

# Visión por Computador Profunda

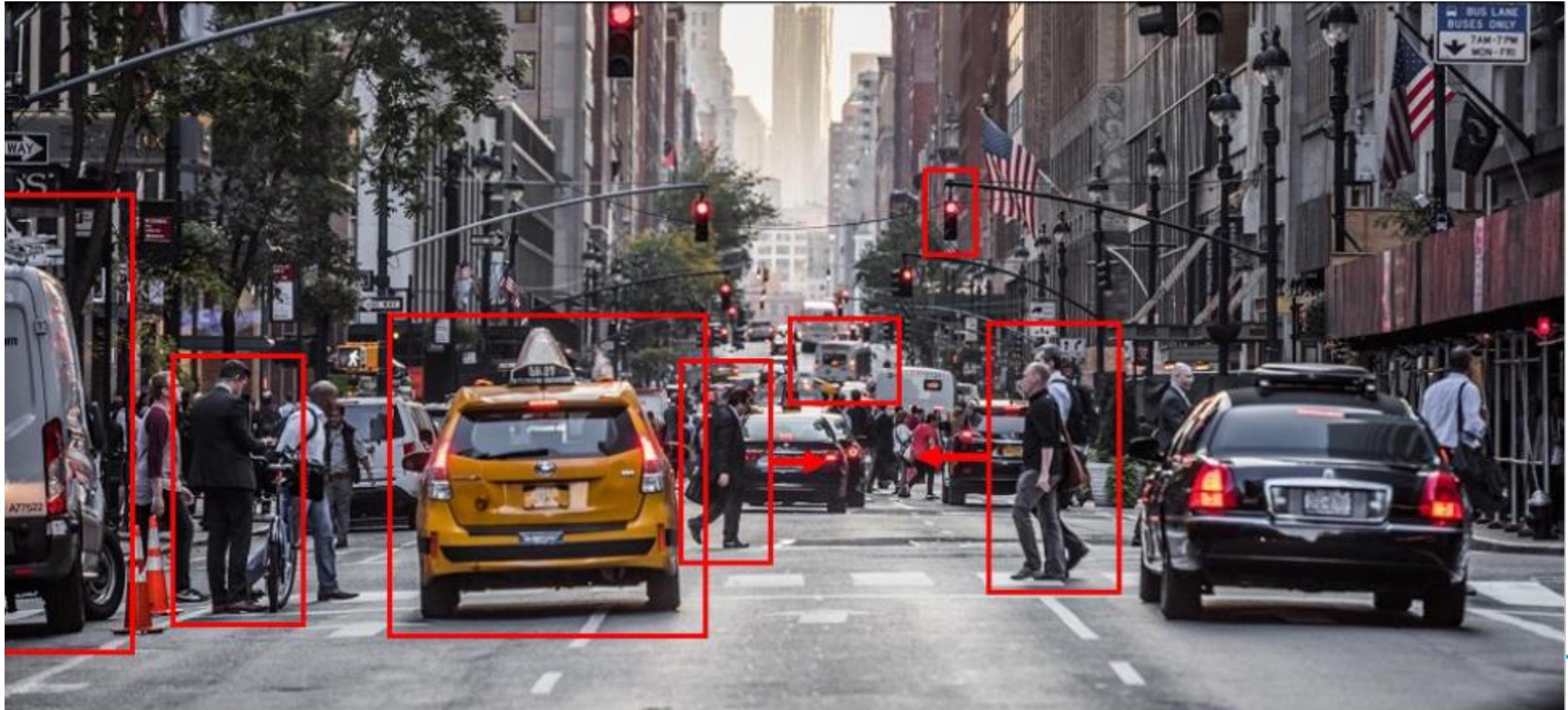




"Conocer qué hay y dónde está, observando."



Descubrir a partir de imágenes qué está presente en el mundo,  
dónde están las cosas, qué acciones están ocurriendo,  
y predecir y anticipar eventos en el mundo.



# El auge e impacto de la visión por computador

Robótica



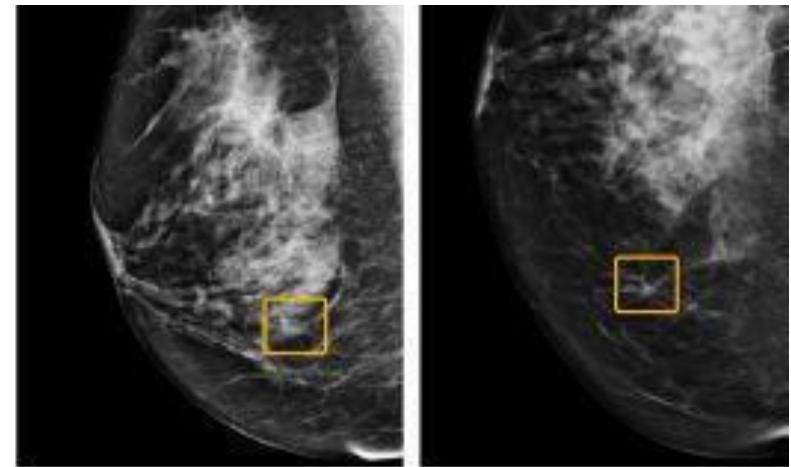
Accesibilidad



Computación móvil



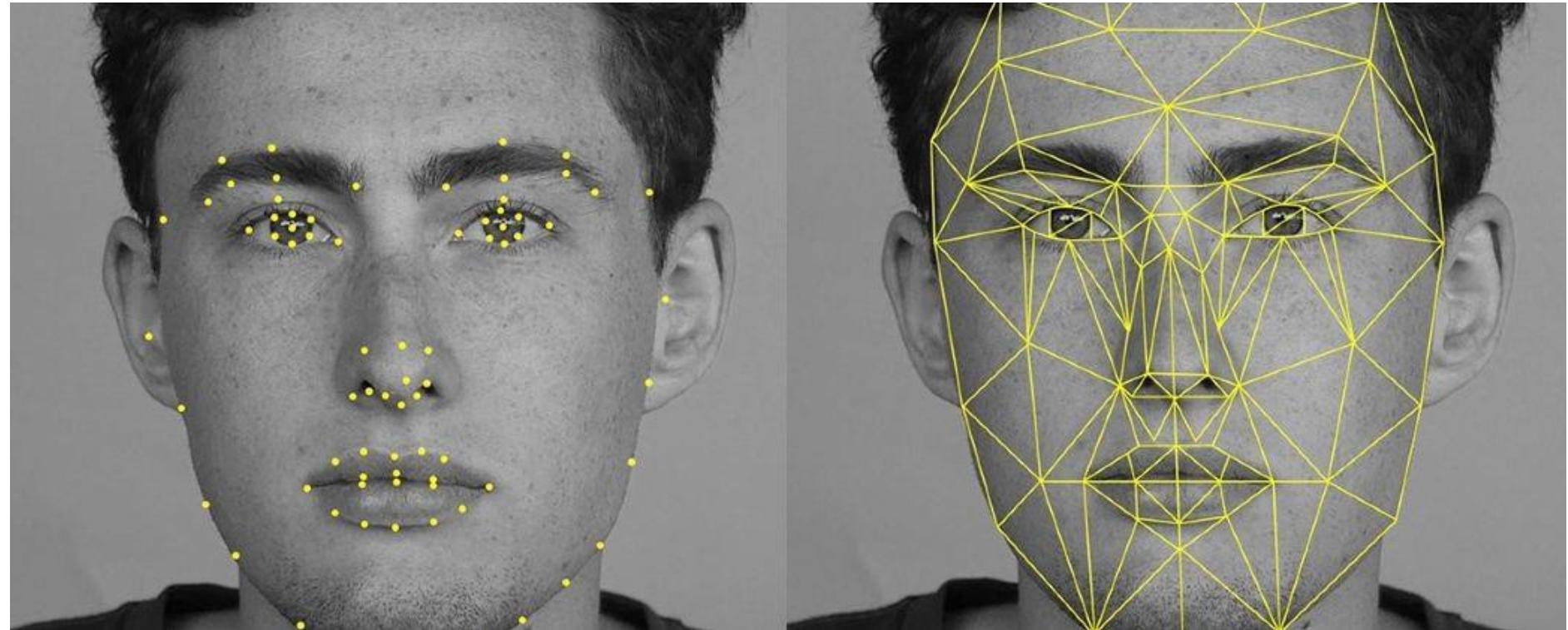
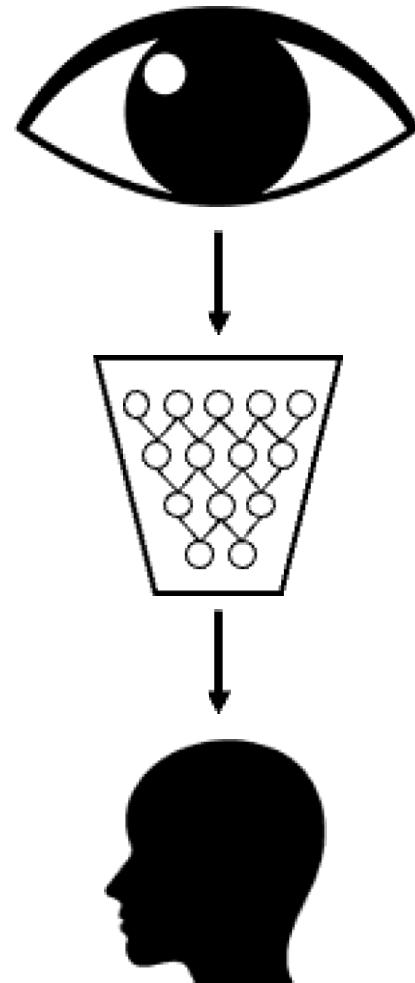
Biología y Medicina



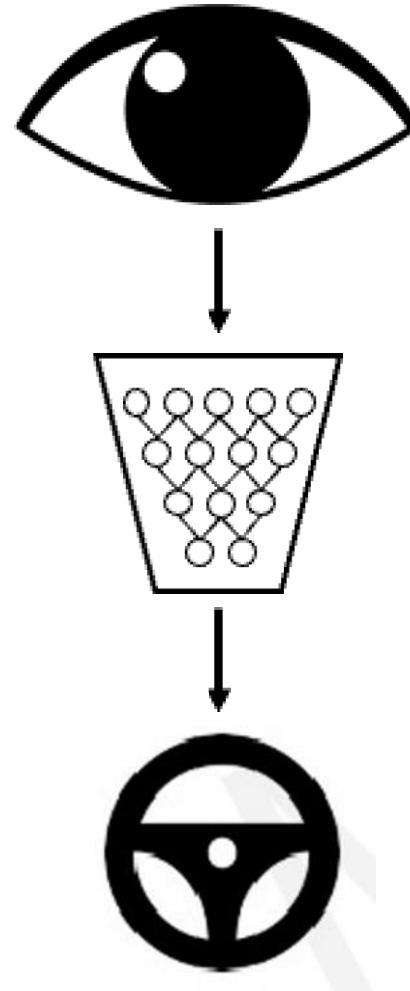
Conducción autónoma



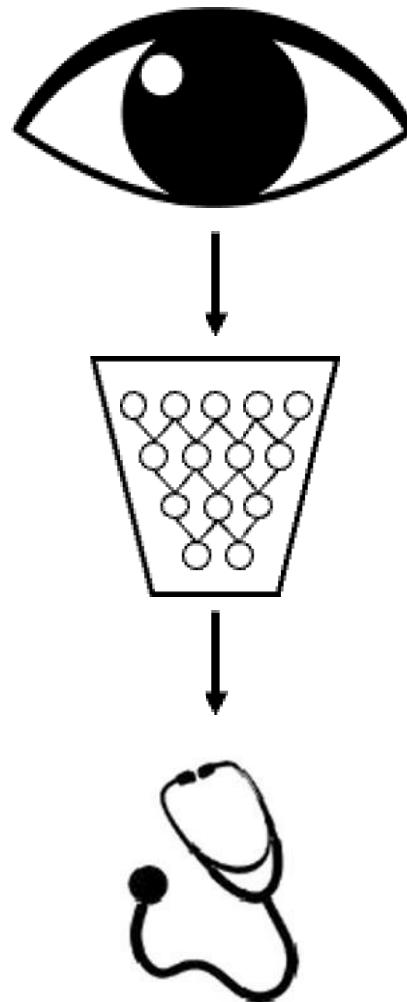
# Impacto: Detección y Reconocimiento Facial



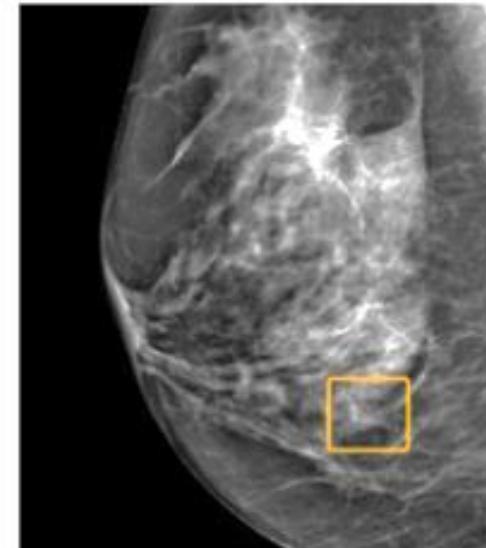
# Impacto: Vehículos Autónomos



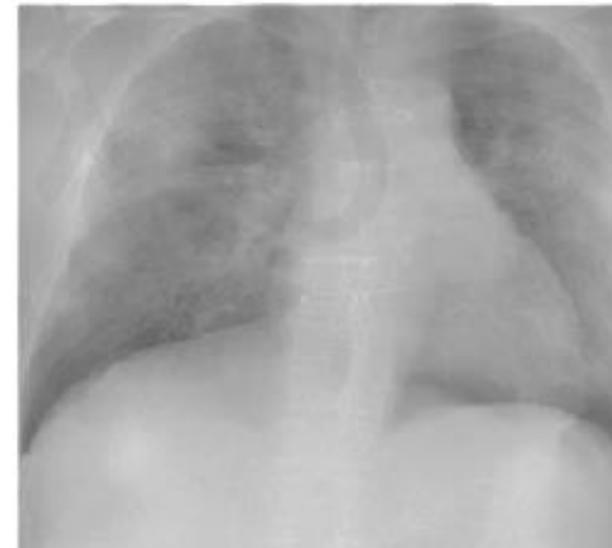
# Impacto: Medicina, Biología, Atención Sanitaria



Cáncer de Mama



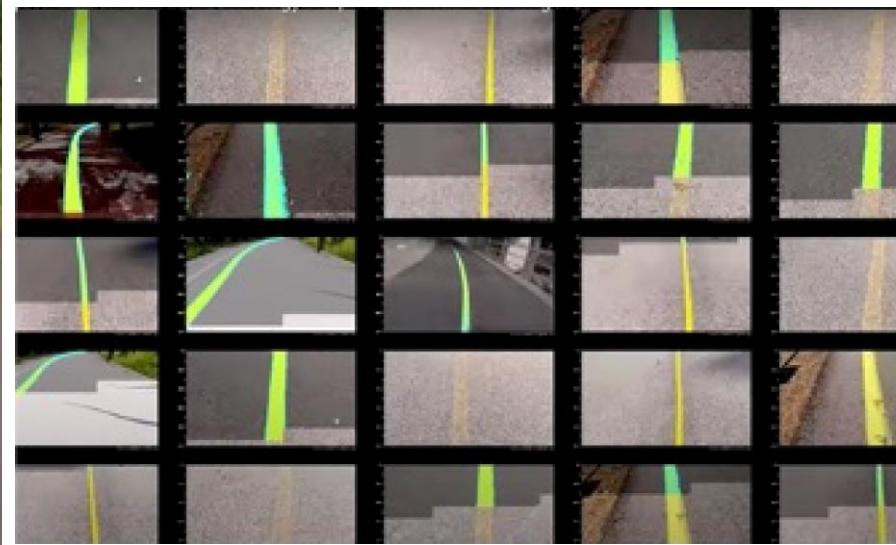
COVID 19



Cáncer de Piel



# Impacto: Accesibilidad



# Lo que las computadoras “ven”



# Las imágenes son números



157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	93	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	95	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	218
187	196	236	75	1	81	47	0	6	217	265	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Lo que las computadoras ven

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	93	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	16	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	103	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	218
187	196	236	75	1	81	47	0	6	217	265	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Una imagen es solo una matriz de números [0, 255], es decir,  
1080×1080×3 para una imagen RGB.



# Tareas en Visión por Computador



Imagen de entrada



147	153	174	148	166	142	179	151	173	147	166	146
186	162	183	78	79	82	83	17	118	218	186	194
180	180	59	14	94	6	10	33	48	198	185	181
206	199	1	134	131	111	120	204	184	16	96	180
194	68	137	26	291	246	239	239	237	87	71	201
172	105	207	238	238	214	220	239	238	98	74	206
188	88	179	208	186	216	211	168	138	76	36	149
189	97	165	84	18	166	134	11	30	63	22	148
199	188	191	198	198	227	178	143	182	198	38	190
205	174	188	202	236	231	149	179	228	49	98	234
190	216	116	149	206	187	86	189	79	39	216	241
180	234	147	108	237	210	127	102	36	107	206	234
160	214	173	66	188	142	96	63	2	108	266	216
187	196	236	76	1	43	47	0	6	217	266	211
189	202	237	148	0	0	12	108	208	198	243	236
190	206	123	207	177	121	123	200	171	13	96	216

Representación de píxeles

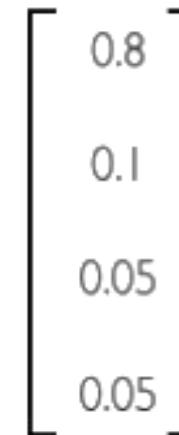
Clasificación

Lincoln

Washington

Jefferson

Obama



**Regresión:** La variable de salida toma un valor continuo.

**Clasificación:** La variable de salida toma una etiqueta de clase. Puede producir la probabilidad de pertenencia a una clase en particular.



# Detección de Características de Alto Nivel

Identifiquemos las características clave en cada categoría de imagen:



Nariz, Ojos



Ruedas, Placa



Puerta, Ventanas



# Extracción Manual de Características

Conocimiento del dominio

Definir características

Detectar características para clasificar

Problemas?



# Extracción Manual de Características

Conocimiento del dominio

Definir características

Detectar características para clasificar

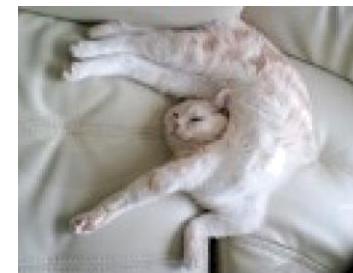
Variación de punto de vista



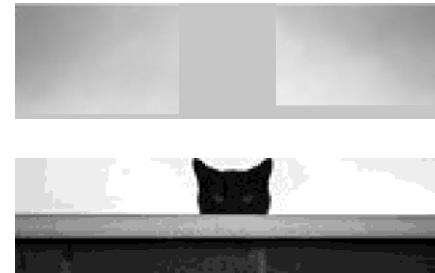
Variación de escala



Deformación



Oclusión



Condiciones de iluminación



Desorden de fondo



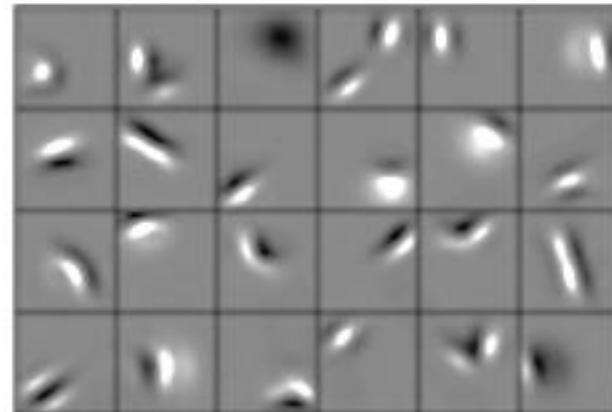
Variación intra-clase



# Aprendizaje de Representaciones de Características

Aprendemos una jerarquía de características directamente a partir de los datos en lugar de diseñarlas manualmente?

Características de bajo nivel:



Bordes, manchas oscuras

Características de nivel medio



Ojos, orejas, nariz

Características de alto nivel:



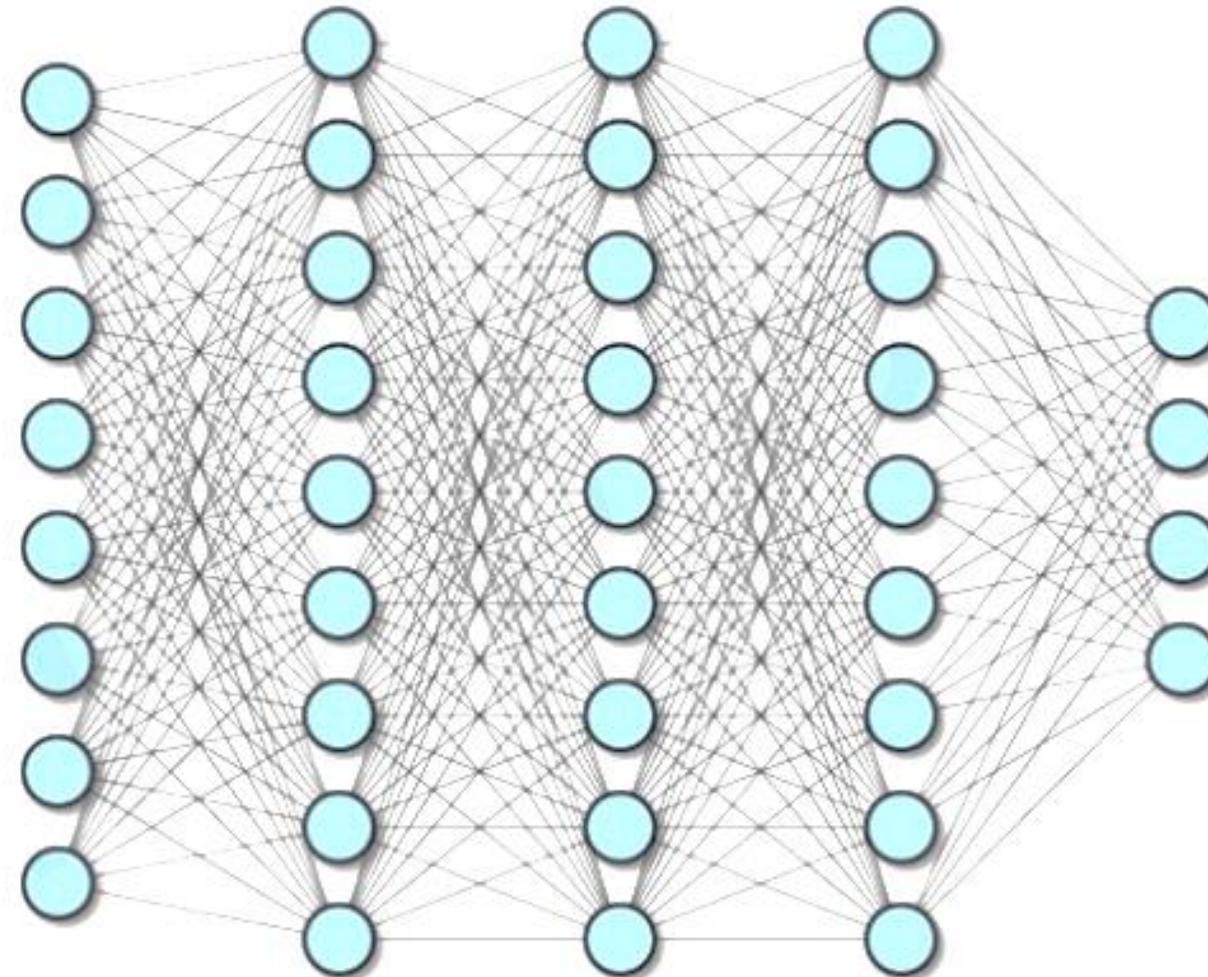
Estructura facial



# Aprendizaje de Características Visuales



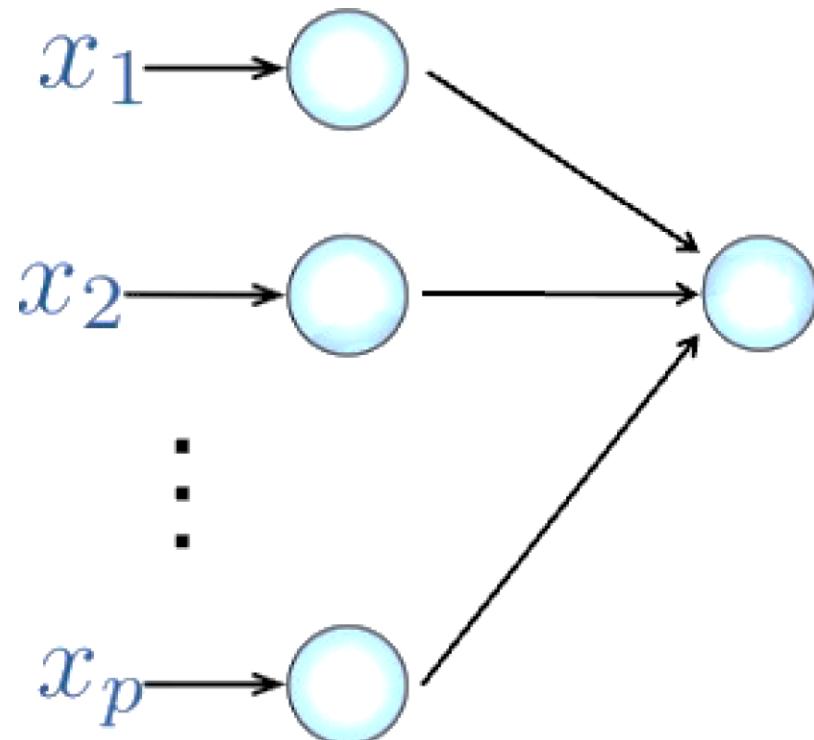
# Red Neuronal Totalmente Conectada



# Red Neuronal Totalmente Conectada

**Entrada:**

- Imagen 2D
- Vector de valores de píxeles



**Totalmente Conectada:**

- Conecta cada neurona en la capa oculta con todas las neuronas en la capa de entrada
- ¡Sin información espacial!
- ¡Y muchísimos parámetros!

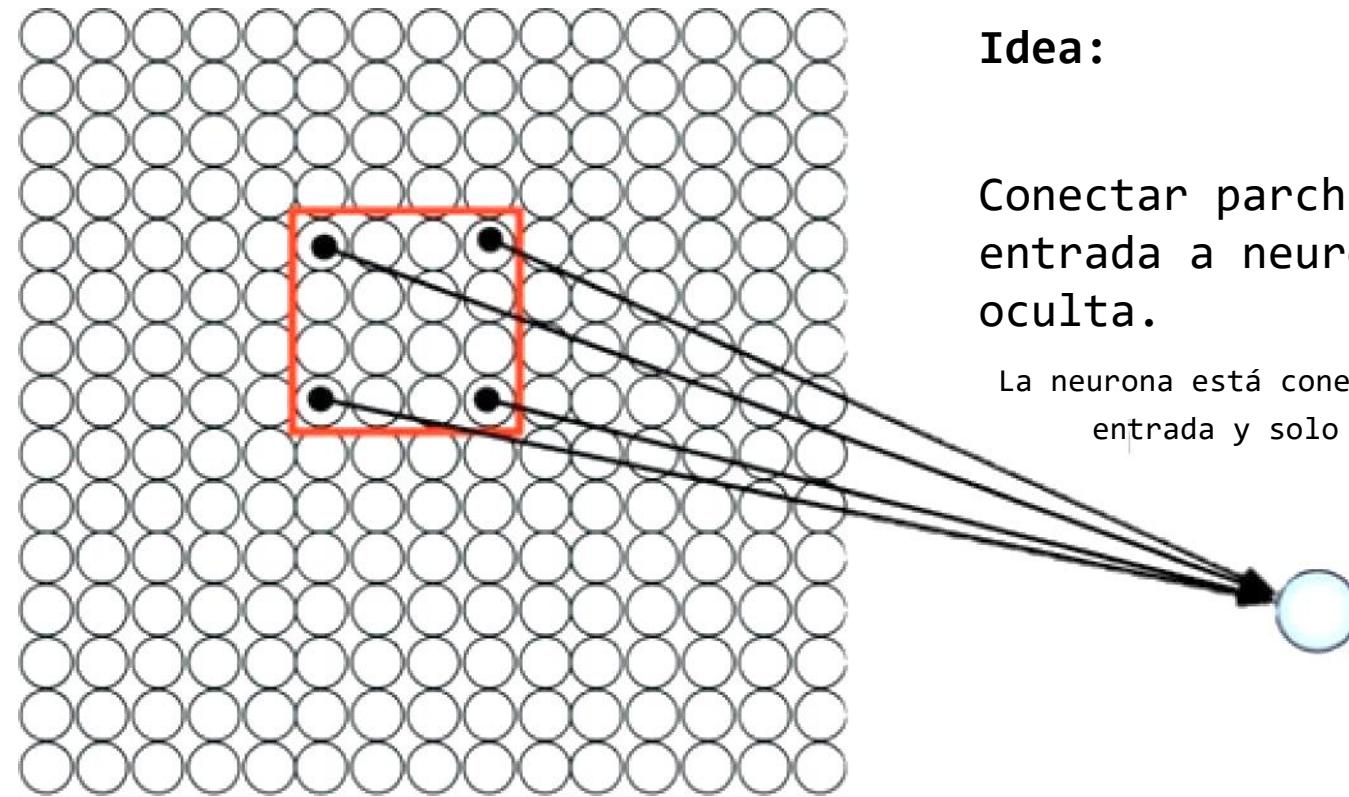
¿Cómo podemos usar la estructura espacial en la entrada para informar la arquitectura de la red?



# Usando la Estructura Espacial

**Entrada:**

Imagen 2D,  
matriz de  
valores de  
píxeles



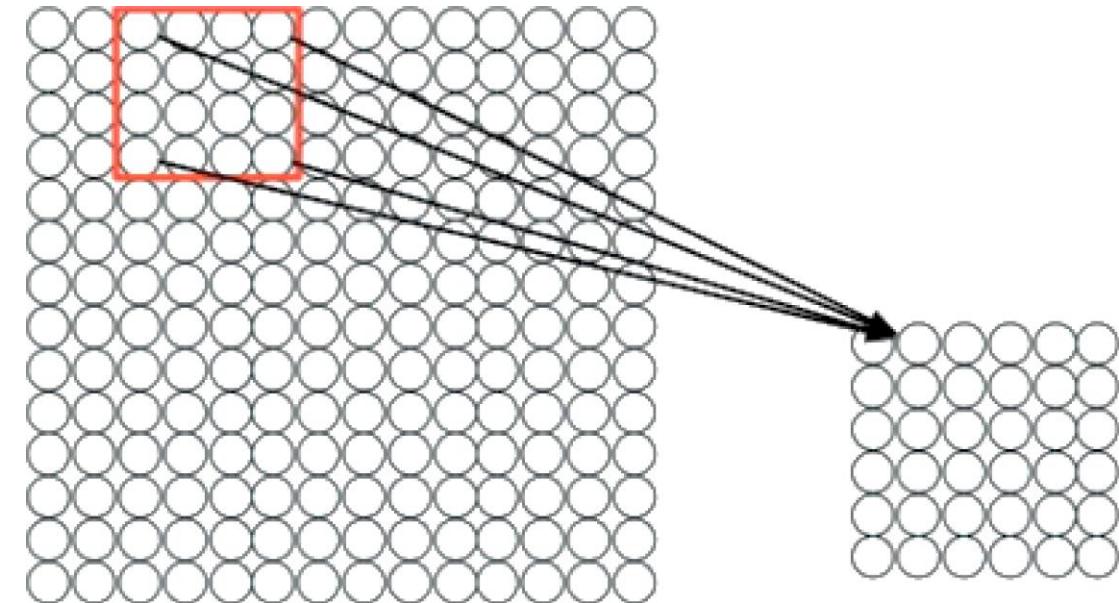
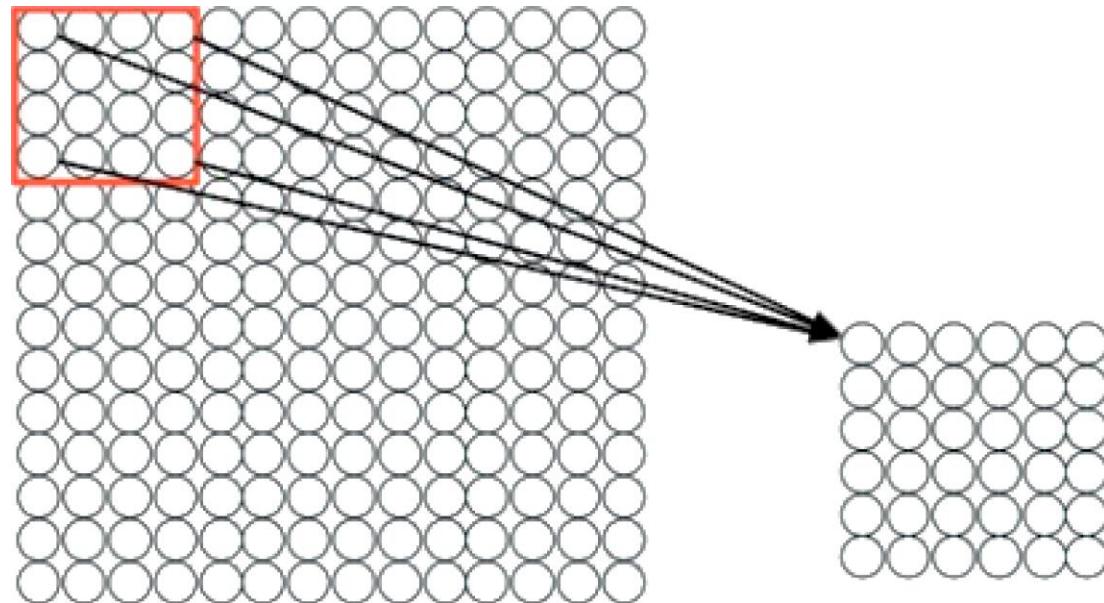
**Idea:**

Conectar parches de la  
entrada a neuronas en la capa  
oculta.

La neurona está conectada a una región de la  
entrada y solo “ve” esos valores.



# Usando la Estructura Espacial



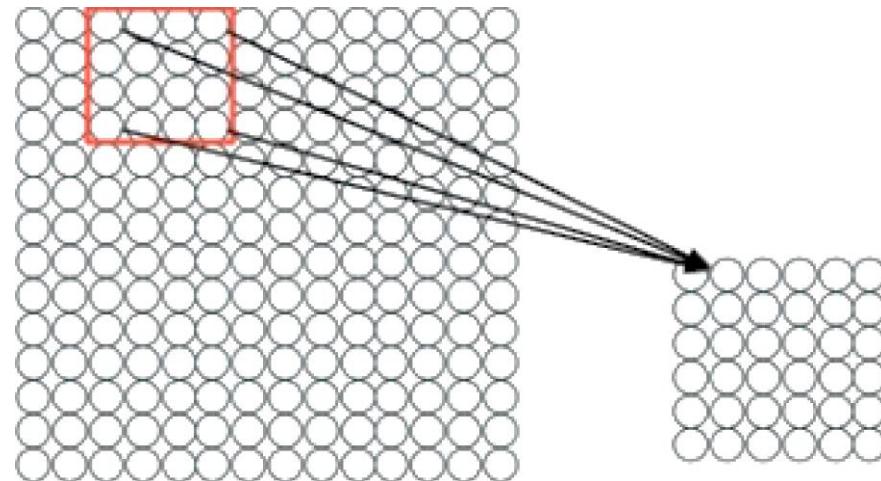
Coneectar un parche en la capa de entrada a una única neurona en la capa siguiente.

Usar una ventana deslizante para definir las conexiones.

¿Cómo podemos ponderar el parche para detectar características particulares?



# Extracción de Características con Convolución



- Filtro de tamaño  $4 \times 4$ : 16 pesos diferentes
- Aplicar este mismo filtro a parches de  $4 \times 4$  en la entrada
- Desplazar 2 píxeles para el siguiente parche

Esta operación por parches es **convolución**

1. Aplicar un conjunto de pesos –un filtro– para extraer características locales
2. Usar múltiples filtros para extraer diferentes características
3. Compartir espacialmente los parámetros de cada filtro



# Extracción de Características y Convolución

## Un Estudio de Caso



# X ó X?

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

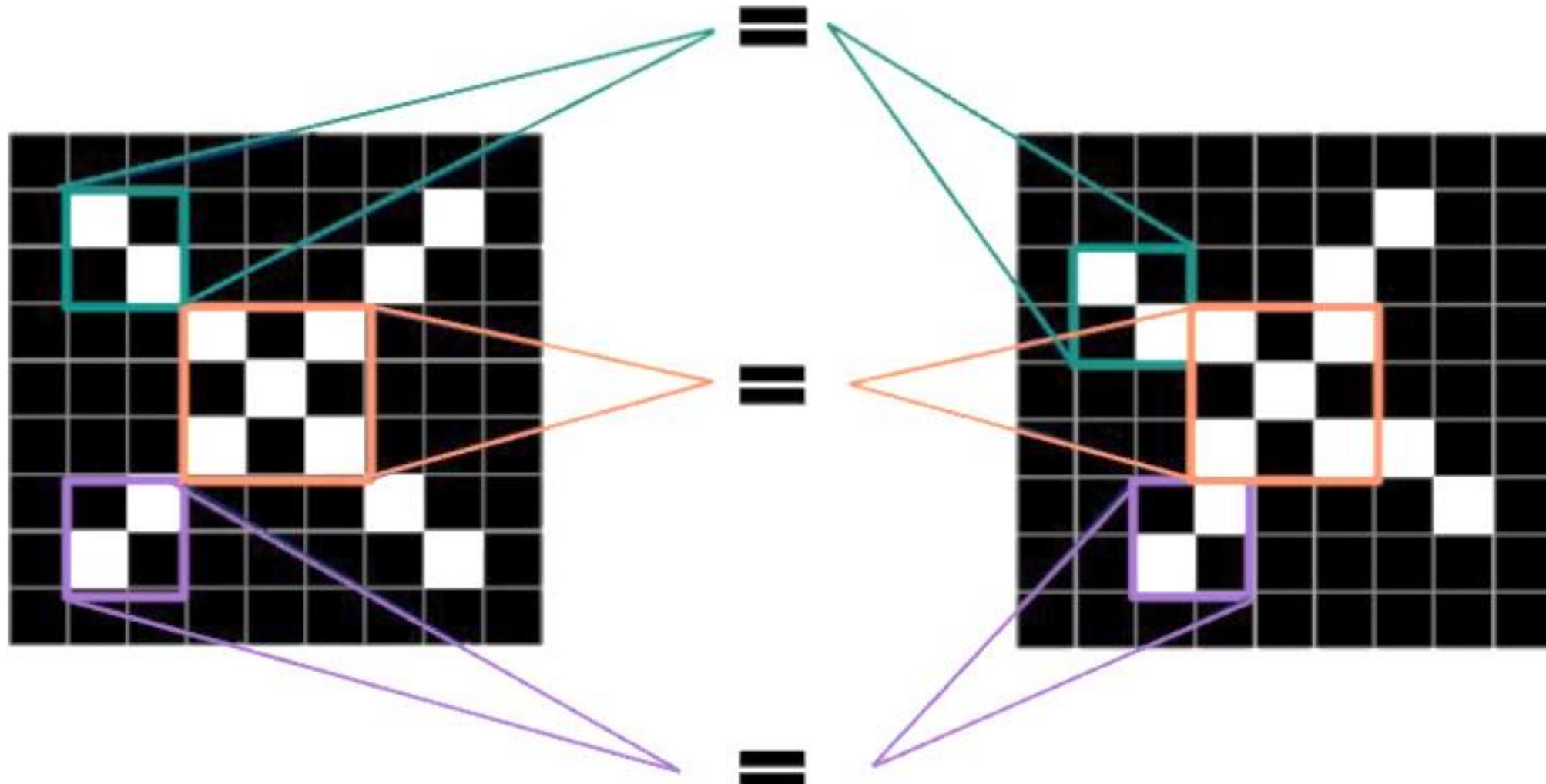


-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	1	-1	-1	-1	-1	-1
-1	-1	1	1	-1	-1	1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	1	-1
-1	-1	-1	1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

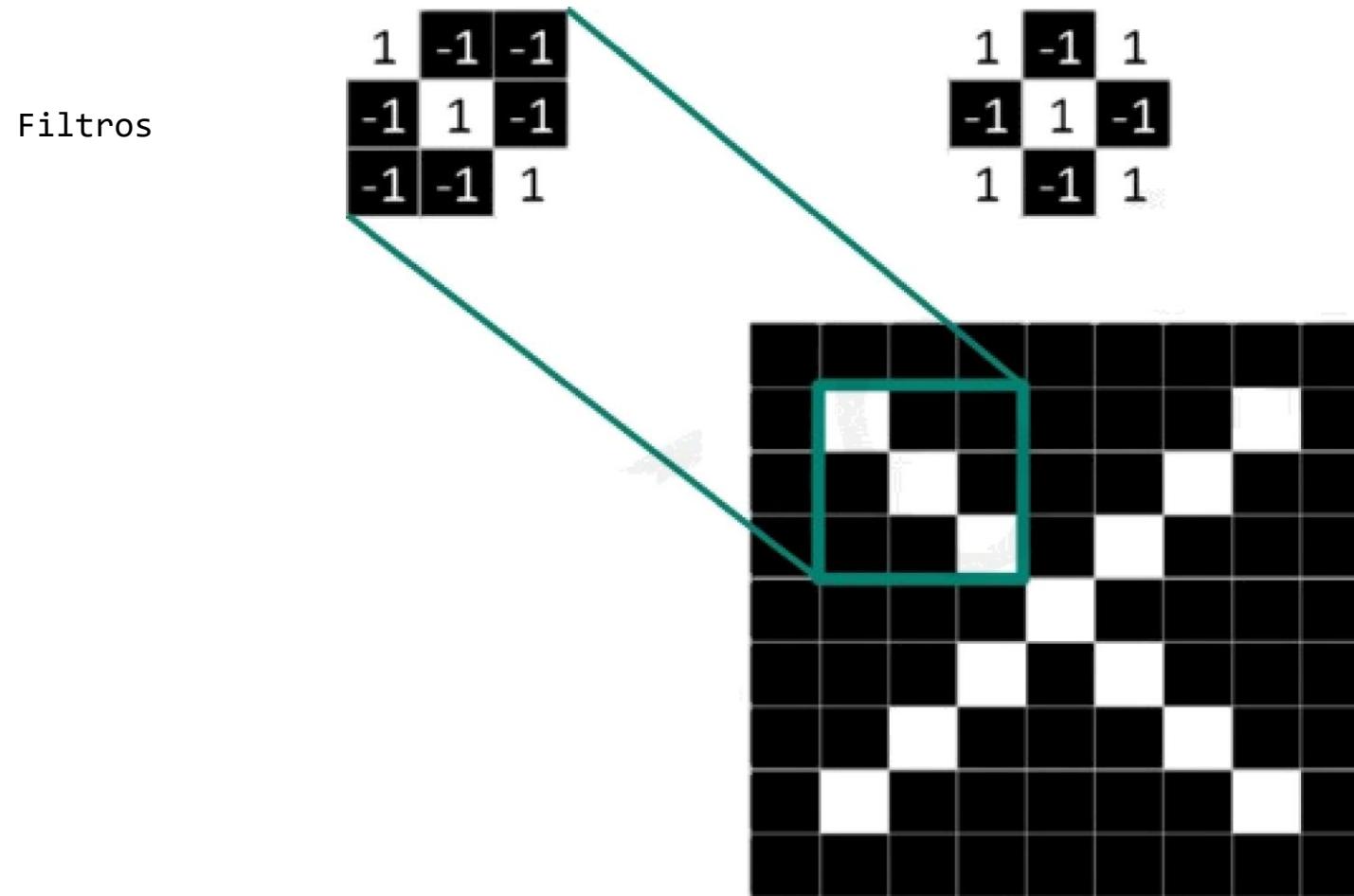
La imagen se representa como una matriz de valores de píxeles... ¡y las computadoras son literales!  
Queremos poder clasificar una X como una X incluso si está desplazada, reducida, rotada o deformada.



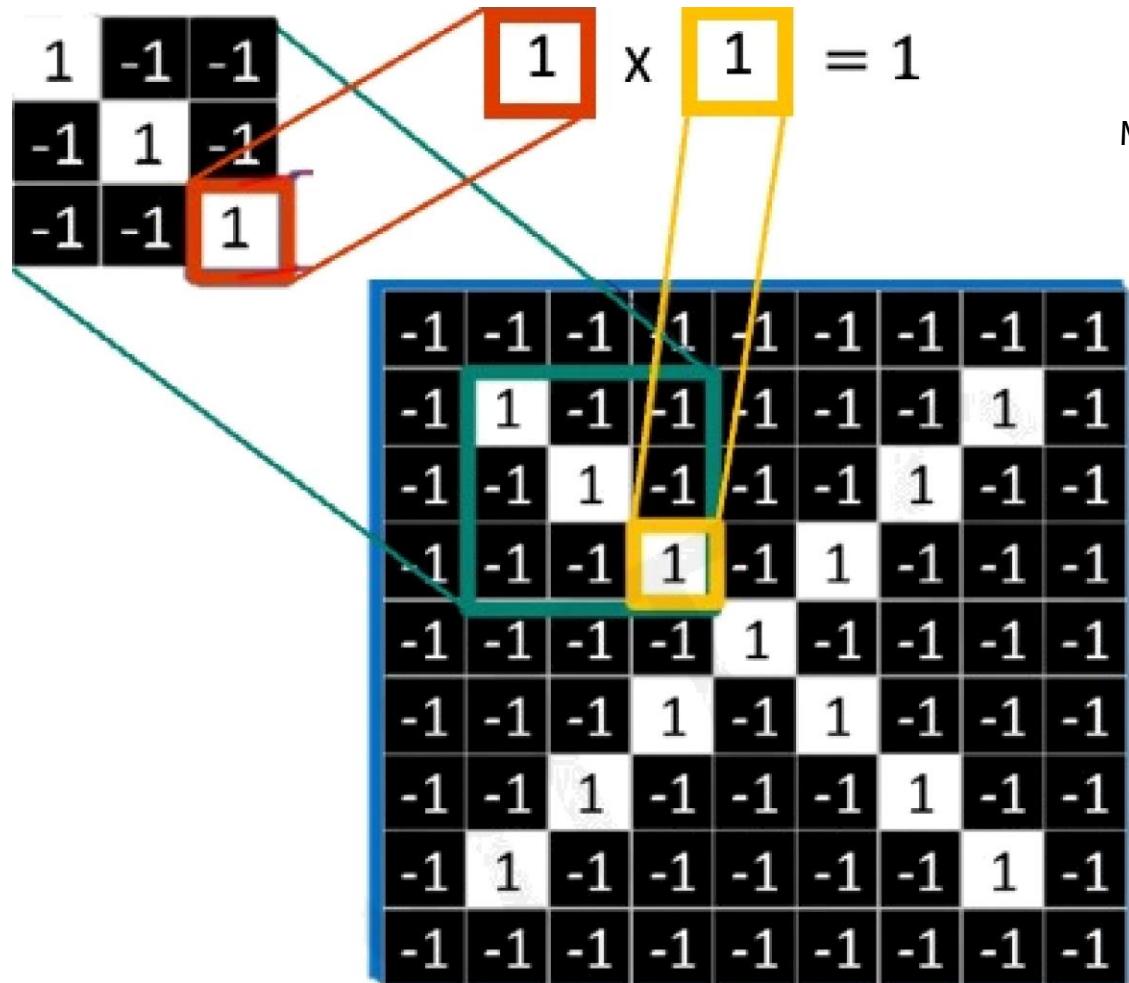
# Características de la X



# Filtros para Detectar Características de la X



# La Operación de Convolución



Sumar salidas

1	1	1
1	1	1
1	1	1

— 9



# La Operación de Convolución

Supongamos que queremos calcular la convolución de una imagen de  $5 \times 5$  y un filtro de  $3 \times 3$ :

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Imagen



1	0	1
0	1	0
1	0	1

Filtro

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas.



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

Filtro



4		

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

Filtro



4	3	

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	$\times 1$	0	$\times 0$	0	$\times 1$
0	1	1	$\times 0$	1	$\times 1$	0	$\times 0$
0	0	1	$\times 1$	1	$\times 0$	1	$\times 1$
0	0	1	1	1	0		
0	1	1	0	0			



1	0	1
0	1	0
1	0	1

Filtro



4	3	4

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	1	0
0 <sub>x0</sub>	0 <sub>x1</sub>	1 <sub>x0</sub>	1	1
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2		

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0
0	0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1
0	0 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2	4	

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	$1 \times 1$	$1 \times 0$
0	0	1	$\times 0$	$1 \times 1$
0	0	1	$\times 1$	$1 \times 0$
0	1	1	0	$\times 1$



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2	4	3

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2	4	3
2		

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	1	0
0	0 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1
0	0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	0
0	1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	0 <small><math>\times 1</math></small>	0



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2	4	3
2	3	

Mapa de  
características



# La Operación de Convolución

Deslizamos el filtro  $3 \times 3$  sobre la imagen de entrada, multiplicamos elemento a elemento y sumamos las salidas:

1	1	1	0	0
0	1	1	1	0
0	0	1	$\times 1$	$\times 0$
0	0	1	$\times 0$	$\times 1$
0	1	1	$\times 1$	$\times 0$



1	0	1
0	1	0
1	0	1

Filtro



4	3	4
2	4	3
2	3	4

Mapa de  
características



# Producción de Mapas de Características



Original



Afilado



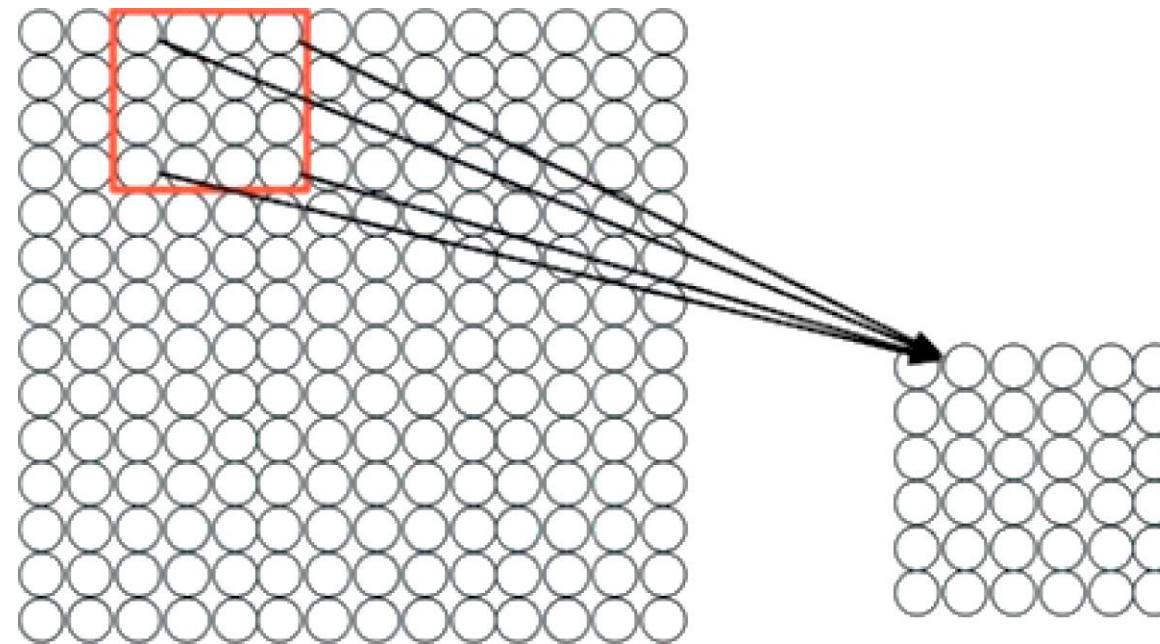
Detección de  
bordes



Detección de  
bordes “fuerte”



# Extracción de Características con Convolución



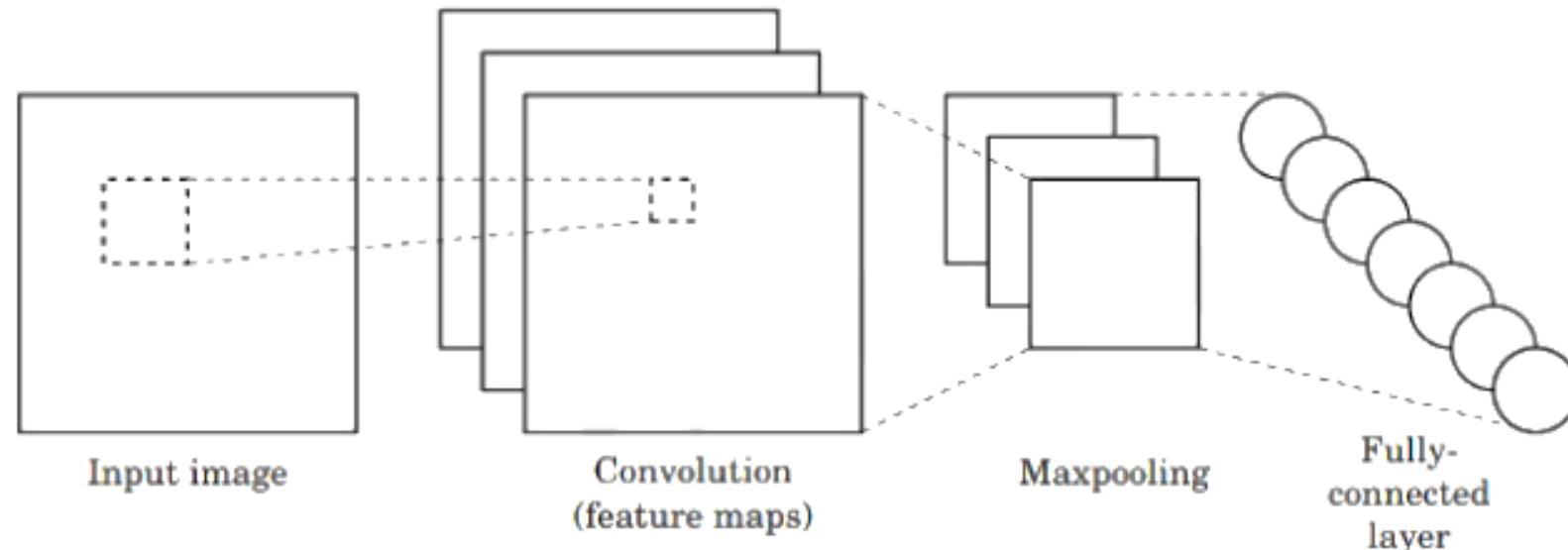
1. Aplicar un conjunto de pesos –un filtro– para extraer características locales
2. Usar múltiples filtros para extraer diferentes características
3. Compartir espacialmente los parámetros de cada filtro.



# Redes Neuronales Convolucionales (CNNs)



# CNNs para Clasificación



`tf.keras.layers.Conv2D`

`tf.keras.activations.*`

`tf.keras.layers.MaxPool2D`

1. **Convolución:** Aplicar filtros para generar mapas de características
2. **No linealidad:** A menudo ReLU.
3. **Pooling:** Operación de reducción de muestreo en cada mapa de características.

`torch.nn.Conv2d`

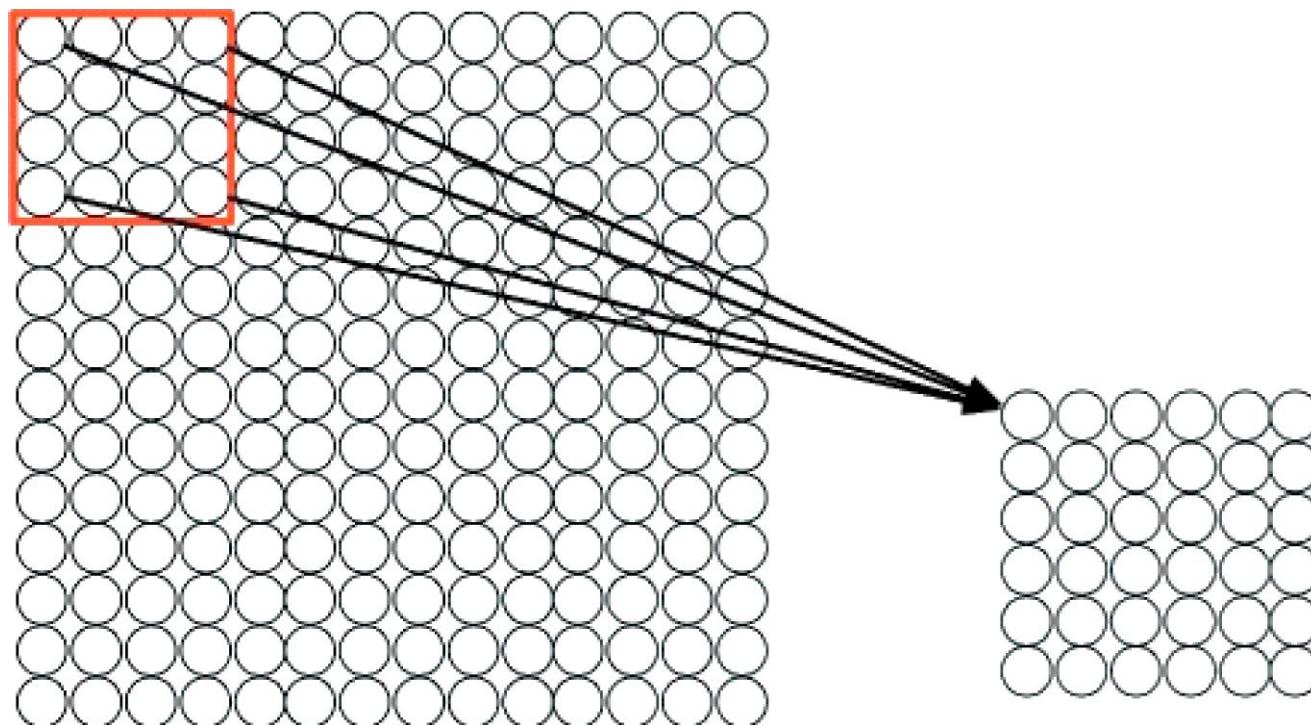
`torch.nn.ReLU,...`

`torch.nn.MaxPool2d`



Entrenar el modelo con datos de imagen.  
Aprender los pesos de los filtros en las capas convolucionales.

# Capas Convolucionales – Conectividad Local



$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

para la neurona  $(p,q)$  en la capa oculta

 `tf.keras.layers.Conv2D`

 `torch.nn.Conv2d`

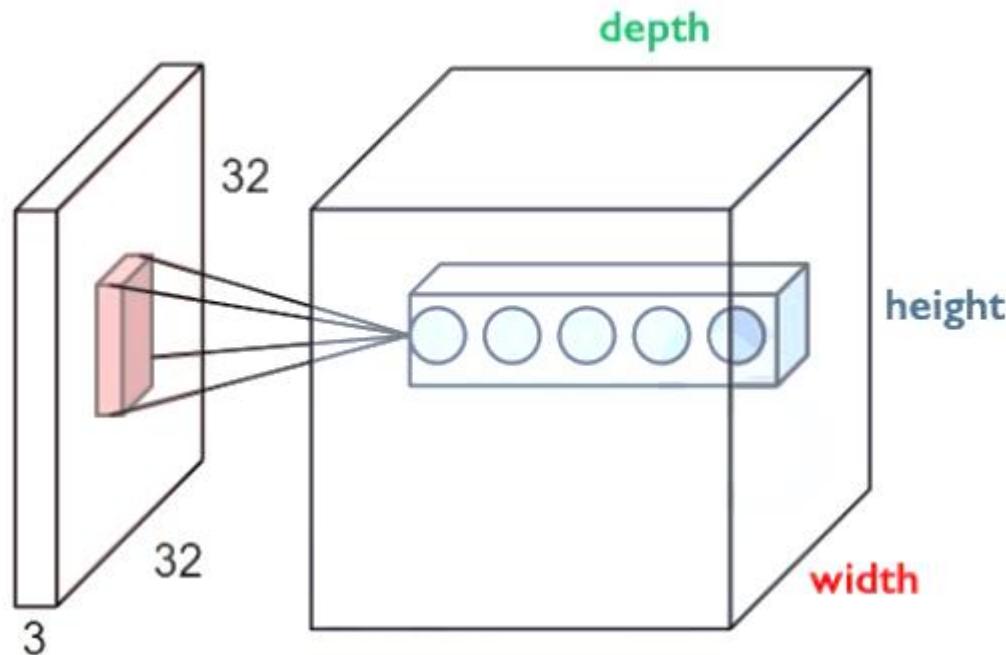
Para una neurona en la capa oculta:

- Tomar entradas de un parche
- Calcular la suma ponderada
- Aplicar sesgo (bias)

1. Aplicar una ventana de pesos
2. Calcular combinaciones lineales
3. Activar con una función no lineal)



# CNNs: Disposición Espacial del Volumen de Salida



Dimensiones de la capa:

$h \times w \times d$

donde  $h$  y  $w$  son dimensiones espaciales y  $d$  (profundidad) = número de filtros

**Stride:**

Tamaño del paso del filtro

**Campo Receptivo:**

Ubicaciones en la imagen de entrada a las que una neurona está conectada



```
tf.keras.layers.Conv2D( filters=d, kernel_size=(h,w), strides=s )
```



```
torch.nn.Conv2d( in_channels=3, out_channels=d, kernel_size=(h,w), stride=s )
```

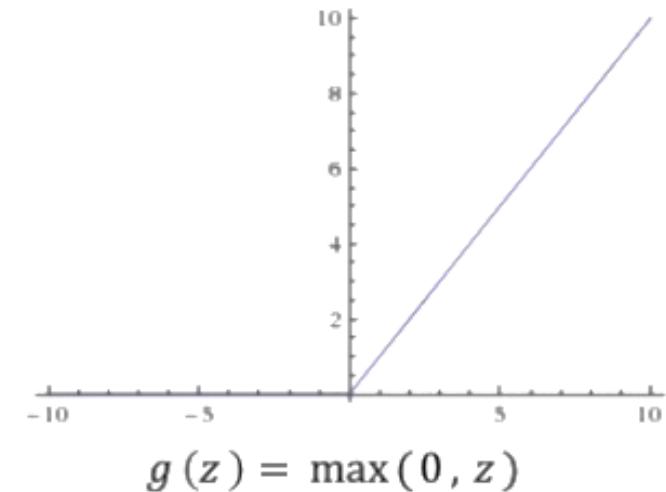
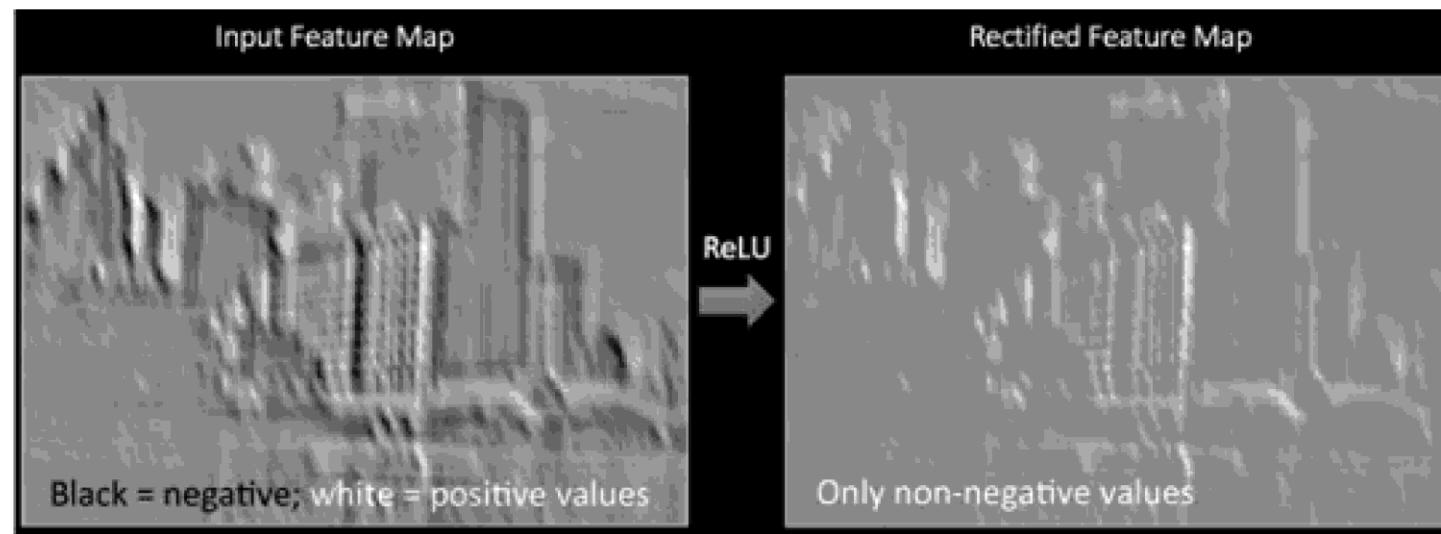
¡Es necesario especificar la dimensionalidad de la entrada!



# Introduciendo No Linealidad

- Aplicar después de cada operación de convolución (es decir, después de las capas convolucionales)
- ReLU: operación pixel a pixel que reemplaza todos los valores negativos por cero. ¡Operación no lineal!

Rectified Linear Unit (ReLU)

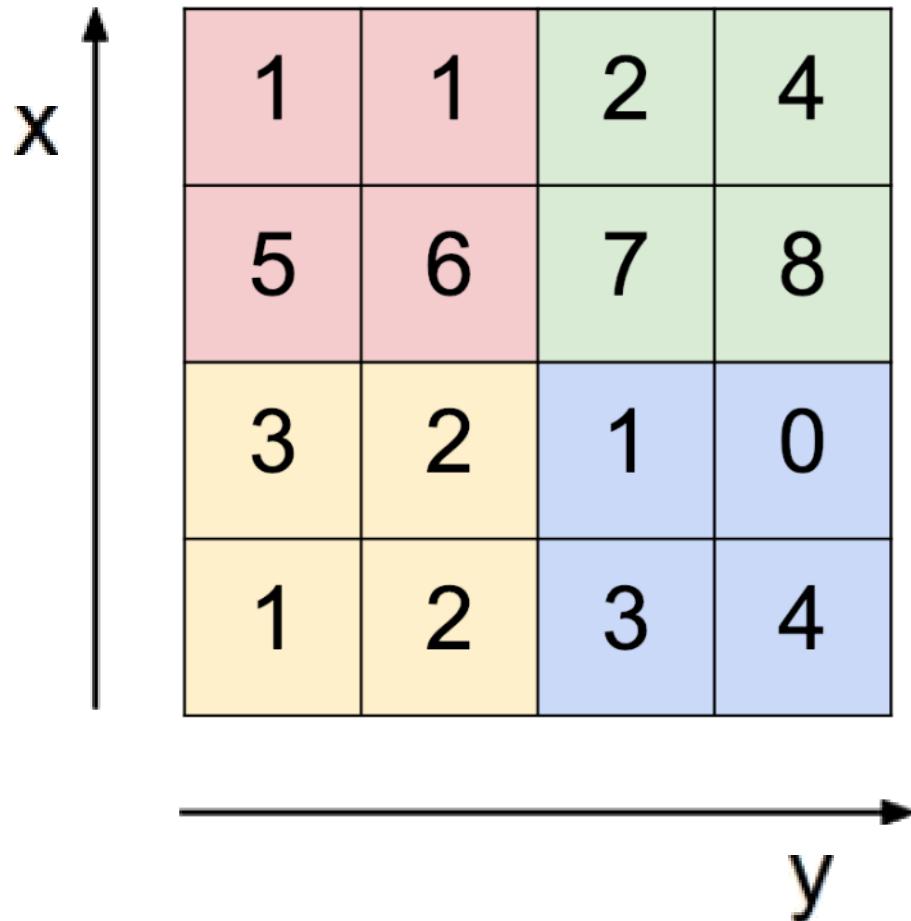


 `tf.keras.layers.ReLU`

 `torch.nn.ReLU`



# Pooling



Max pooling con filtros de  
2×2 y stride 2

tf.keras.layers.MaxPool2D(  
pool\_size=(2, 2),  
strides=2  
)



torch.nn.MaxPool2d(  
kernel\_size=(2, 2),  
stride=2  
)



6	8
3	4

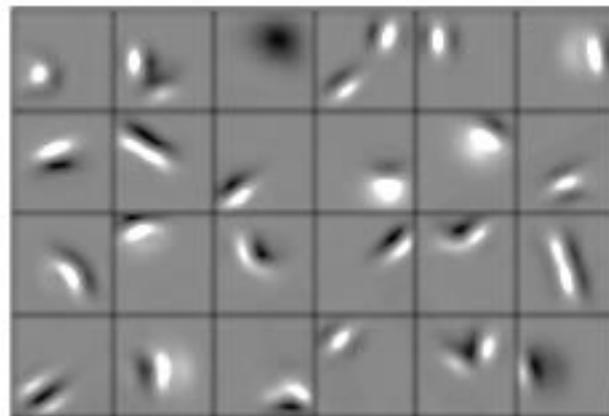
1. Reducción de dimensionalidad
2. Invariancia espacial

¿Cómo más podemos reducir la resolución y preservar la invariancia espacial?



# Aprendizaje de Representaciones en CNNs Profundas

Características de bajo nivel:



Bordes, manchas oscuras

Conv Layer 1

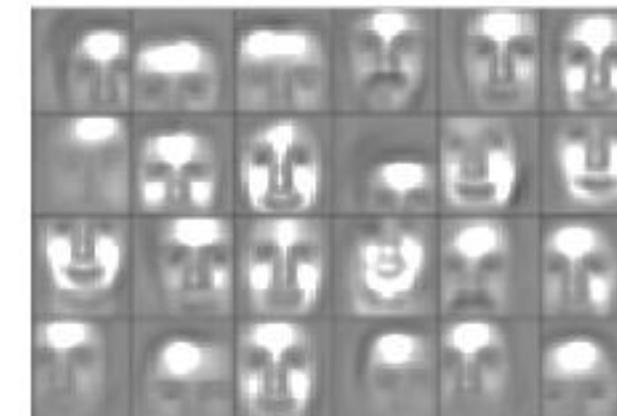
Características de nivel medio



Ojos, orejas, nariz

Conv Layer 2

Características de alto nivel:

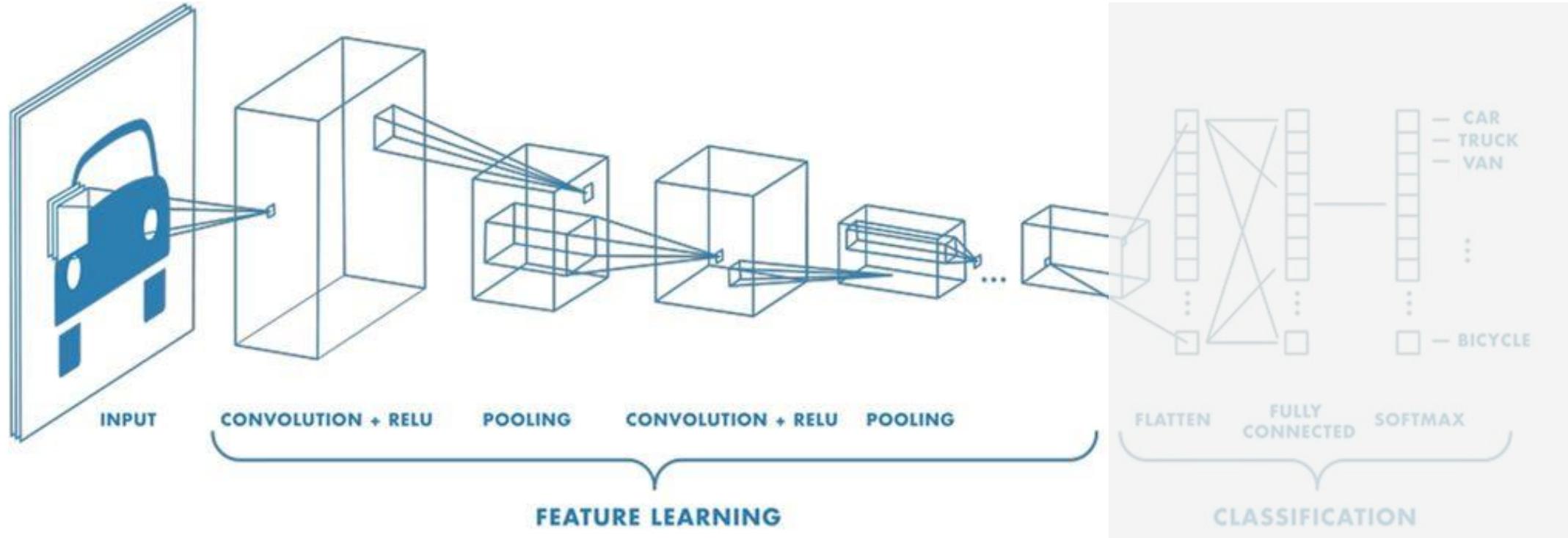


Estructura facial

Conv Layer 3



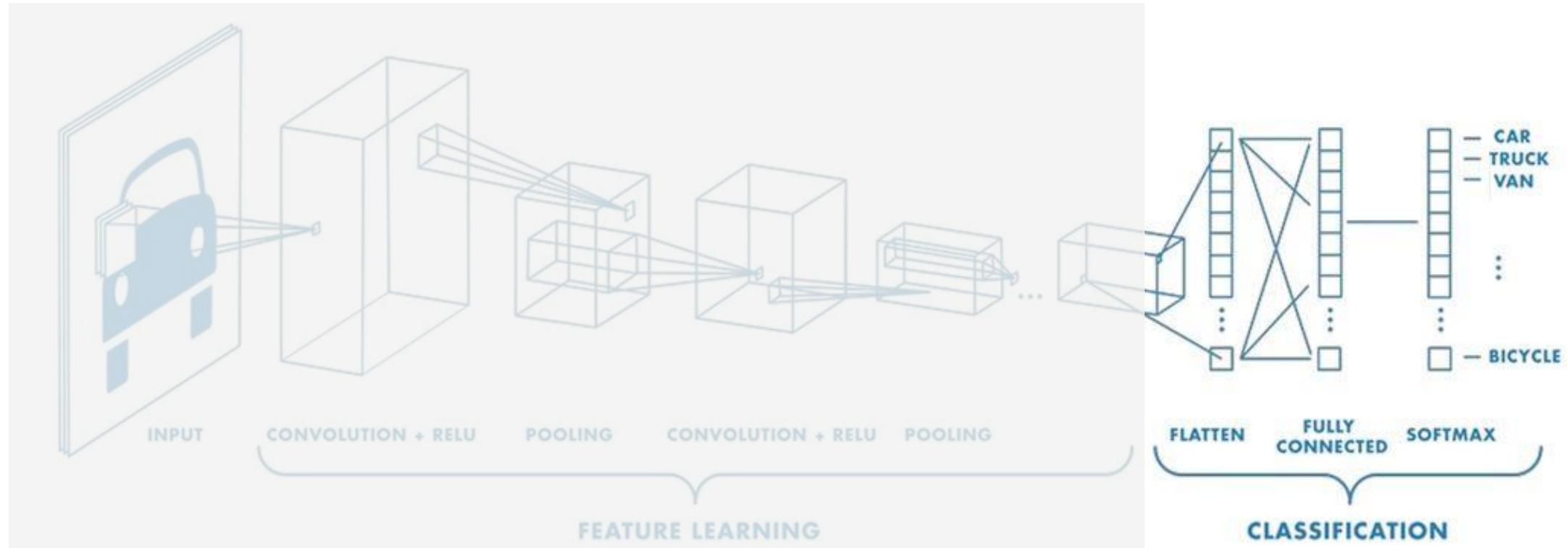
# CNNs para Clasificación: Aprendizaje de Características



1. Aprender características en la imagen de entrada mediante convolución
2. Introducir no linealidad mediante una función de activación (¡los datos del mundo real no son lineales!)
3. Reducir la dimensionalidad y preservar la invariancia espacial con pooling



# CNNs para Clasificación: Probabilidades de Clase



- Las capas CONV y POOL generan características de alto nivel de la entrada.
- La capa totalmente conectada usa estas características para clasificar la imagen de entrada.
- El resultado se expresa como la probabilidad de que la imagen pertenezca a una clase en particular

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$



# Poniéndolo todo junto: CNN en TensorFlow



```
import tensorflow as tf

def generate_model():
    model = tf.keras.Sequential([
        # first convolutional layer
        tf.keras.layers.Conv2D(32, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),

        # second convolutional layer
        tf.keras.layers.Conv2D(64, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),

        # fully connected classifier
        tf.keras.layers.Flatten(),
        tf.keras.layers.Dense(1024, activation='relu'),
        tf.keras.layers.Dense(10, activation='softmax')  # 10 outputs
    ])
    return model
```



# Poniéndolo todo junto: CNN en PyTorch



```
import torch

def generate_model():
    model = nn.Sequential([
        # first and second convolutional layer
        torch.nn.Conv2d(in_channels=3, out_channels=32, filter_size=3),
        torch.nn.ReLU(),
        torch.nn.MaxPool2d(kernel_size=2, stride=2),

        torch.nn.Conv2d(in_channels=32, out_channels=64, filter_size=3),
        torch.nn.ReLU(),
        torch.nn.MaxPool2d(kernel_size=2, stride=2),

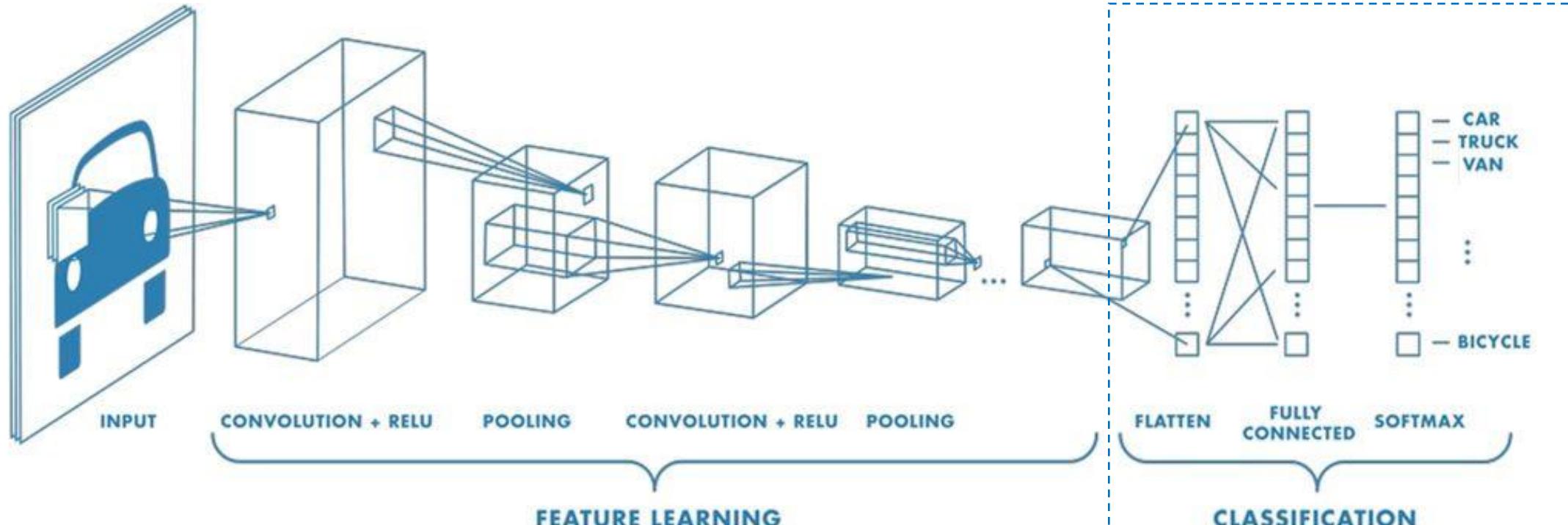
        # fully connected classifier
        torch.nn.Flatten(),
        torch.nn.Linear(64*6*6, 1024), # flattened dim after 2 conv layers
        torch.nn.ReLU(),
        torch.nn.Linear(1024, 10), # 10 outputs
    ])
    return model
```



# Una Arquitectura para Muchas Aplicaciones



# Una Arquitectura para Muchas Aplicaciones

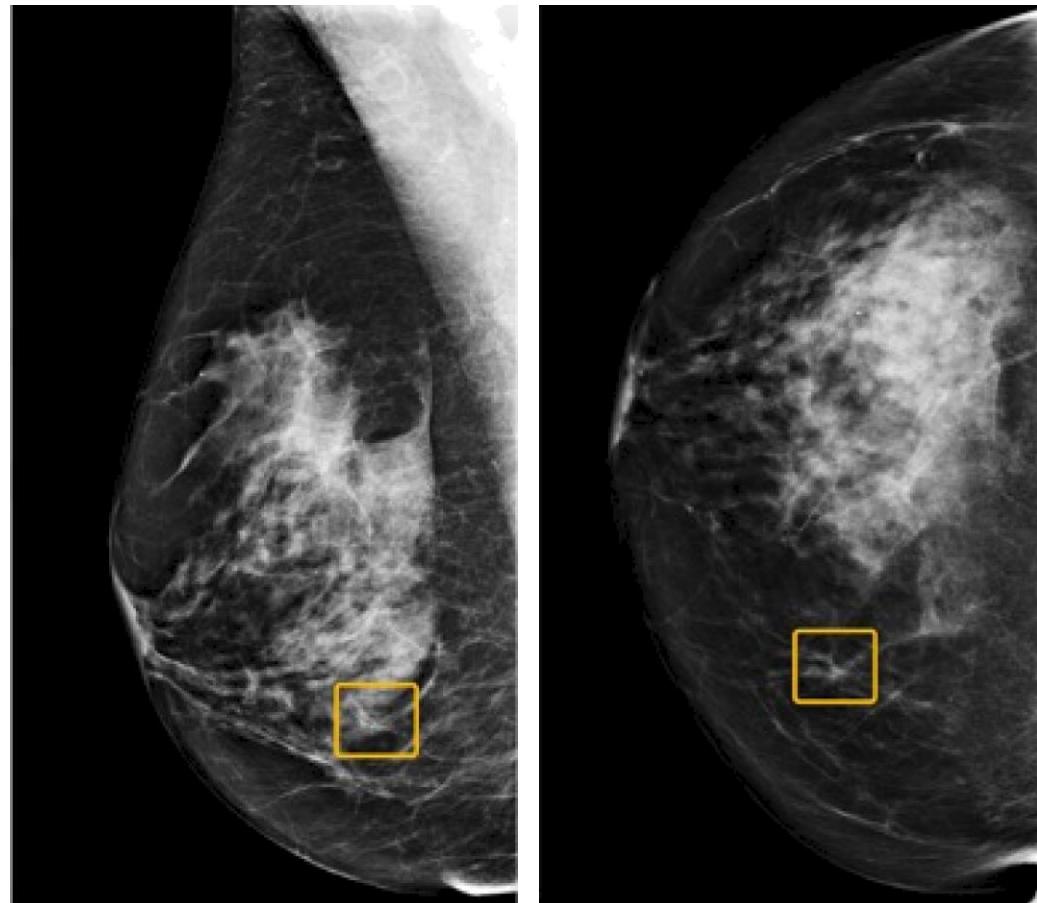
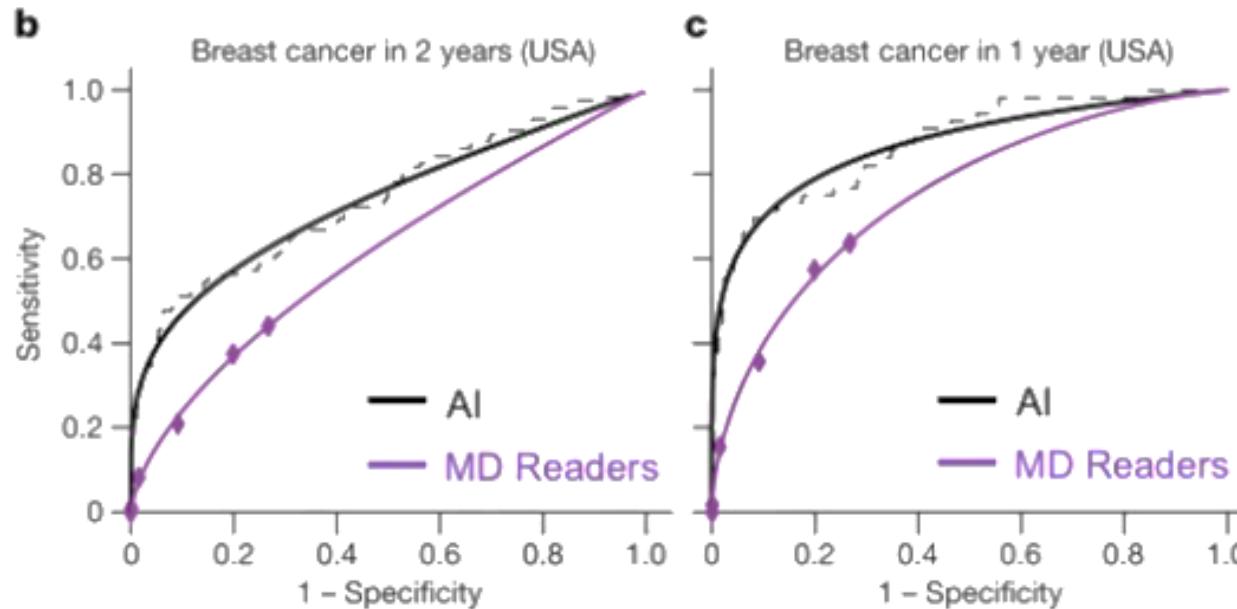


Clasificación  
Detección de objetos  
Segmentación  
Control probabilístico



# Clasificación: Detección de Cáncer de Mama

Evaluación internacional de un sistema de IA para la detección de cáncer de mama



El sistema basado en CNN superó a radiólogos expertos en la detección de cáncer de mama a partir de mamografías.

Caso de cáncer de mama no detectado por el radiólogo pero sí detectado por la IA.



# Detección de Objetos

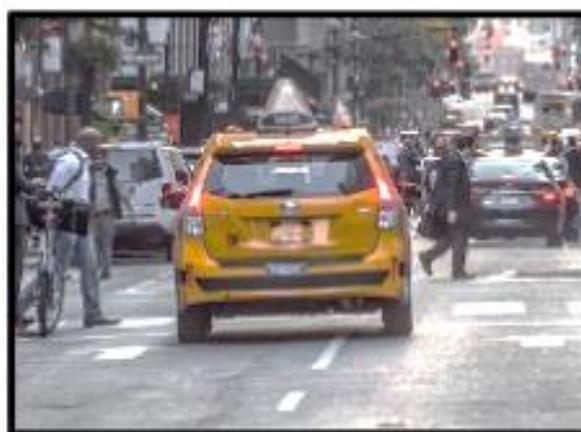
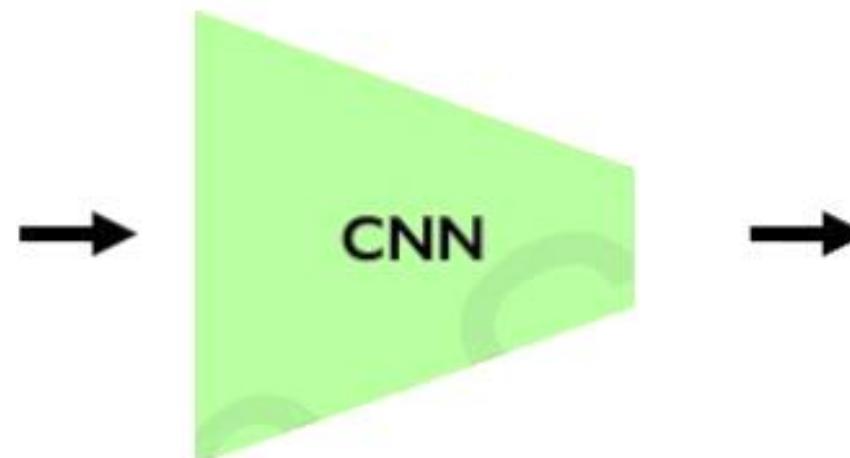


Image x



Taxi

Class label y

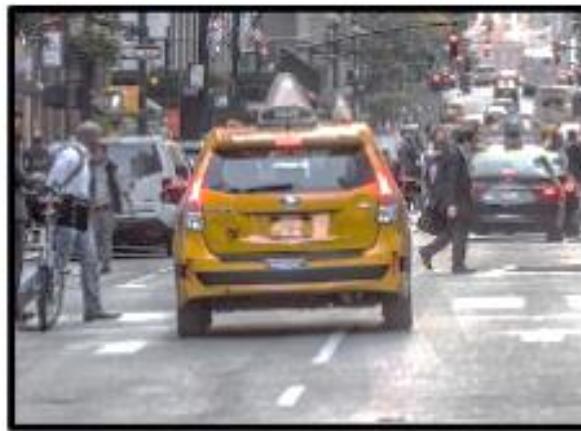
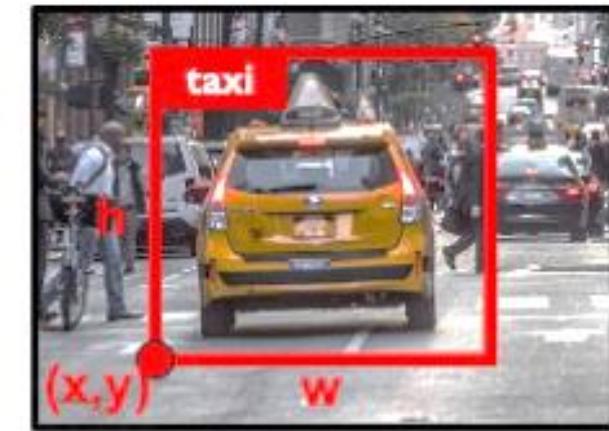
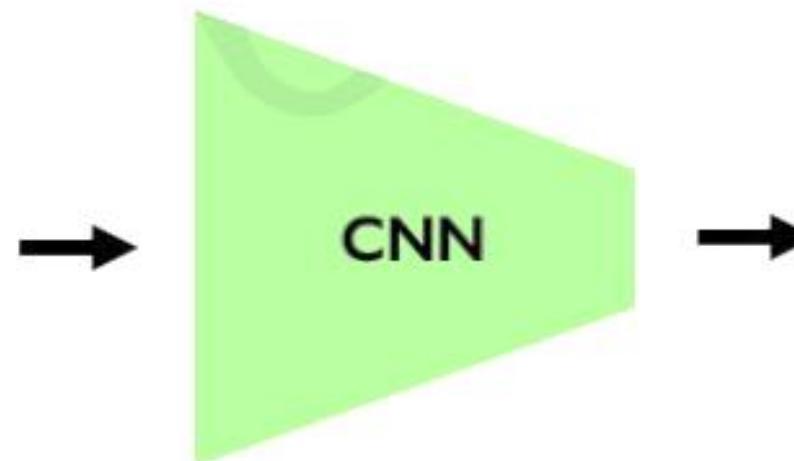


Image x



Label (x, y, w, h)



# Detección de Objetos



Image X

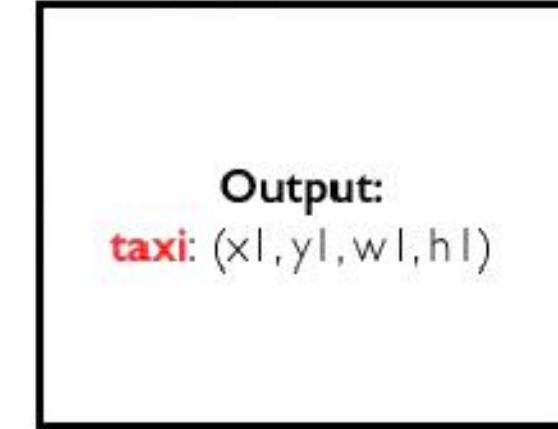
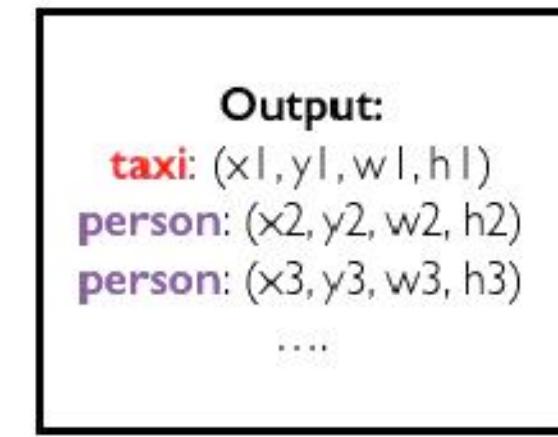


Image X



# Solución Ingenua para la Detección de Objetos

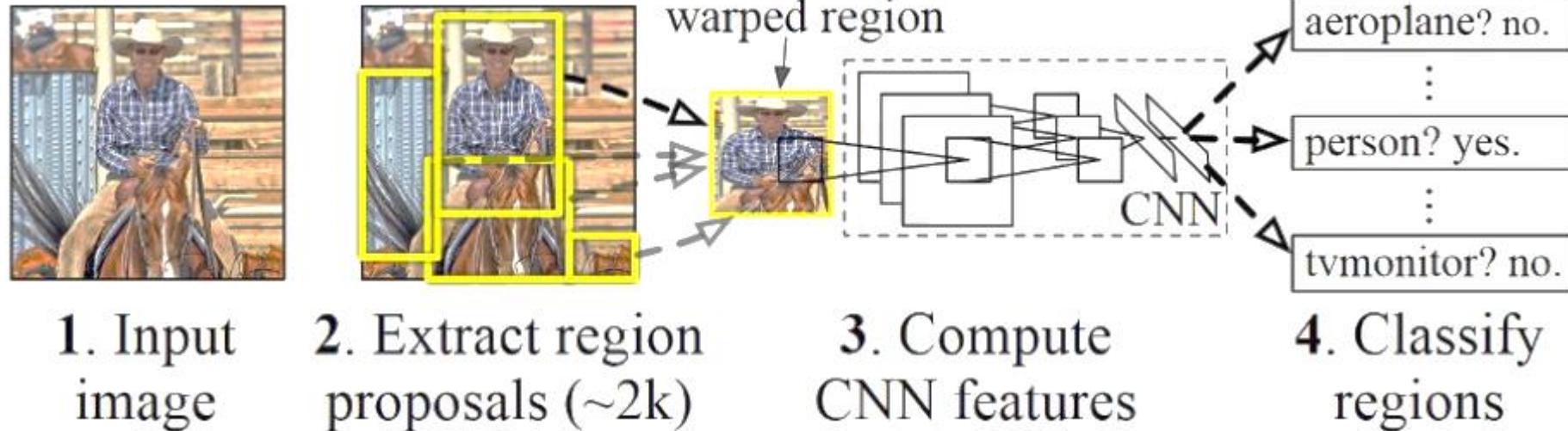


Problema: ¡Demasiadas entradas! Esto da como resultado demasiadas escalas, posiciones y tamaños



# Detección de Objetos con R-CNNs

Algoritmo R-CNN: encontrar regiones que creemos que contienen objetos. Usar una CNN para clasificarlas.



Problemas:

1. ¡Lento! Muchas regiones; inferencia que consume mucho tiempo.
2. Frágil: propuestas de regiones definidas manualmente



# Faster R-CNN Aprende Propuestas de Regiones

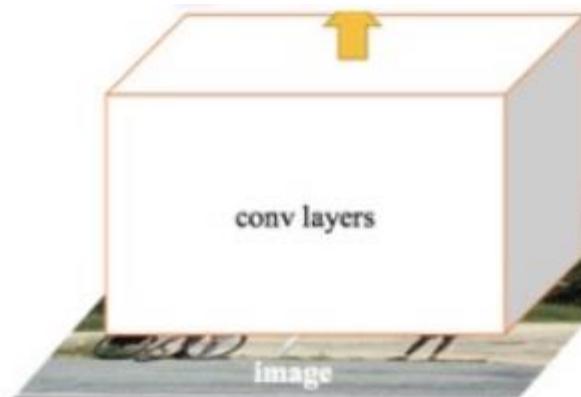
Extracción de características sobre las regiones propuestas

Clasificación de regiones → detección de objetos

La imagen de entrada va directamente al extractor de características convolucional.

¡Rápido! Solo se ingresa la imagen una vez.

Red de propuestas de regiones para aprender regiones candidata  
Aprendido, basado en dato



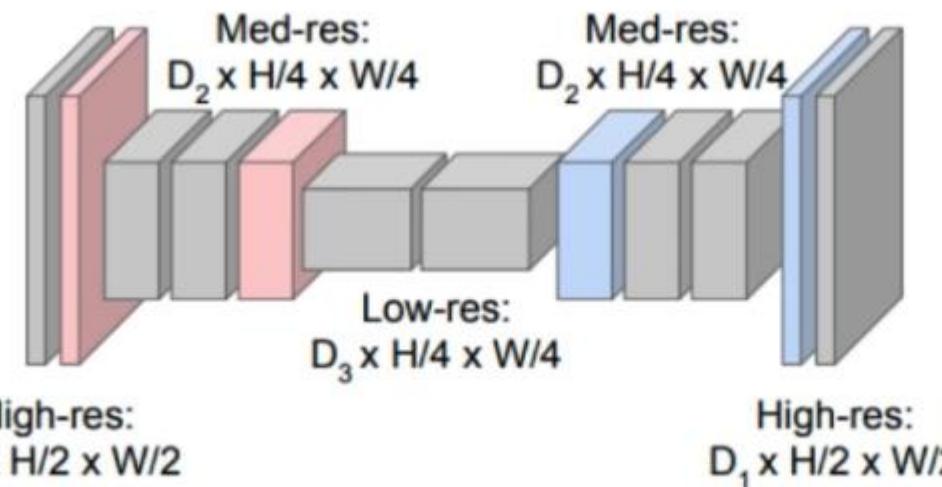
# Segmentación Semántica: Redes Totalmente Convolucionales

FCN: Fully Convolutional Network.

Red diseñada con todas las capas convolucionales, con operaciones de downsampling y upsampling.



Input:  
 $3 \times H \times W$



Predictions:  
 $H \times W$

 `tf.keras.layers.Conv2DTranspose`

 `torch.nn.ConvTranspose2d`



# Control Continuo: Navegación a partir de la Visión

Raw Perception  
 $I$   
(ex. camera)



Coarse Maps  
 $M$   
(ex. GPS)

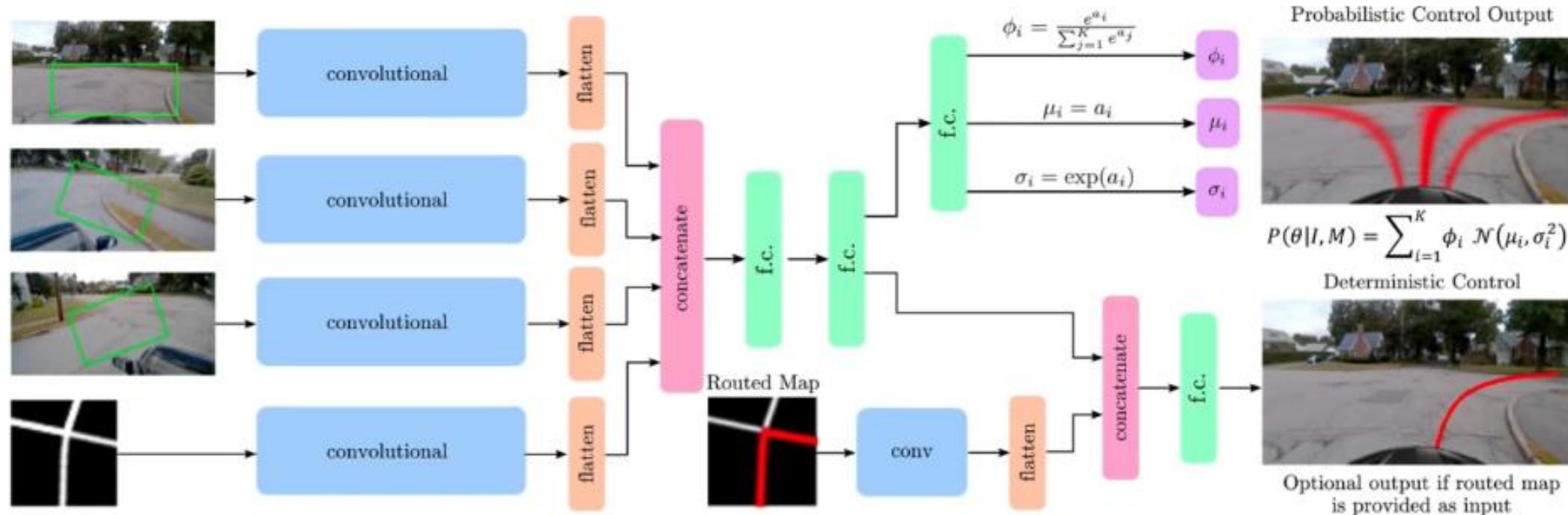


Possible Control Commands



# Marco de Extremo a Extremo para Navegación Autónoma

Todo el modelo se entrena de extremo a extremo sin ninguna etiquetación o anotación humana.



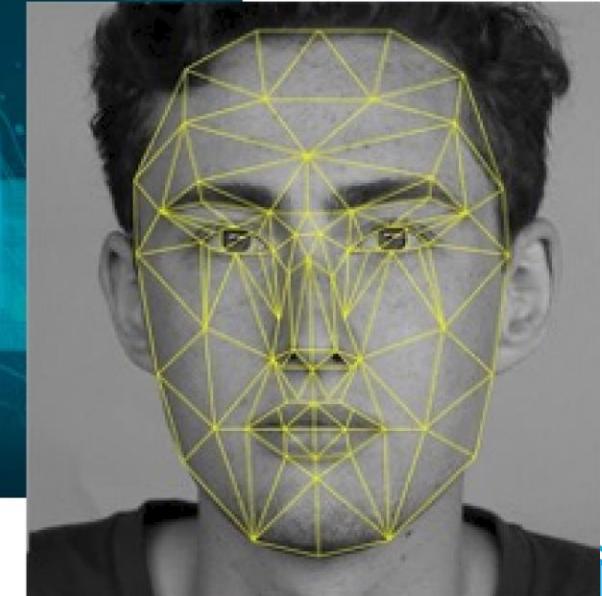
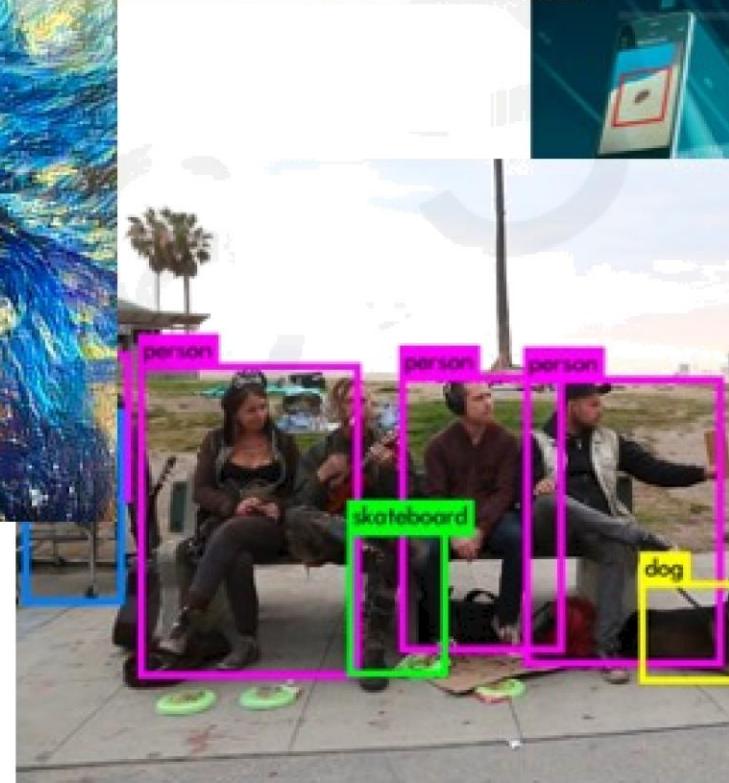
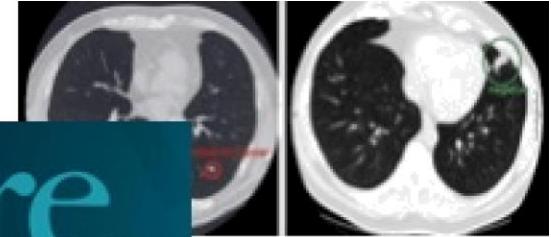


# Auto ON

Navigation and Localization



# Aprendizaje Profundo para Visión por Computador: Impacto



# Aprendizaje Profundo para Visión por Computador: Resumen

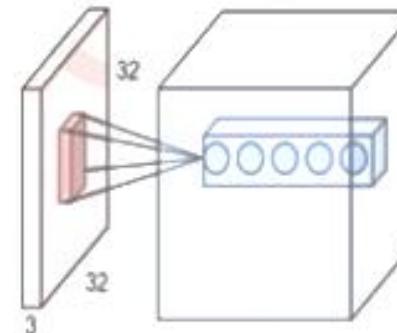
## Fundamentos

- ¿Por qué visión por computador?
- Representación de imágenes
- Convoluciones para extracción de características



## CNNs

- Arquitectura de CNN
- Aplicación a clasificación
- ImageNet



## Aplicaciones

- Segmentación, descripción de imágenes, control
- Seguridad, medicina, robótica

