejidad VS performance

Model	top-5 classification error on ILSVRC-2012 (%)	
	validation set	test set
16-layer	7.5%	7.4%
19-layer	7.5%	7.3%
model fusion	7.1%	7.0%

Model	VOC-2007 (mean AP, %)	VOC-2012 (mean AP, %)	Caltech-101 (mean class recall, %)	Caltech-256 (mean class recall, %)
16-layer	89.3	89.0	91.8±1.0	85.0±0.2
19-layer	89.3	89.0	92.3±0.5	85.1±0.3
model fusion	89.7	89.3	92.7±0.5	86.2±0.3

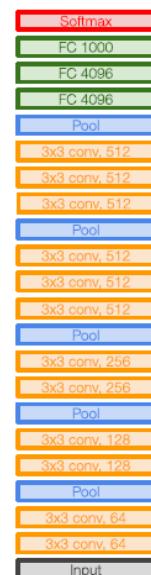
ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
FC-4096	FC-4096	FC-4096	FC-4096	FC-4096	FC-4096
FC-1000	FC-1000	FC-1000	FC-1000	FC-1000	FC-1000
soft-max	soft-max	soft-max	soft-max	soft-max	soft-max

Mejor

Máster en Big Data y Data Science

Arquitecturas populares en clasificación: VGG

INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150K$ params: 0 (not counting biases)
 CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 3) \times 64 = 1,728$
 CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 64) \times 64 = 36,864$
 POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800K$ params: 0
 CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 64) \times 128 = 73,728$
 CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 128) \times 128 = 147,456$
 POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400K$ params: 0
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 128) \times 256 = 294,912$
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
 POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200K$ params: 0
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: 0
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25K$ params: 0
 FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$
 FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$
 FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$



VGG16

TOTAL memory: $24M * 4 \text{ bytes} \approx 96\text{MB} / \text{image}$ (for a forward pass)

Fuente: <http://cs231n.stanford.edu/>

14

Arquitecturas populares en clasificación: GoogleLeNet

➤ GoogleLeNet

- Desarrollado por Google
- Publicado en 2015
- Int. Conf. on Computer Vision and Pattern Recognition
- +19K citas (a Ene 2020)

Disponible en
<https://github.com/zy314099208/CVPR-2015-Going-Deeper-with-Convolutions>

Going Deeper with Convolutions

Christian Szegedy¹, Wei Liu², Yangqing Jia¹, Pierre Sermanet¹, Scott Reed³, Dragomir Anguelov¹, Dumitru Erhan¹, Vincent Vanhoucke¹, Andrew Rabinovich⁴

¹Google Inc. ²University of North Carolina, Chapel Hill ³University of Michigan, Ann Arbor ⁴Magic Leap Inc.

¹{szegedy, jia, yq, sermanet, dragomir, dumitru, vanhoucke}@google.com
²wliu@cs.unc.edu, ³reedscott@umich.edu, ⁴arabinovich@microsoft.com

Abstract

We propose a deep convolutional neural network architecture codenamed Inception that achieves the new state of the art for classification and detection in the ImageNet Large-Scale Visual Recognition Challenge 2014 (ILSVRC14). The main hallmark of this architecture is the improved utilization of the computing resources inside the network. By a carefully crafted design, we increased the depth and width of the network while keeping the computational budget constant. To optimize quality, the architectural decisions were based on the Hebbian principle and the intuition of multi-scale processing. One particular incarnation used in our submission for ILSVRC14 is called GoogLeNet, a 22 layers deep network, the quality of which is assessed in the context of classification and detection.

ger and bigger deep networks, but from the synergy of deep architectures and classical computer vision, like the R-CNN algorithm by Girshick et al [6]. Another notable factor is that with the ongoing traction of mobile and embedded computing, the efficiency of our algorithms – especially their power and memory use – gains importance. It is noteworthy that the considerations leading to the design of the deep architecture presented in this paper included this factor rather than having a sheer fixation on accuracy numbers. For most of the experiments, the models were designed to keep a computational budget of 1.5 billion multiply-adds at inference time, so that they do not end up being a purely academic curiosity, but could be put to real world use, even on large datasets, at a reasonable cost.

In this paper, we will focus on an efficient deep neural network architecture for computer vision, codenamed Inception, which derives its name from the Network in net-

Arquitecturas populares en clasificación: Google

- Incremento sustancial de capas
- 22 capas en la red
- “Network-in-network”

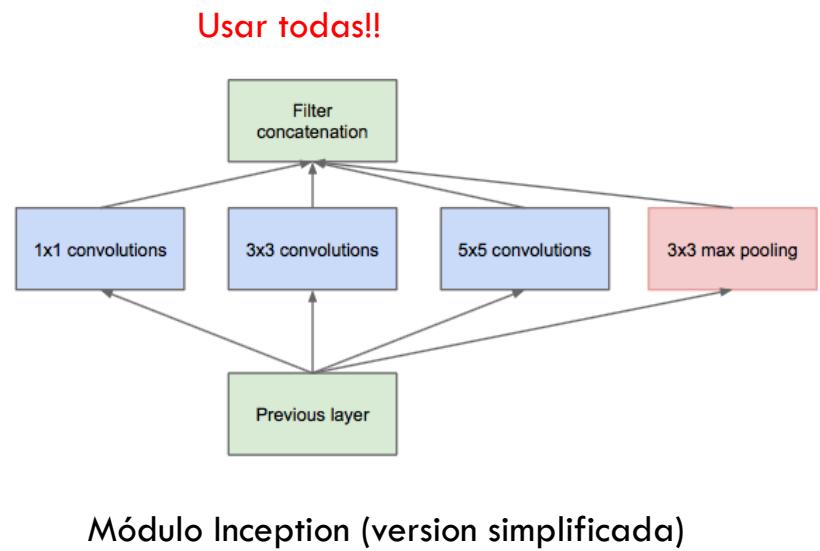
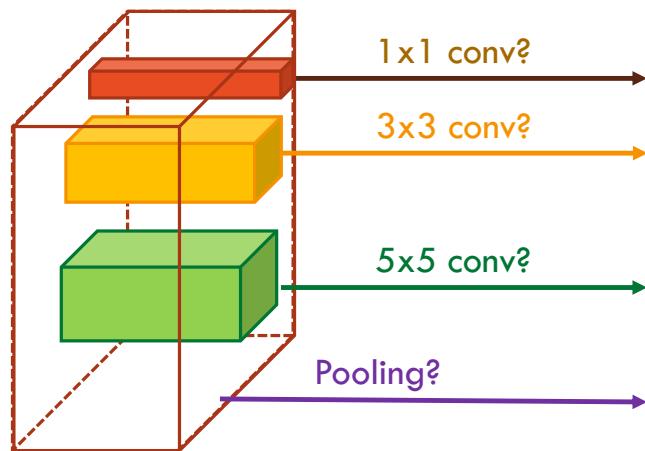


Arquitecturas populares en clasificación: GoogleLeNet

- Problemas derivados de incrementar la profundidad (*challenges of going deeper*)
 - Incremento del número de parámetros
 - Incremento del coste computacional

Arquitecturas populares en clasificación: GoogleLeNet

➤ Convoluciones?



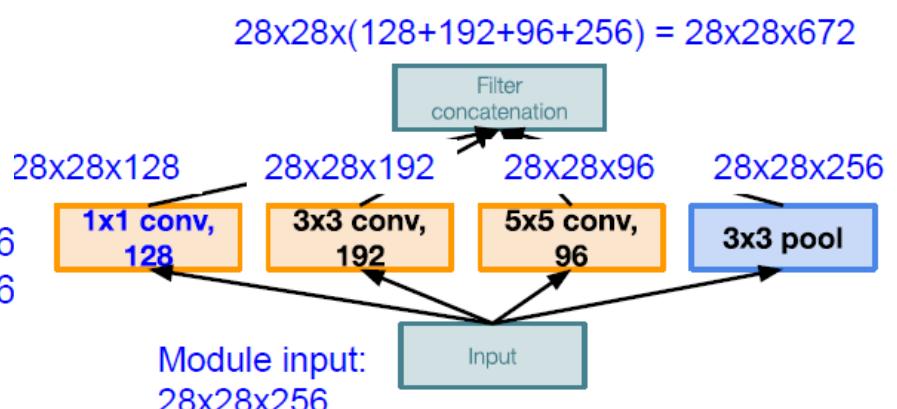
Arquitecturas populares en clasificación: GoogleLeNet

➤ ¿Cómo es el volumen de salida en cada módulo?

Conv Ops:

[1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
[3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$
[5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$

Total: 854M ops

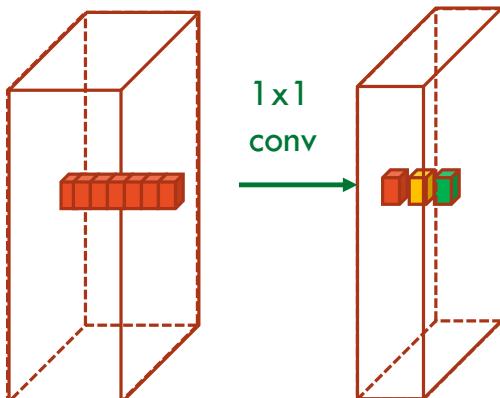


Naive Inception module

Arquitecturas populares en clasificación: GoogleLeNet

➤ Idea

- ¿Tienen sentido convoluciones 1x1?
- ¿Podemos reducir dimensionalidad (depth)?



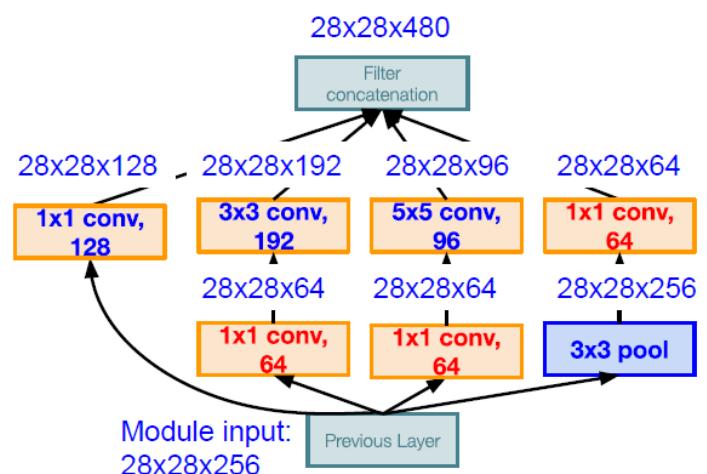
(c) Mohsen Ghafoorian

Arquitecturas populares en clasificación: GoogleLeNet

➤ Solución: convoluciones 1x1

Conv Ops:

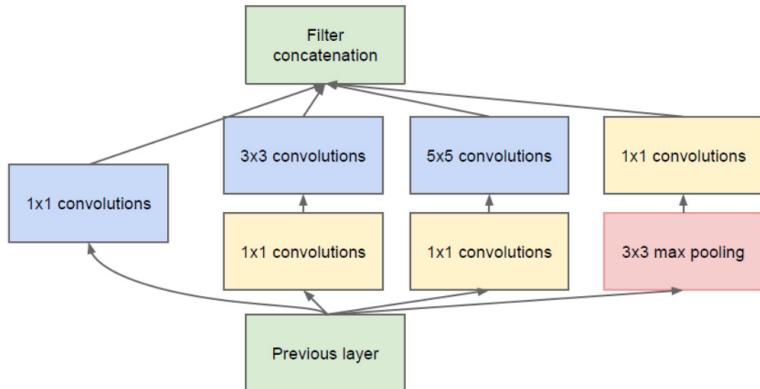
[1x1 conv, 64] 28x28x64x1x256
[1x1 conv, 64] 28x28x64x1x256
[1x1 conv, 128] 28x28x128x1x256
[3x3 conv, 192] 28x28x192x3x3x64
[5x5 conv, 96] 28x28x96x5x5x64
[1x1 conv, 64] 28x28x64x1x256
Total: 358M ops



Inception module with dimension reduction

Arquitecturas populares en clasificación: GoogleLeNet

➤ Arquitectura completa

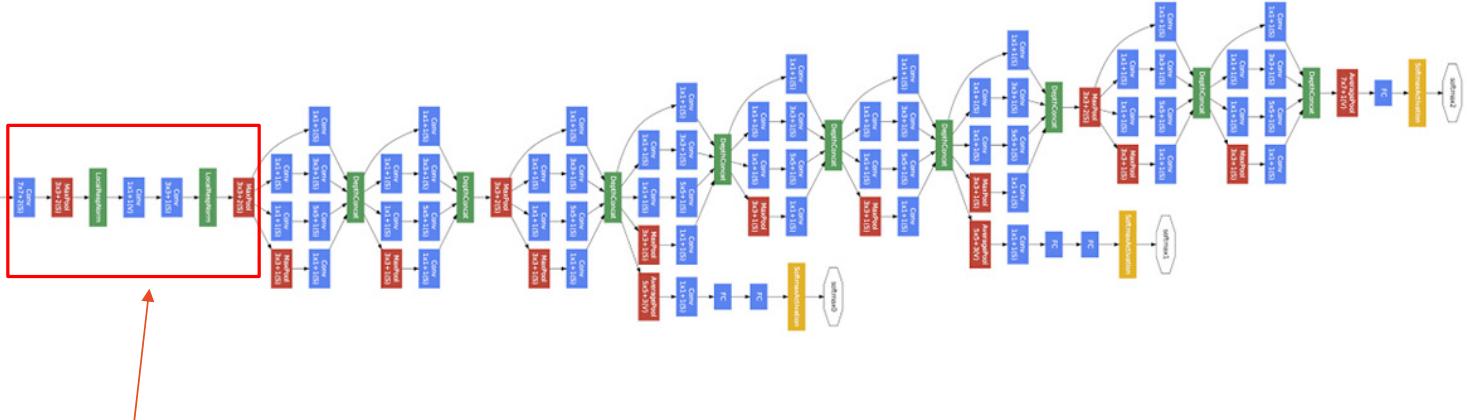


9 módulos *inception*



Arquitecturas populares en clasificación: GoogleLeNet

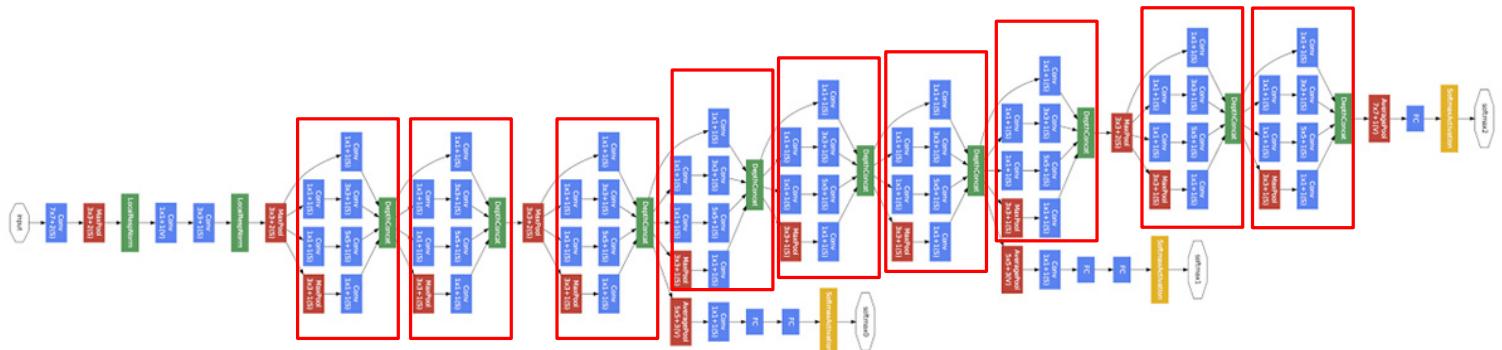
➤ Arquitectura completa



ENTRADA: Conv+MaxPool+Norm+2xConv+Norm+MaxPool

Arquitecturas populares en clasificación: GoogleLeNet

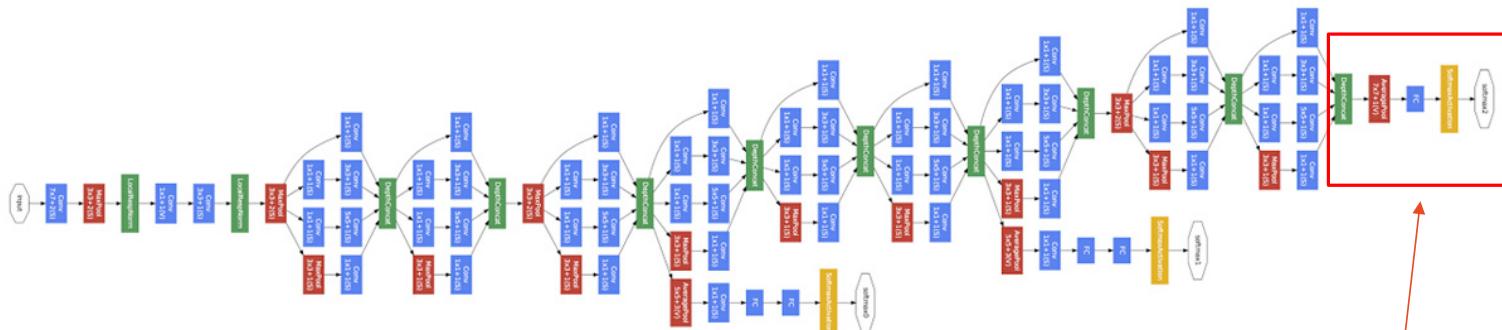
➤ Arquitectura completa



9 módulos *inception*

Arquitecturas populares en clasificación: GoogleLeNet

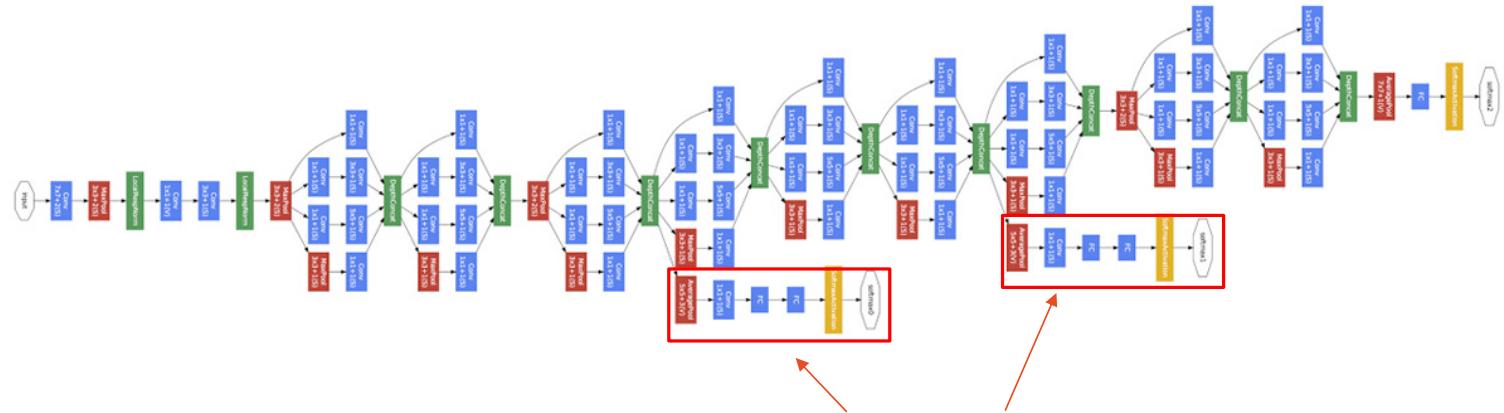
➤ Arquitectura completa



SALIDA: MaxPool + FC + Softmax

Arquitecturas populares en clasificación: GoogleLeNet

➤ Arquitectura completa



Arquitecturas populares en clasificación: GoogleLeNet

➤ Parámetros

type	patch size/ stride	output size	depth	#1x1	#3x3 reduce	#3x3	#5x5 reduce	#5x5	pool proj	params	ops
convolution	7x7/2	112x112x64	1							2.7K	34M
max pool	3x3/2	56x56x64	0								
convolution	3x3/1	56x56x192	2		64	192				112K	360M
max pool	3x3/2	28x28x192	0								
inception (3a)		28x28x256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28x28x480	2	128	128	192	32	96	64	380K	304M
max pool	3x3/2	14x14x480	0								
inception (4a)		14x14x512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14x14x512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14x14x512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14x14x528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14x14x832	2	256	160	320	32	128	128	840K	170M
max pool	3x3/2	7x7x832	0								
inception (5a)		7x7x832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7x7x1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7x7/1	1x1x1024	0								
dropout (40%)		1x1x1024	0								
linear		1x1x1000	1							1000K	1M
softmax		1x1x1000	0								

Arquitecturas populares en clasificación: GoogleLeNet

➤ Parámetros de VGG Net

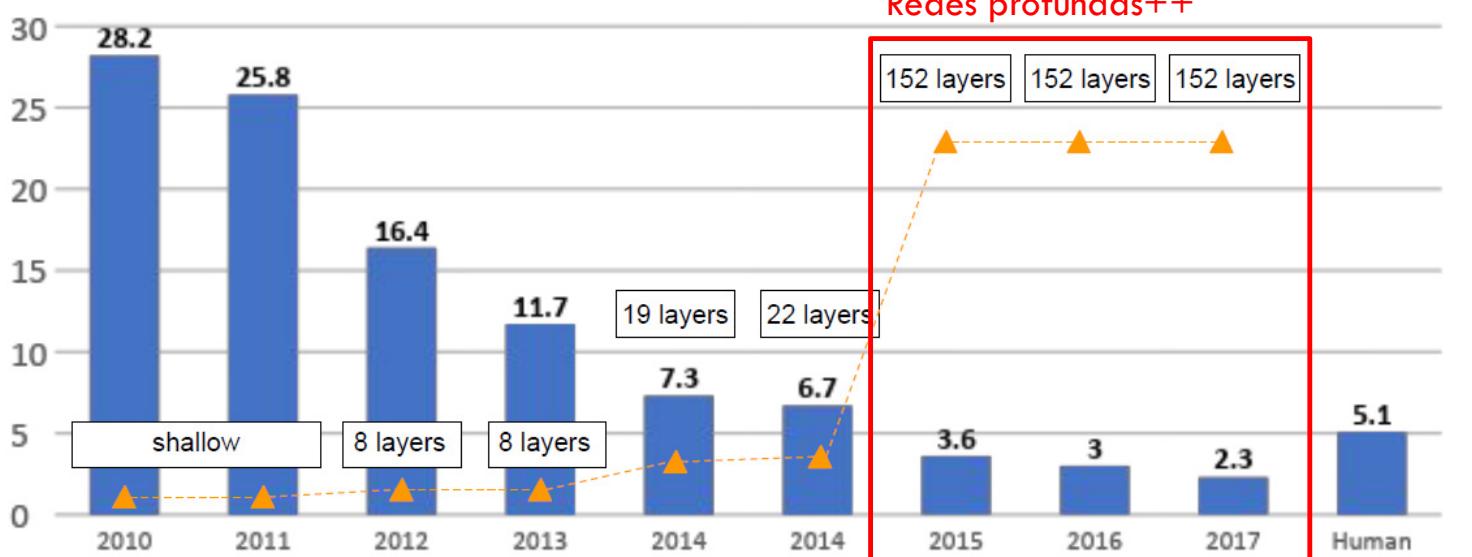
INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150K$ params: 0 (not counting biases)
 CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 3) \times 64 = 1,728$
 CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 64) \times 64 = 36,864$
 POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800K$ params: 0
 CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 64) \times 128 = 73,728$
 CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 128) \times 128 = 147,456$
 POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400K$ params: 0
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 128) \times 256 = 294,912$
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
 CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
 POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200K$ params: 0
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: 0
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
 POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25K$ params: 0
 FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$
 FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$
 FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	conv3-64
conv3-64	conv3-64	conv3-64	conv3-64
maxpool			
conv3-128	conv3-128	conv3-128	conv3-128
conv3-128	conv3-128	conv3-128	conv3-128
maxpool			
conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv1-256	conv3-256
		conv3-256	conv3-256
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv1-512	conv3-512
	conv3-512	conv3-512	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv1-512	conv3-512
	conv3-512	conv3-512	conv3-512
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

Arquitecturas populares en clasificación

➤ ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

Redes profundas++



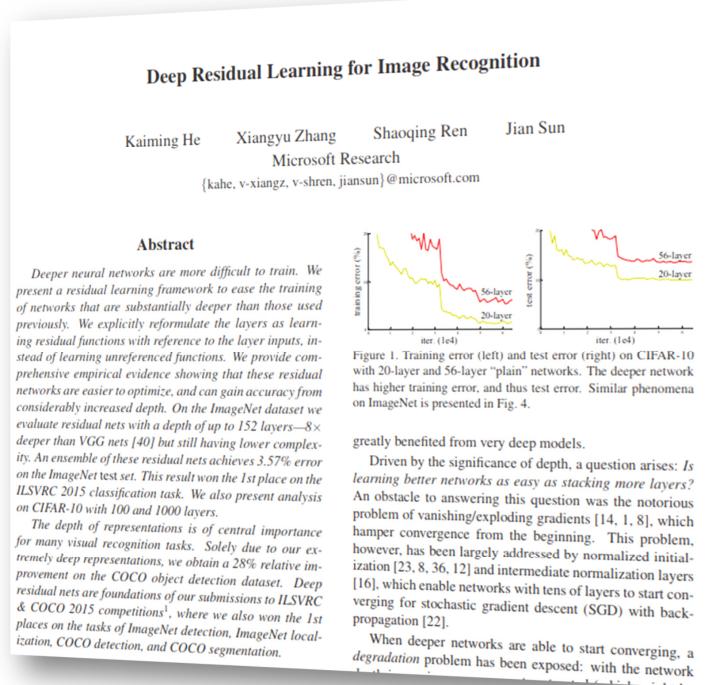
Arquitecturas populares en clasificación: ResNet

➤ ResNet

- Desarrollado por Microsoft
- Publicado en 2016
- International Conference on Computer Vision and Pattern Recognition
- +40K citas (a Ene 2020)

Disponible at

<https://github.com/KaimingHe/deep-residual-networks>



Máster en Big Data y Data Science

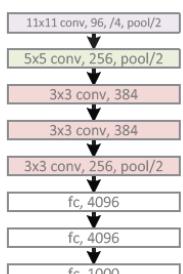
Indexación, búsqueda y análisis en repositorios multimedia

30

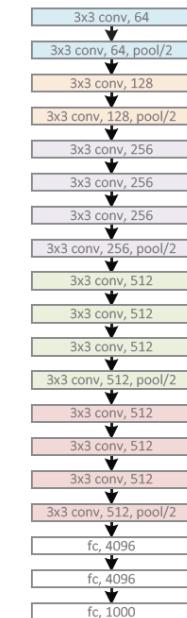
Arquitecturas populares en clasificación: ResNet

➤ Revolución de profundidad (going deeper again)

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, 152 layers
(ILSVRC 2015)



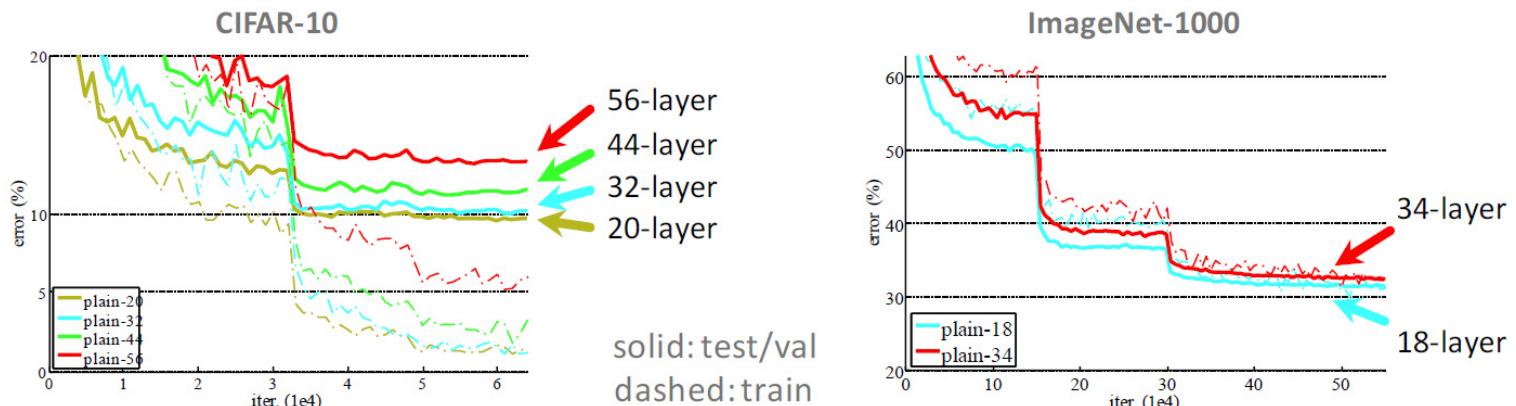
Máster en Big Data y Data Science

Indexación, búsqueda y análisis en repositorio



Arquitecturas populares en clasificación: ResNet

➤ Revolución de profundidad (*curse of dimensionality*)

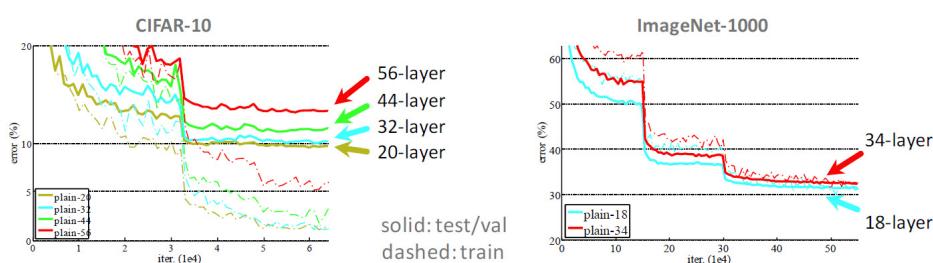


¿Mejor aprendizaje simplemente añadiendo más capas?

K. He, et al. [Deep residual learning for image recognition](#). arXiv preprint arXiv:1512.03385, CVPR 2016

Arquitecturas populares en clasificación: ResNet

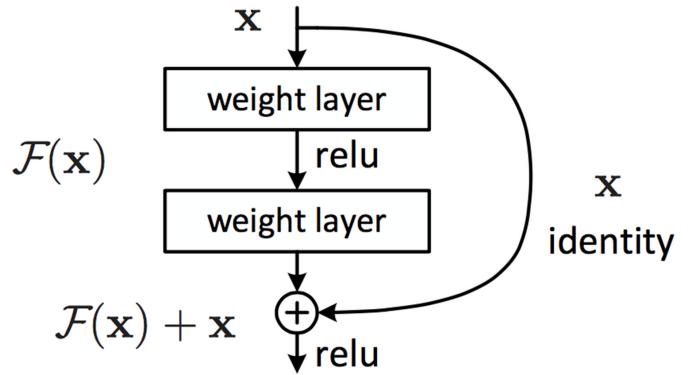
?



¿Por qué la diferencia entre el error de entrenamiento/test es más alto para CIFAR-10 que para el dataset de imagenet-1000?

Arquitecturas populares en clasificación: ResNet

- Ganador ILSVRC 2015
- Bloques residuales (*Residual blocks*)
 - Adelantar el input
 - Permite entrenar redes con muchas capas

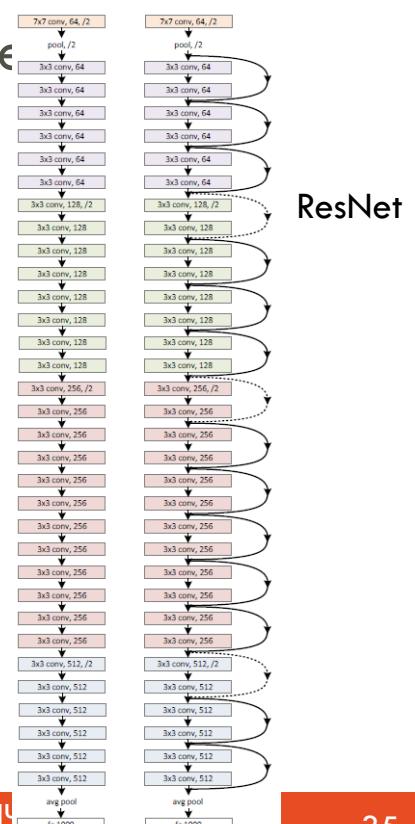


K. He, et al. [Deep residual learning for image recognition](#). arXiv preprint arXiv:1512.03385, CVPR 2016

Arquitecturas populares en clasificación: ResNet

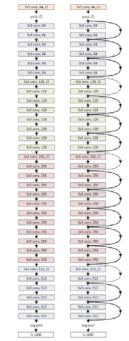
- Ganador ILSVRC 2015
- Bloques residuales (*Residual blocks*)
- Diseño simple
 - Casi todo son 3×3 conv
 - Sin Fully Connected ni dropout

Red
“plana”

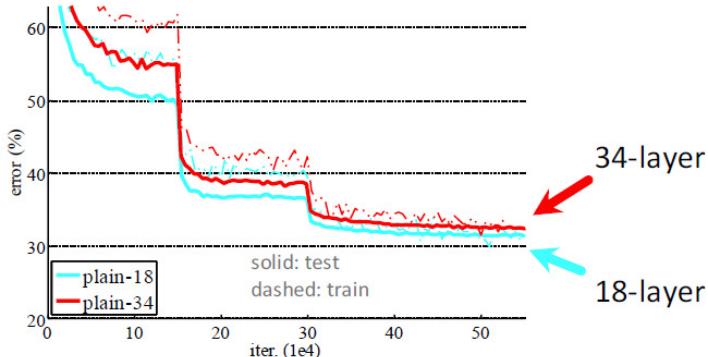


Arquitecturas populares en clasificación: ResNet

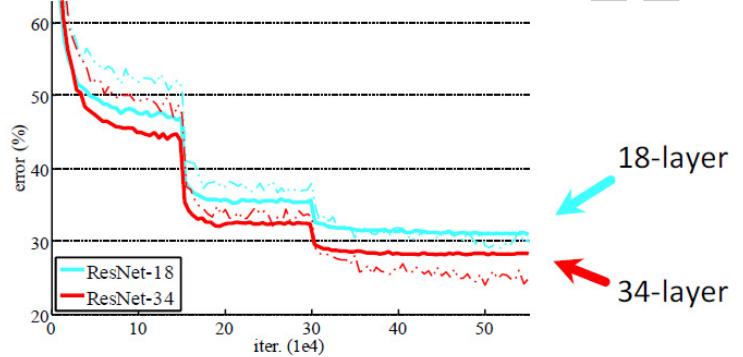
- Deep ResNets pueden entrenarse mejor



ImageNet plain nets



ImageNet ResNets



K. He, et al. [Deep residual learning for image recognition](#). arXiv preprint arXiv:1512.03385, CVPR 2016

Agenda: Multimedia (imagen, video)

Temas avanzados

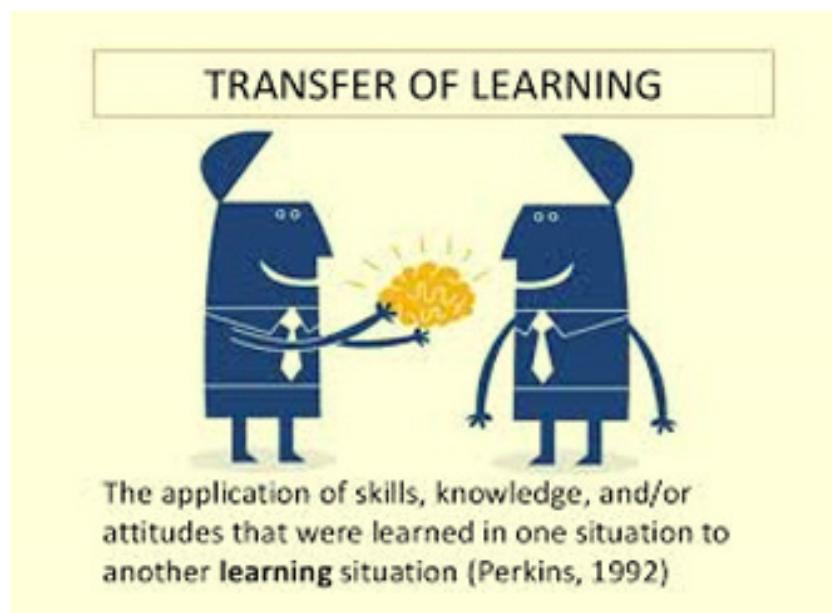
- Arquitecturas populares en clasificación
- Transfer Learning
 - Introducción
 - Impacto
 - Concepto
 - *Transfer learning* en CNNs
 - Estrategias
 - Casos de estudio
 - Ejemplos

INTRODUCCION



<https://youtu.be/0YhJxJZOWBw>

INTRODUCCION



INTRODUCCION

ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.



5000 object classes and concepts
16M images

<https://www.vision.ee.ethz.ch/webvision/>

	VGGNet	DeepVideo	GNMT
Used For	Identifying Image Category	Identifying Video Category	Translation
Input	Image	Video	English Text
Output	1000 Categories	47 Categories	French Text
Parameters	140M	~100M	380M
Data Size	1.2M Images with assigned Category	1.1M Videos with assigned Category	6M Sentence Pairs, 340M Words
Dataset	ILSVRC-2012	Sports-1M	WMT'14

INTRODUCCION

Se necesitan grandes cantidades de datos para
entrenar/utilizar redes neuronales convolucionales (CNNs)

¿Datasets pequeños?



1500 imágenes entrenamiento / 1500 imágenes test

INTRODUCCION

Se necesitan grandes cantidades de datos para entrenar/utilizar redes neuronales convolucionales (CNNs)

¿Datasets pequeños?



¿Existe alternativa?

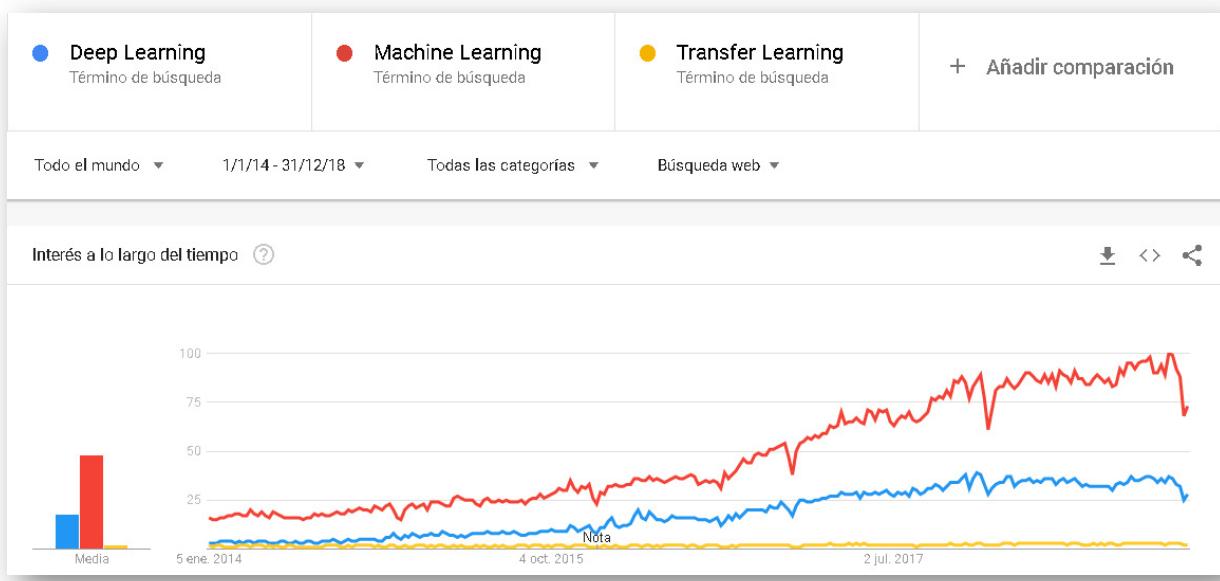


Si, transferencia del conocimiento (*transfer learning*)

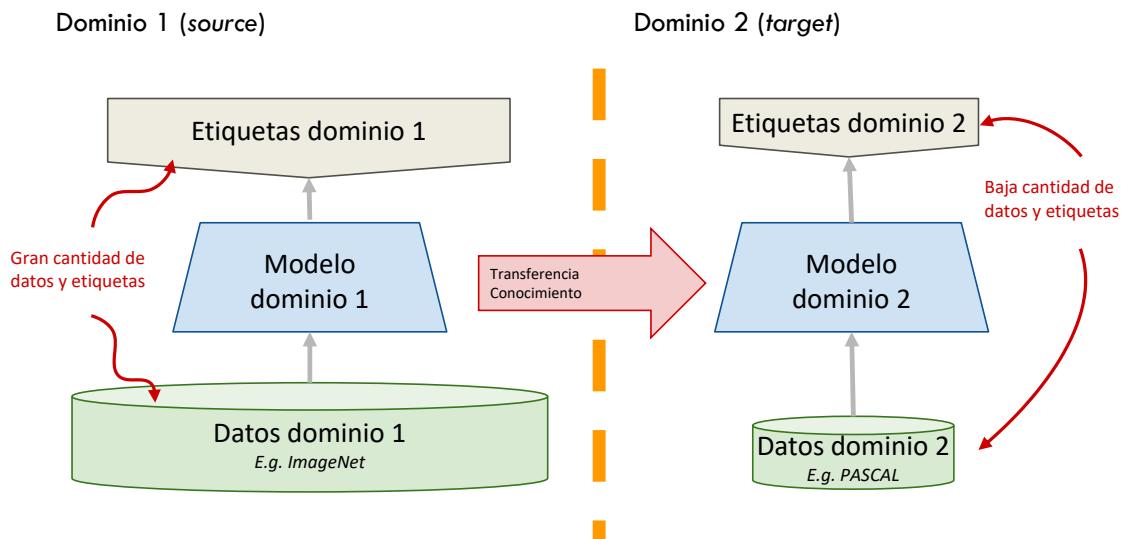


Utilizar CNNs como extractores de características

IMPACTO

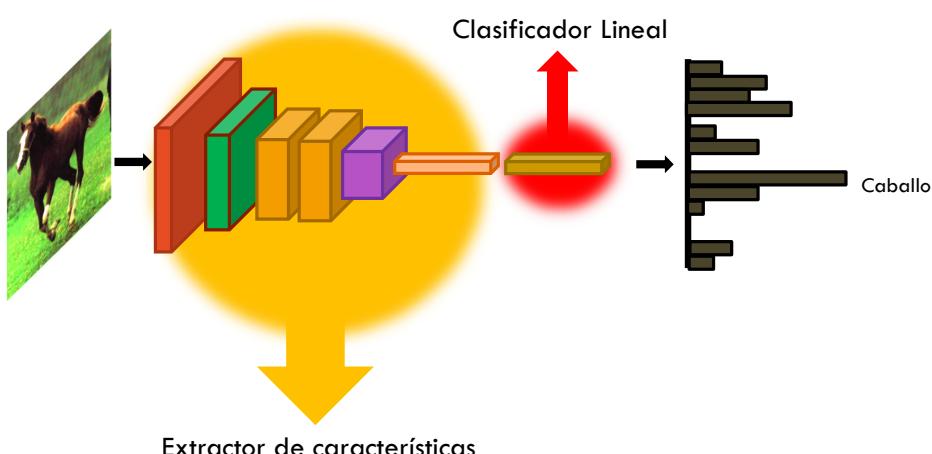


CONCEPTO



TRANSFER LEARNING EN CNNs

- Descomposición funcional de una red CNN



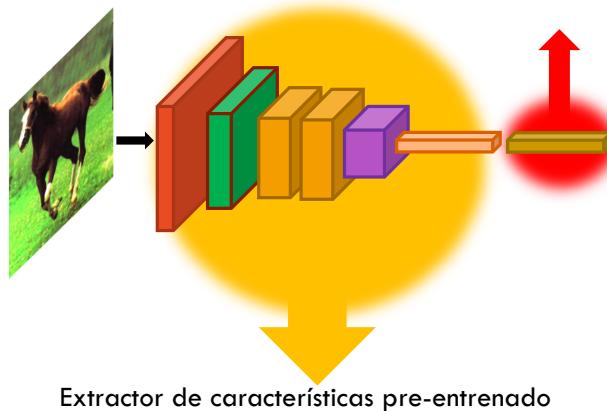
Fuente: <http://www.cs.cornell.edu/courses/cs4670/2019sp/>

TRANSFER LEARNING EN CNNs

¿Qué transferir a un nuevo problema de clasificación?

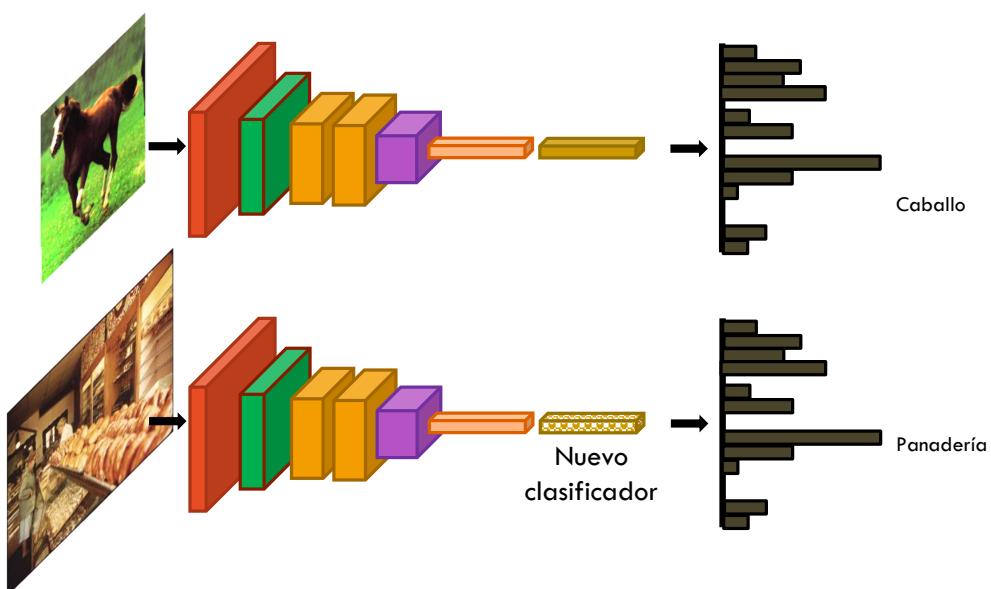
➤ Idea:

- Congelar (la mayoría de) parámetros del extractor de características
- Reentrenar el clasificador



Fuente: <http://www.cs.cornell.edu/courses/cs4670/2019sp/>

TRANSFER LEARNING EN CNNs



ESTRATEGIAS

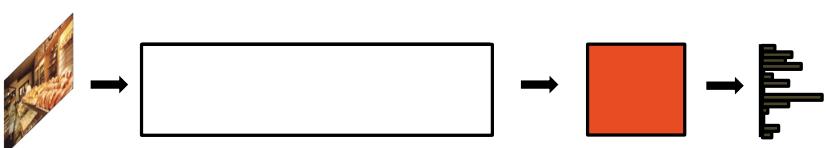
Estrategia 1
Entrenar la red completa



Estrategia 2
Entrenar clasificador y parte del extractor de características



Estrategia 3
Entrenar solamente el clasificador



Parámetros fijos



Parámetros a entrenar



ESTRATEGIAS

Tamaño del target dataset



Estrategia 2
Entrenar clasificador y parte del resto



Estrategia 2
Entrenar clasificador y parte del resto



Estrategia 3
Entrenar solamente clasificador



Similitud source-target datasets

Indexación, búsqueda y análisis
en repositorios multimedia

CASO DE ESTUDIO

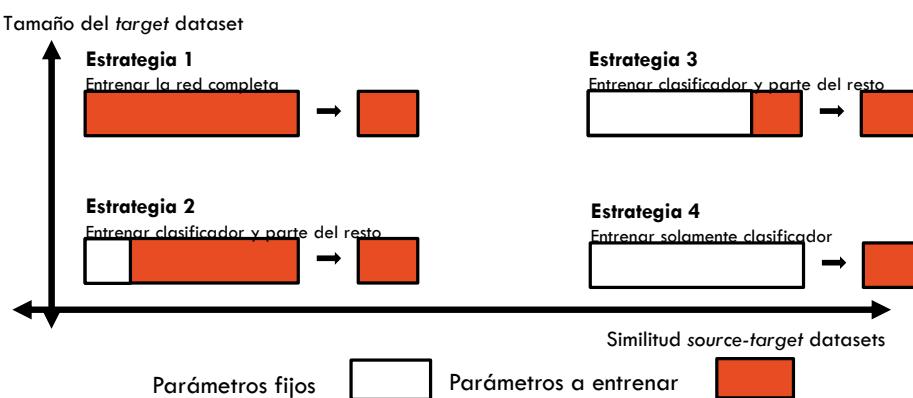
- Reconocimiento de escenas:
 - 15 categorías
 - 1500 imágenes entrenamiento / 1500 imágenes test



02:00

- Sea el dataset de escenas mencionado con la partición mencionada en la diapositiva anterior, que estrategia de *Transfer Learning* sería más adecuada aplicar:

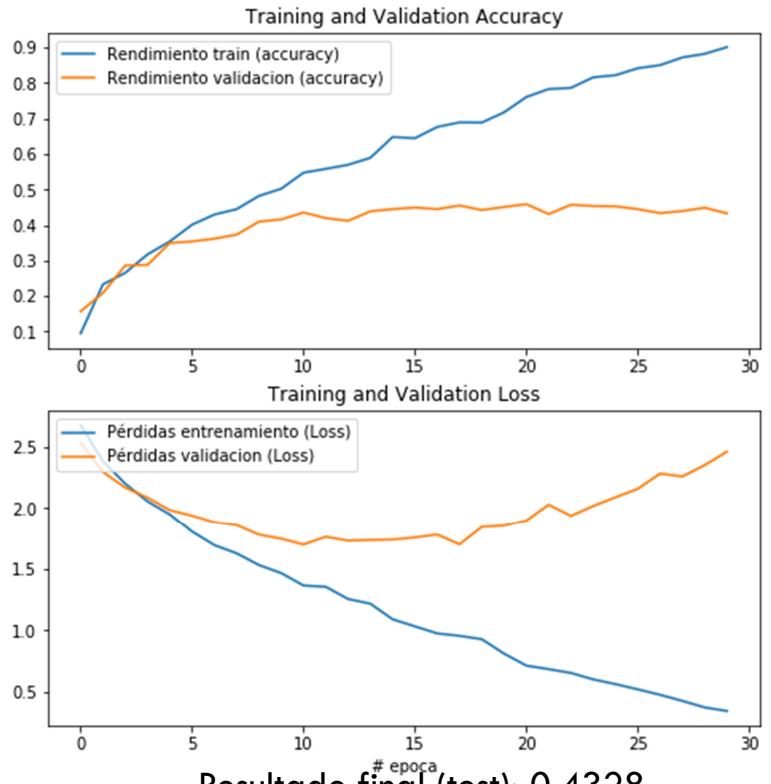
- a) Estrategia 1 b) Estrategia 2 c) Estrategia 3 d) Estrategia 4



CASO ESTUDIO

➤ Resultados iniciales

```
Model: "Red ejemplo"
Layer      Output_shape Param #
=====
conv2d      (None, 30, 30, 6) 168
avg_pooling2d (None, 15, 15, 6) 0
conv2d_1     (None, 13, 13, 16) 880
avg_pooling2d (None, 6, 6, 16) 0
flatten      (None, 576) 0
dense        (None, 120) 69240
dense_1      (None, 84) 10164
dense_2      (None, 15) 1275
=====
Total params: 81,727
Trainable params: 81,727
Non-trainable params: 0
```



Resultado final (test): 0.4328

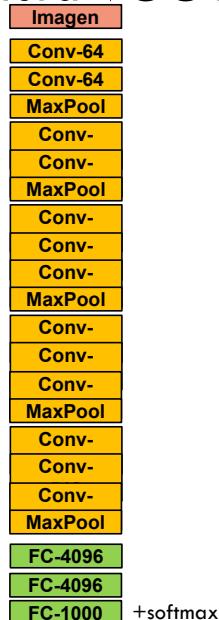
Indexación, búsqueda y análisis
en repositorios multimedia

52

Máster en Big Data y Data Science

CASO ESTUDIO

➤ Arquitectura VGG16 pre-entrenada en Imagenet



ImageNet Challenge

IMAGENET



- 1,000 object classes (categories).

- Images:
 - 1.2 M train
 - 100k test.

92.7% top-5 rendimiento (accuracy) en test

143,667,240 parámetros

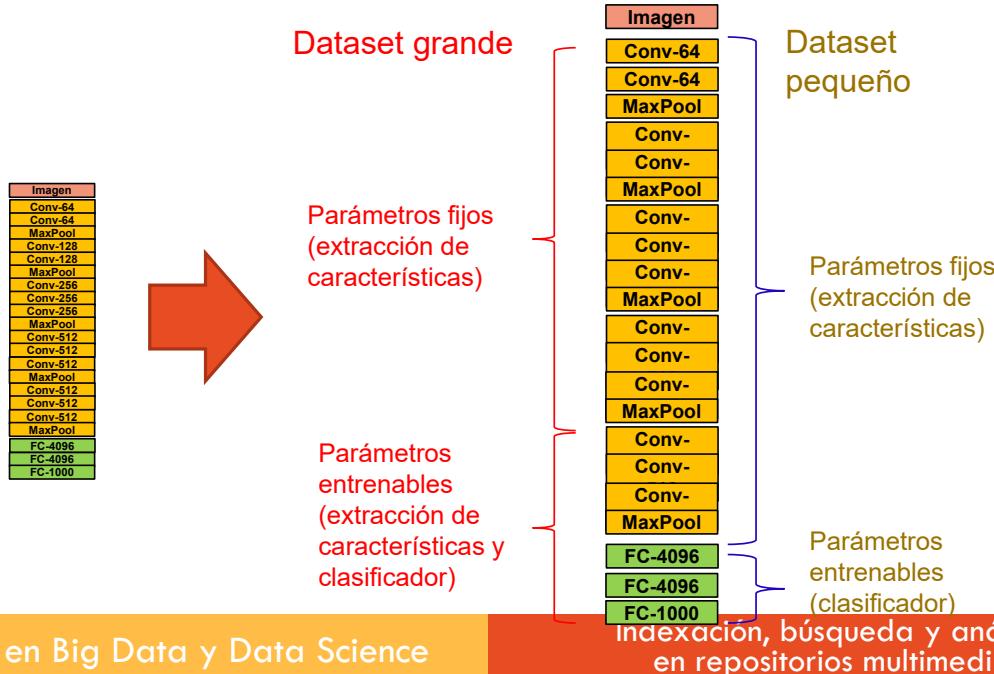
Máster en Big Data y Data Science

Indexación, búsqueda y análisis
en repositorios multimedia

53

CASO ESTUDIO

➤ Opciones para transferencia con datasets similares



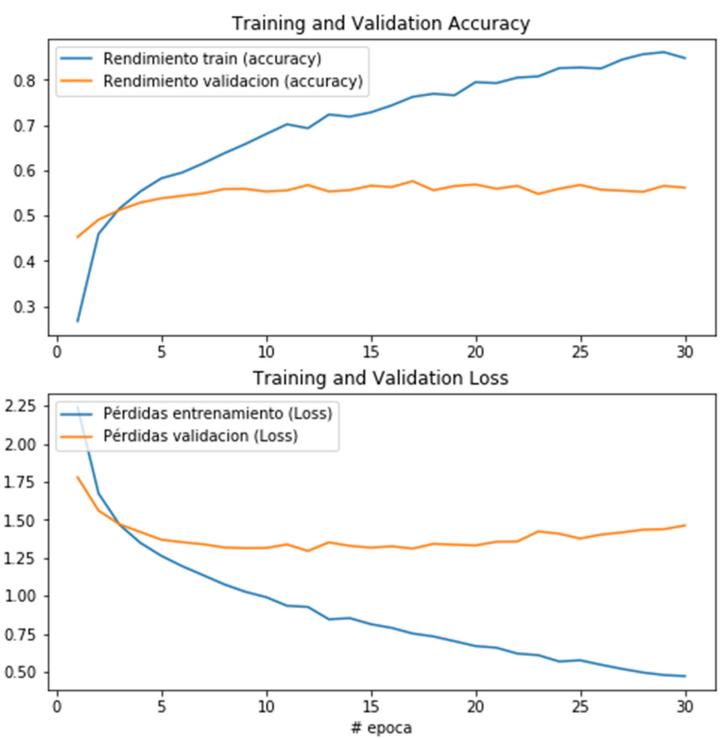
CASO ESTUDIO

➤ Resultados finales 1

```
Model: "Adaptada VGG19 a Scenes15"
Layer      Output_shape Param #
=====
vgg16*      (None, 1, 1, 512)
14714688

dense_1      (None, 120) 61560
dense_2      (None, 15) 1815
=====
Total params: 14,778,063
Trainable params: 63,375
Non-trainable params: 14,714,688
=====
*vgg16 sin FC-4096, FC-4096 y FC-1000
```

+30% rendimiento
con Transfer Learning



Resultado final (test): 0.5615

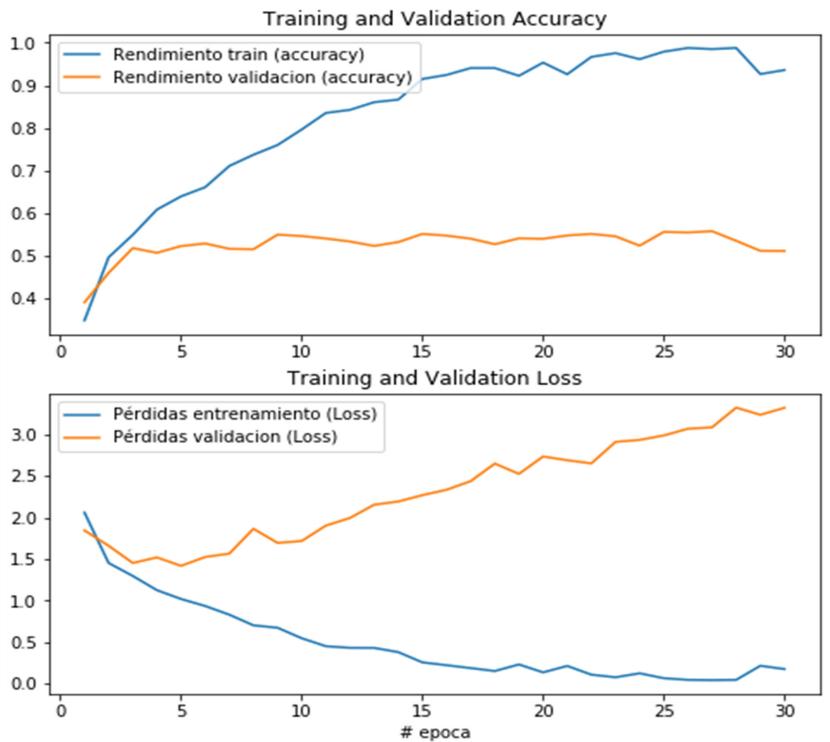
CASO ESTUDIO

➤ Resultados finales 2

Model: "Adaptada VGG19 a Scenes15"

Layer	Output_shape	Param #
vgg16*	(None, 1, 1, 512)	14714688
dense_1	(None, 4096)	2101248
dense_2	(None, 4096)	16781312
dense_3	(None, 15)	61455
Total params:	33,658,703	
Trainable params:	18,944,015	
Non-trainable params:	14,714,688	

*vgg16 sin FC-4096, FC-4096 y FC-1000



Resultado final (test): 0.5114

Indexación, búsqueda y análisis
en repositorios multimedia

56

Máster en Big Data y Data Science

EJEMPLOS

➤ Transfer Learning para detección de cáncer de piel:

Stanford | News

JANUARY 25, 2017

Deep learning algorithm does as well as dermatologists in identifying skin cancer

In hopes of creating better access to medical care, Stanford researchers have trained an algorithm to diagnose skin cancer.

BY TAYLOR KUBOTA

It's scary enough making a doctor's appointment to see if a strange mole could be cancerous. Imagine, then, that you were in that situation while also living far away from the nearest doctor, unable to take time off work and unsure you had the money to cover the cost of the visit. In a scenario like this, an option to receive a diagnosis through your smartphone could be lifesaving.

Universal access to health care was on the minds of computer scientists at Stanford when they set out to create an artificially intelligent diagnosis algorithm for skin cancer. They made a database of nearly 130,000 skin disease images and trained their algorithm to visually diagnose potential cancer. From the very first test, it performed with inspiring accuracy.

"We realized it was feasible, not just to do something well, but as

A photograph shows a person's arm being examined by a dermatologist using a handheld dermatoscope.

<https://news.stanford.edu/2017/01/25/artificial-intelligence-used-identify-skin-cancer/>

Máster en Big Data y Data Science

Indexación, búsqueda y análisis
en repositorios multimedia

57

EJEMPLOS

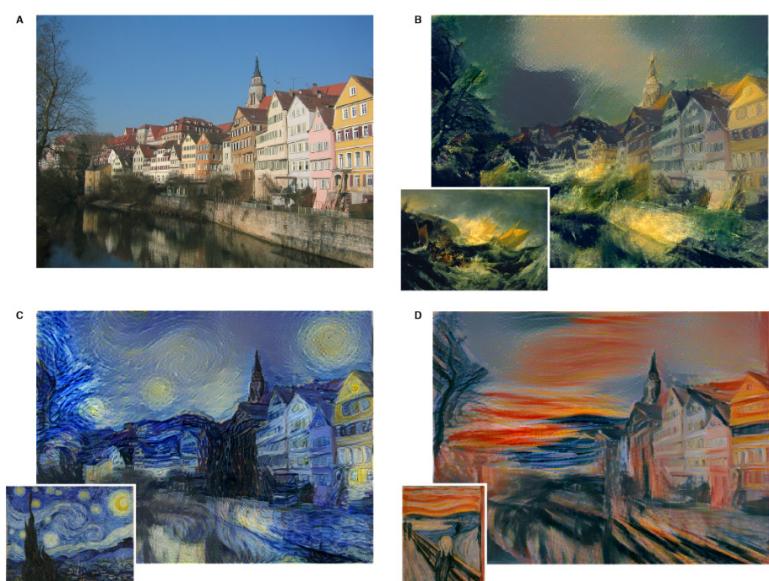
- Transfer Learning para mejora de imágenes



<https://github.com/alexjc/neural-enhance>

EJEMPLOS

- Transfer Learning para copiar el “estilo artístico”



<https://bit.ly/2N7xtqh>

EJEMPLOS

- Transfer Learning para emular “jugadores” de videojuegos

Artificial Intelligence

DeepMind has finally thrashed humans at StarCraft for real

In January, DeepMind's *Starcraft*-playing AI roundly beat humans players, but it had some serious advantages. Now it's learned to play fair

By ALEX LEE
Wednesday 30 October 2019

<https://www.wired.co.uk/article/deepmind-starcraft-alphastar>



Máster en Big Data y Data Science

Indexación, búsqueda y análisis en repositorios multimedia

60

REFERENCIAS

- Básico:
 - Tamoghna Ghosh, Raghav Bali, Dipanjan Sarkar, “Hands-On Transfer Learning with Python”, Packt 2018 (accessible desde el campus UAM en <https://learning.oreilly.com/library/view/hands-on-transfer-learning>)

Máster en Big Data y Data Science

Indexación, búsqueda y análisis en repositorios multimedia

61