

# Image Classification

Aram Karimi

LT2318 H21 Artificial Intelligence: Cognitive Systems

November 15th, 2020

- ▶ What is "image classification"?
- ▶ Why is "image classification" important?
- ▶ How does "image classification" work?
- ▶ "Image classification" methods Using machine
- ▶ The need for AI to understand image data
- ▶ Convolutional Neural Network (CNN)
- ▶ When do we use pre-trained image features?
- ▶ Hands-on Tutorial

# What is "image classification"?

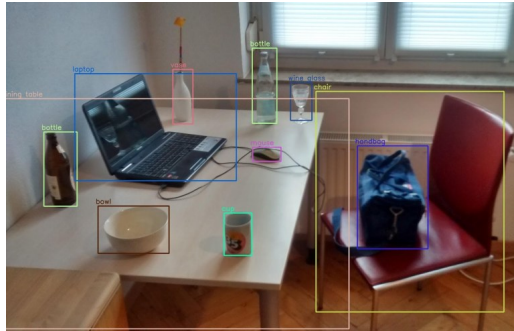


Image classification is a basic task for human, but still one of the most important tasks that computer vision engineers can tackle.

# What makes image classification a very important task?



```

0 2 15 0 0 11 10 0 0 0 0 9 9 0 0 0
0 0 0 4 60 157 236 255 255 177 95 61 32 0 0 29
0 10 16 119 238 255 244 245 243 250 249 255 222 103 10 0
0 14 170 255 255 244 254 255 253 245 255 249 253 261 124 1
2 98 255 228 255 251 254 211 141 116 122 215 251 238 255 45
13 217 243 255 155 33 226 52 2 0 10 13 232 255 255 36
16 229 252 254 49 12 0 0 7 7 0 70 237 252 235 62
6 141 245 255 217 25 11 9 3 0 115 236 243 255 137 0
0 87 252 250 248 215 60 0 1 121 252 255 248 144 6 0
0 13 113 255 255 245 255 182 181 248 252 242 208 36 0 15
1 0 5 117 251 255 241 255 247 255 241 162 17 0 7 0
0 0 0 4 58 251 255 246 254 253 255 120 11 0 1 0
0 0 4 97 255 255 255 248 252 255 244 255 182 10 0 4
0 22 206 252 246 251 241 100 24 113 255 245 255 194 9 0
0 111 255 242 255 158 24 0 0 6 39 255 232 230 56 0
0 218 251 250 137 7 11 0 0 0 2 62 255 250 125 3
0 173 255 255 101 9 20 0 13 3 13 182 251 245 61 0
0 107 251 241 255 230 98 55 19 118 217 248 253 255 52 4
0 18 146 250 255 247 255 255 255 249 255 240 255 129 0 5
0 0 23 113 215 255 250 248 255 255 248 248 118 14 12 0
0 0 6 1 0 52 153 233 255 252 147 37 0 0 4 1
0 0 5 5 0 0 0 0 0 14 1 0 6 6 0 0
  
```

An image is a large grid of numbers between [0, 255]

An image can be of size 800 x 600 pixels, each pixel is represented via three numbers, which provide values of RGB (red, green, blue) channels

## What the computer sees



10	45	78	99	01	52	77	32
44	69	00	18	35	41	00	02
35	60	02	00	75	35	88	09
11	33	50	70	00	08	01	37
43	57	60	06	12	41	69	05
78	91	09	00	66	31	24	60
80	88	71	07	11	10	00	38
76	98	00	00	10	31	33	48
76	54	55	56	00	02	20	73

Image  
classification

82% cat

15% dog

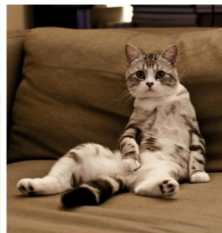
2% hat

1% mug

- ▶ Viewpoint
- ▶ Illumination



► Deformation



## ► Occlusion





## Challenges

### ► Clutter



## Challenges

- ▶ Intraclass variation



## How does Image classification work in machines?

- ▶ In digital image processing, image classification is done by automatically grouping pixels into specified categories, so-called “classes.”
- ▶ The algorithms separate the image into different classes based on their prominent features or specific patterns

Image classification techniques are mainly divided into two categories:

- ▶ Supervised
- ▶ Unsupervised

Where is it used?

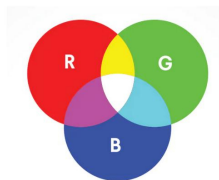
- ▶ Robotics, Computer Vision, NLP, information retrieval, etc.

## Representing images with “features”

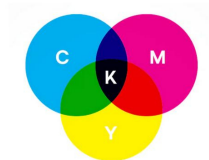
- ▶ CV vs NLP features
- ▶ Visual features
  - ▶ Color, size, center, orientation, etc.
  - ▶ Invariant to transformations
- ▶ Lexical features / semantic classes
  - ▶ Labels, context
- ▶ Learned vs. pre-engineered features
  - ▶ We need to choose how to represent an image

## Visual feature: color

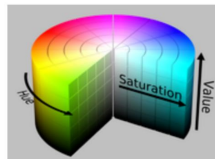
RGB



CMYK



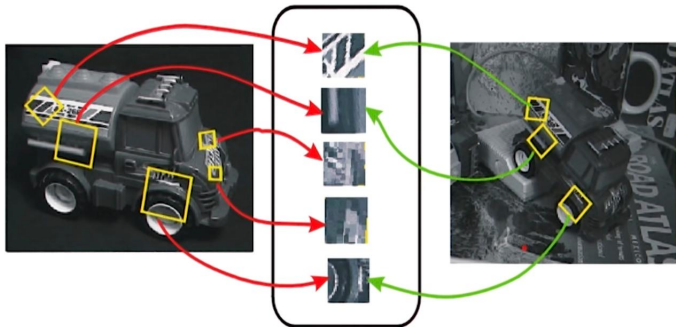
HUE



## Visual features: SIFT, HOG, SURF

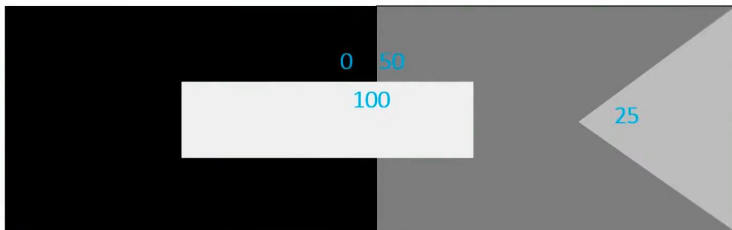
- ▶ **SIFT**: Scale-Invariant Image Transform
  - ▶ Commonly used in CV
  - ▶ Extract invariant (not changeable) image features
  - ▶ Applied to grayscale images
  - ▶ Mathematically complicated, computationally heavy
  - ▶ Based on histogram of gradients, e.g. computing the gradients of each pixel in the image takes a lot of time
  - ▶ Quite slow compared to SURF
  - ▶ Does not work well with lighting changes and blur

## Visual features: SIFT



## Visual features: HOG(Histogram of Oriented Gradients)

- ▶ **HOG::** compute centered horizontal and vertical gradients
  - ▶ Tries to extract contrasts in various image parts
  - ▶ Computes gradients magnitudes and their directions



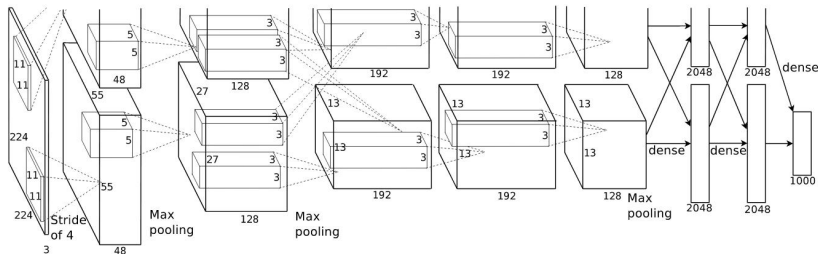
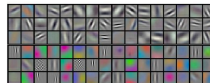


## The need for AI to understand image data

- ▶ Supervise learning task
  - ▶ **Convolutional Neural Networks** as feature extractors  
hierarchical layer-wise representation learning
  - ▶ Extract invariant (not changeable) image features
  - ▶ CNN is a hierarchical deep learning model which is able to  
model data at more and more abstract representations
  - ▶ CNN features are highly adaptive, they are trained end-to-end
  - ▶ CNN can learn features similar to SIFT and HOG from training  
examples alone, which is quite cool. Therefore, **using CNNs  
minimizes feature engineering**

## Convolutional Neural Network (CNN)

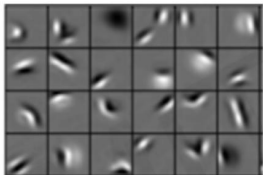
- ▶ Used for object detection, image classification, image captioning, etc.
- ▶ Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet classification with deep convolutional neural networks. In Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1 (NIPS'12). Curran Associates Inc., Red Hook, NY, USA, 1097–1105.



## Learning Feature Representation

Can we learn a hierarchy of features directly from the data instead of hard engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

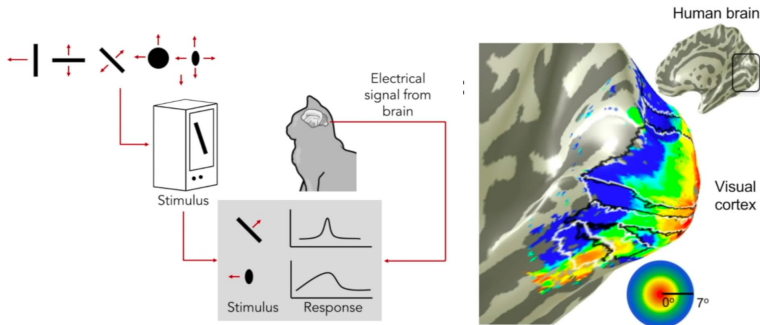
High level features



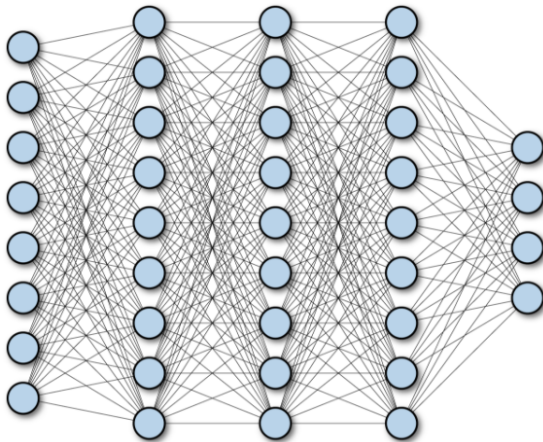
Facial structure

## CNNs: inspiration

- Hubel and Wiesel (1959, 1962, 1968): cat's visual cortex maps information in a structured and hierarchical way



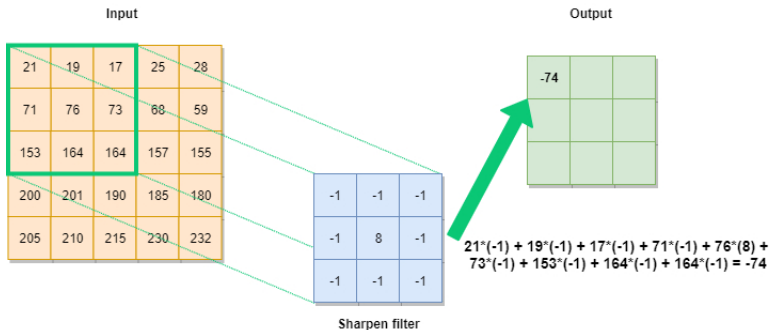
# Fully Connected Layer



## Using Spatial Structure

**Input:** 2D image.

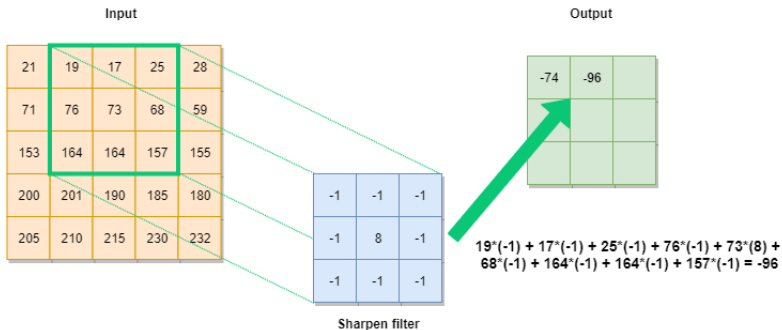
**Idea:** Connect patches of input to neurons in hidden layer.  
 (Neuron connected to region of input only sees these values)



## Using Spatial Structure

**Input:** 2D image.

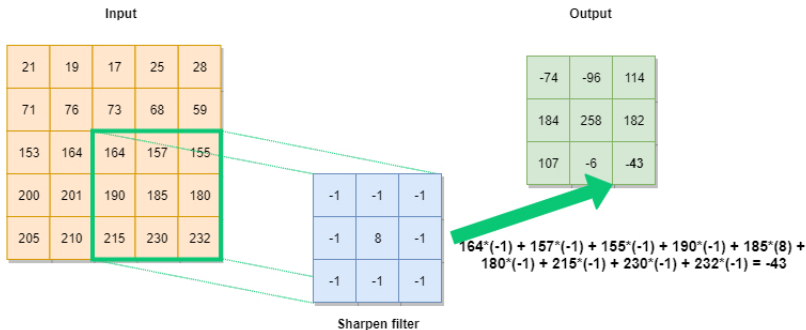
**Idea:** Connect patches of input to neurons in hidden layer.  
 (Neuron connected to region of input only sees these values)



## Using Spatial Structure

**Input:** 2D image.

**Idea:** Connect patches of input to neurons in hidden layer.  
 (Neuron connected to region of input only sees these values)

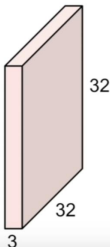




# CNN: we want to preserve spatial structure

## Convolution Layer

32x32x3 image



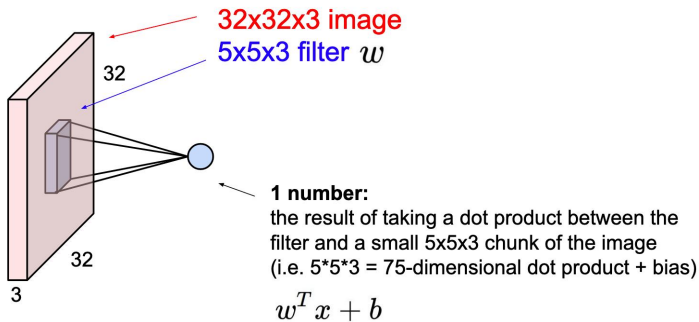
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

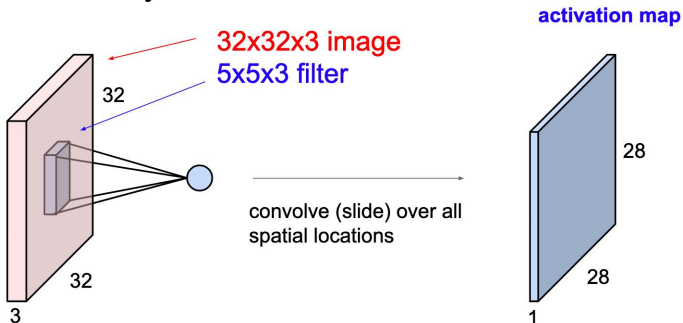
# CNN: apply filter to the convolution layer

## Convolution Layer



# CNN: use filter to get an activation map

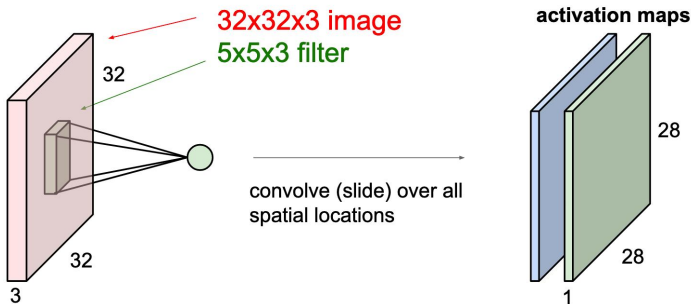
## Convolution Layer



# CNN: maps per filter

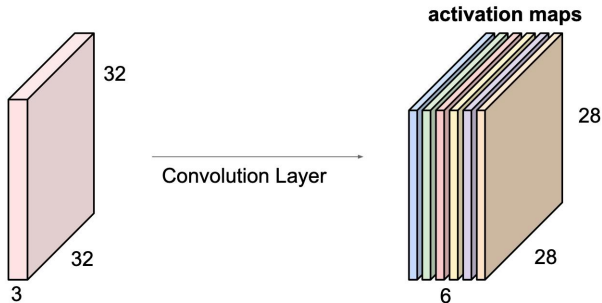
## Convolution Layer

consider a second, **green** filter



# CNN: activation maps

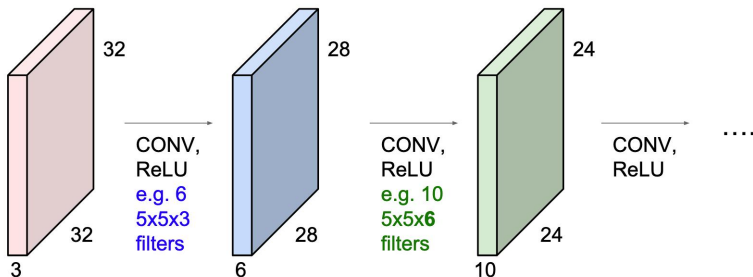
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



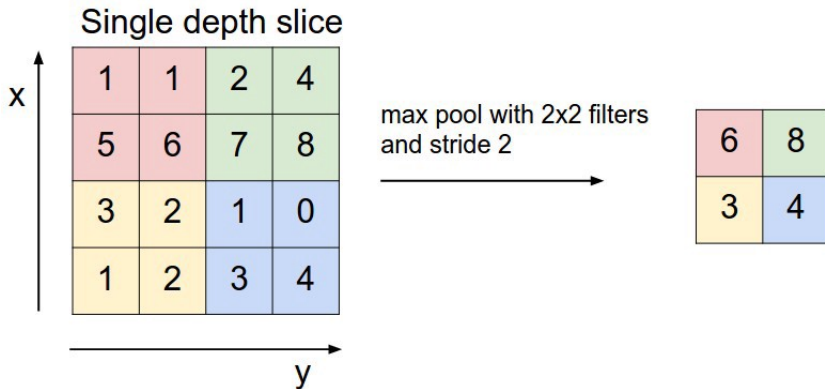
We stack these up to get a "new image" of size 28x28x6!

# CNN: a stack of CONV, FC, POOL + activations

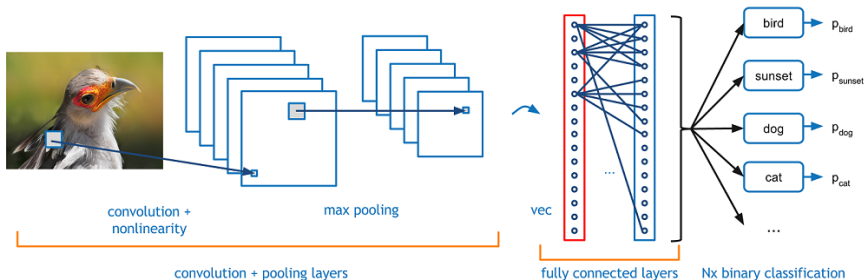
**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



# Pooling layer



Finally, the raw values which are predicted output by network are converted to probabilistic values with use of soft max function.

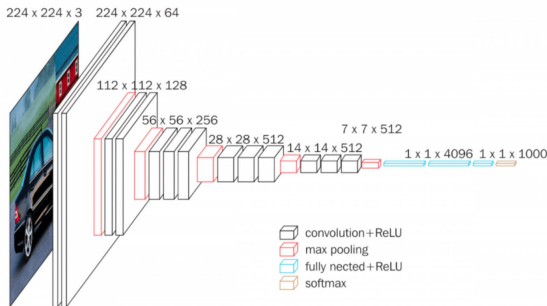




## CNN: conclusion

- ▶ Smaller filters, deeper architectures
- ▶ Tend to remove POOL and FC, keep CONV only

### Example CNN network structure, VGG16



## Useful Links

- ▶ PyTorch Image Classification Tutorial:  
[https://pytorch.org/tutorials/beginner/blitz/cifar10\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html)
- ▶ TensorFlow Image Classification Tutorial:  
<https://www.tensorflow.org/tutorials/images/classification>
- ▶ More recent work on CNNs:  
[https://github.com/matterport/Mask\\_RCNN](https://github.com/matterport/Mask_RCNN)  
<https://github.com/facebookresearch/detectron2>
- ▶ Accuracy scores for published CNNs: [Accuracy scores for published CNNs:](#)