
Artificial Neural Networks

Esra Suel

CASA0006: Data Science for Spatial Systems

Slides adapted from Ender Konukoglu at ETH Zurich & Huanfa Chen at UCL

Outline

- Deep learning for images: Convolutional Neural Networks (CNNs)
 - CNN basics
 - CNNs for geospatial research
 - Potential Problems with CNNs
- Deep learning for text data

Deep learning for images

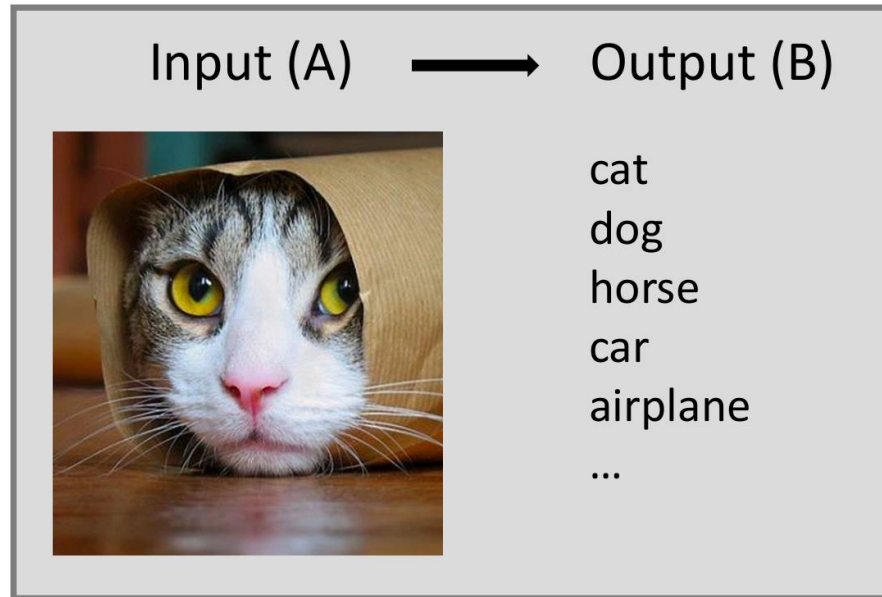
Convolutional Neural Networks (CNNs)

Deep learning for images (computer vision)

Applications

1. Image Classification
 2. Image Classification with Localization
 3. Object Detection
 4. Semantic Segmentation
 5. Instance Segmentation
-

Image Classification



airplane

automobile

bird

cat

deer

dog

frog

horse

ship

truck

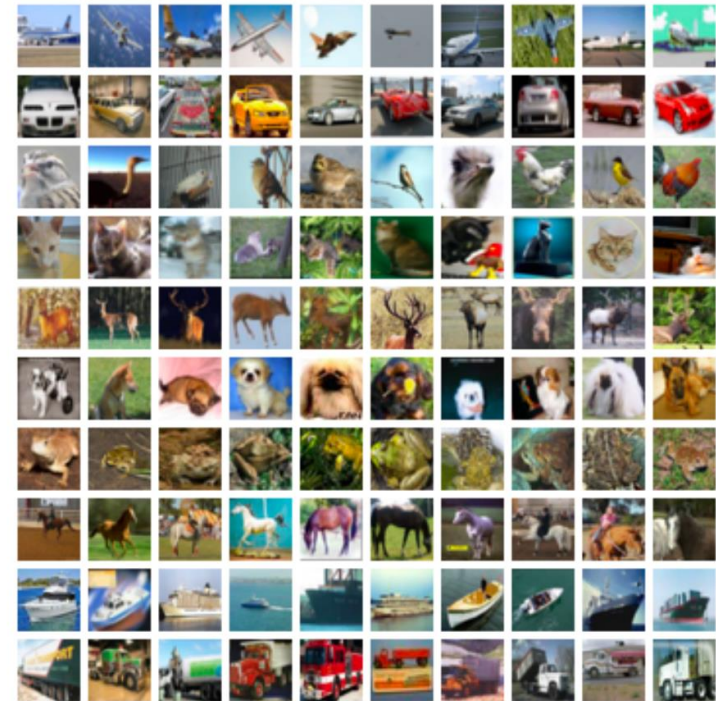
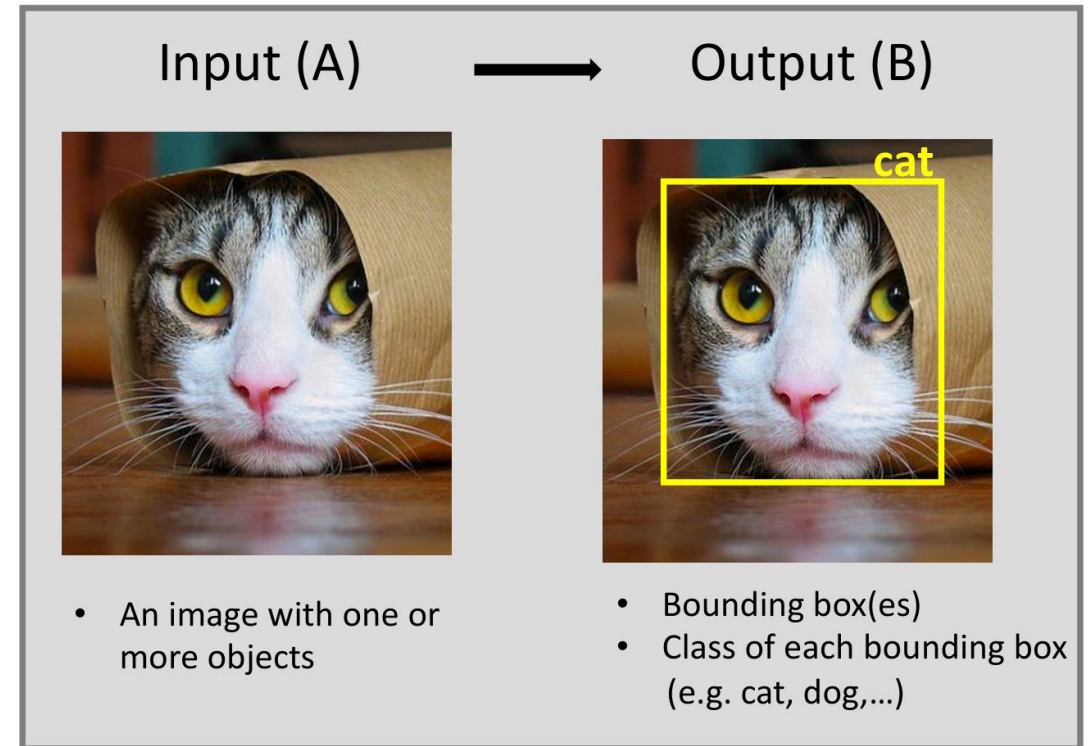


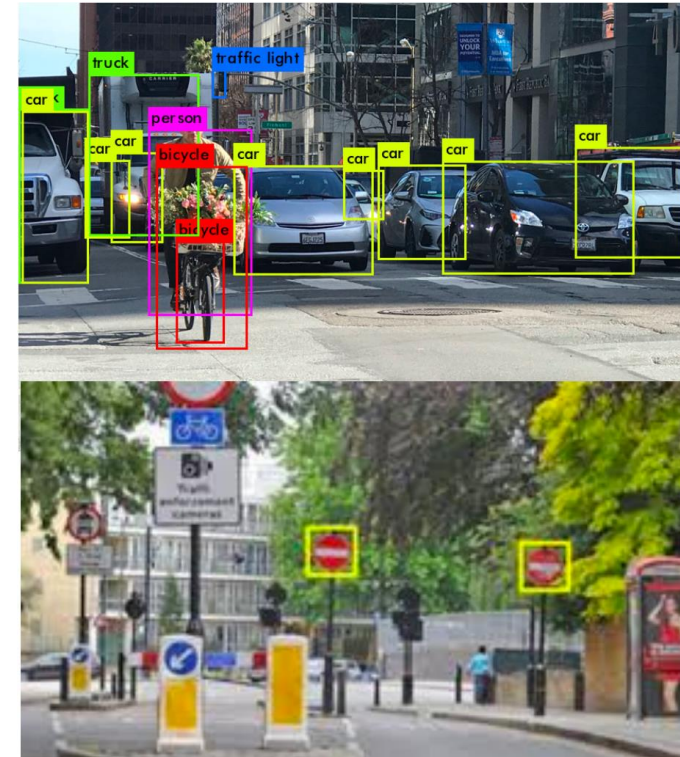
Image Classification with Localization

- Given an image we want to learn the class of the image and the class location in the image. We want a class and a rectangle of where that object is.
- Usually only one object is presented.



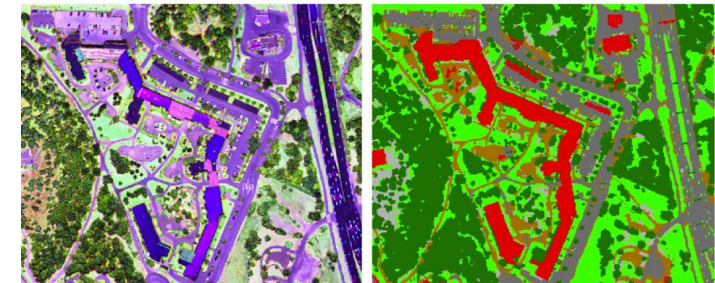
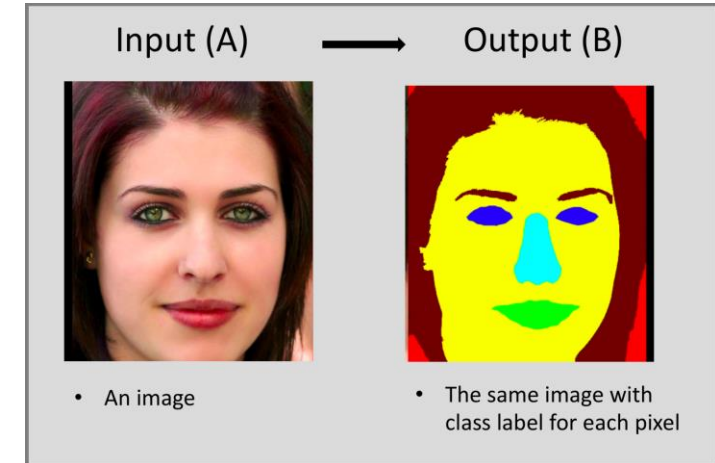
Object Detection

- We want to detect all objects in the image that belong to a specific class and give their boundary box.
- An image can contain more than one object with different classes.



Semantic Segmentation

- We want to label each pixel in the image with a category label.
- Semantic Segmentation don't differentiate instances, only care about pixel labels.
- It detects no objects just pixels – it does not tell you #cars in the image.
- If two objects of the same class are intersected, we won't be able to separate them.

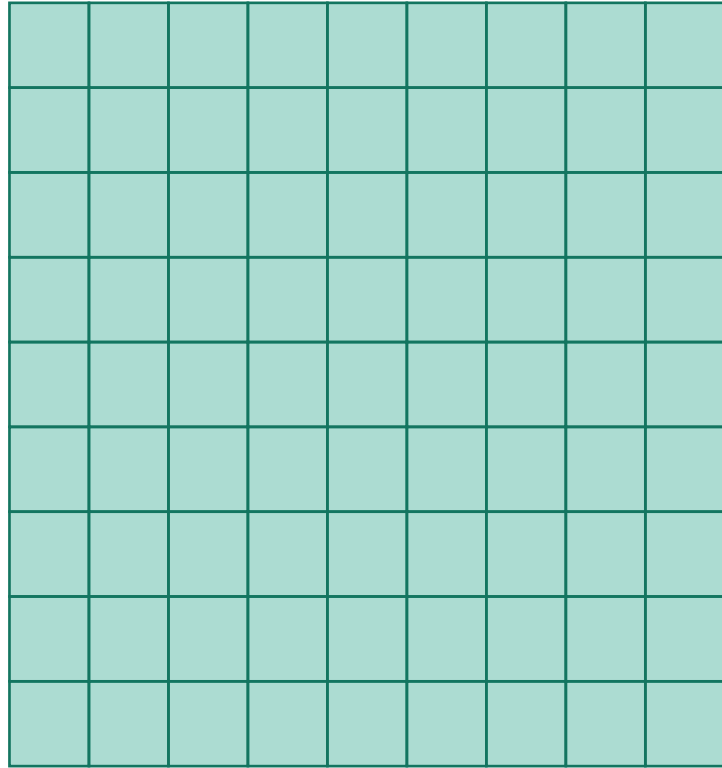


Instance Segmentation

- This looks like the full problem and the most challenging.
- We want to identify and distinguish all objects.



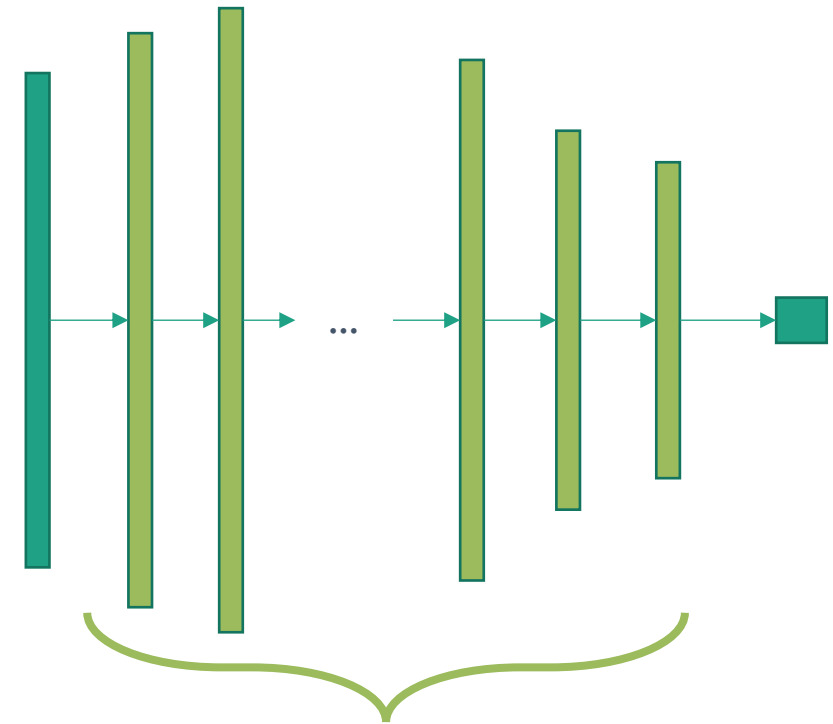
Naïve approach



image



vectorized
image



Hidden layers

Convolutional Neural Networks (CNNs) for images

Images can be large. CNN is fast and accurate to work with images.

E.g. Comparing FCNN and CNN with similar number of output features (~ 1000)

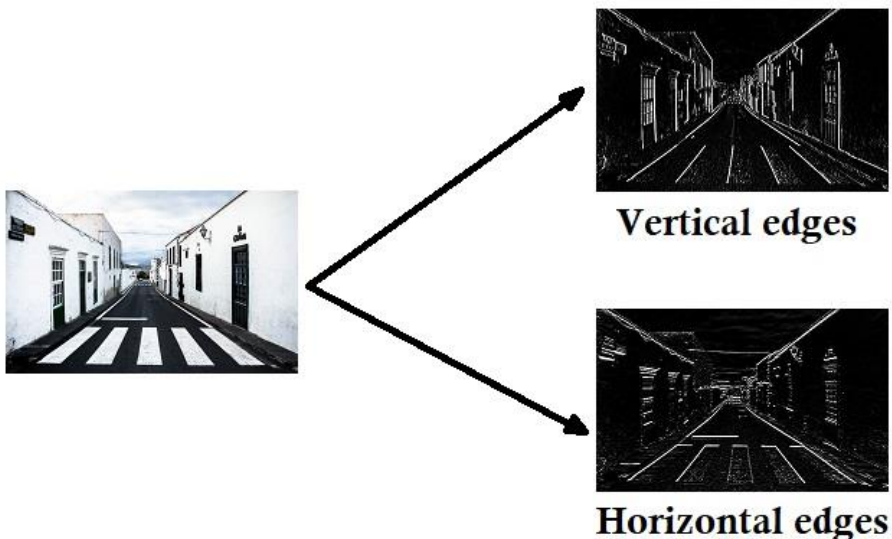
Given a 1000×1000 image (having 10^6 pixels/features), if a FCNN is used (1 hidden layer with 1000 neurons), then there are 10^9 parameters to learn. Very challenging.

In contrast, if we use a CNN with 967×967 filter, the output has 33×33 features (close to 1000) and only 10^6 parameters to learn.

CNNs for images

Parameter sharing

A filter is actually a feature detector. A filter that is useful in one part of the image is probably useful in another part of the image.

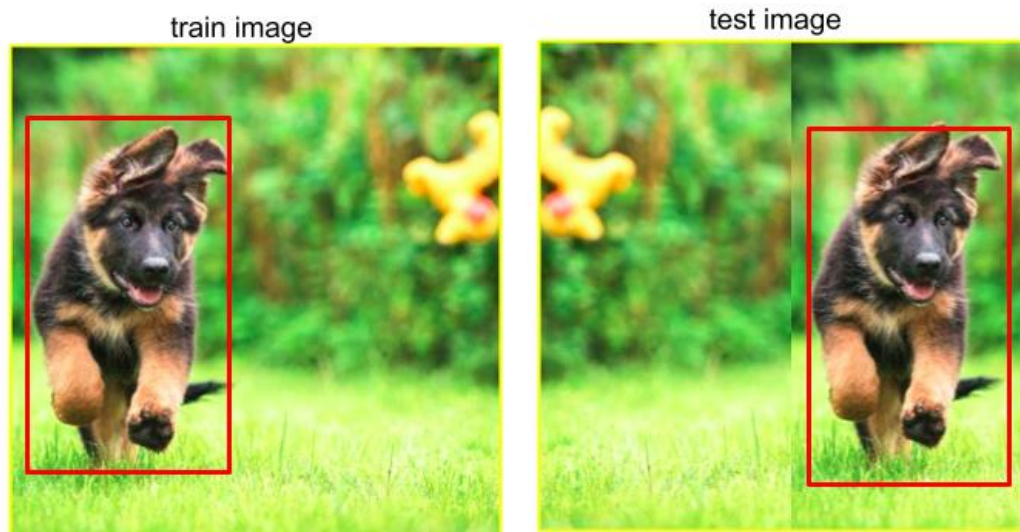


E.g. a vertical edge filter is useful in many parts of this image

CNNs for images

Translation invariance

In each layer, each output value depends only on a small number of inputs, which makes it translation invariance. In other words, if a pattern/object is translated, it is still identifiable by CNN.



CNN would detect the 'dog' object regardless of its location. This is called translation invariance.

The non-linearity will remain the same

Fully connected architecture

$$a_{l,k} = \sum_j w_{l,kj} h_{l-1,j} + b_{l,k}$$

$$h_{l,k} = \sigma \left(\sum_j w_{l,kj} h_{l-1,j} + b_{l,k} \right)$$

Convolutional architecture

$$a_{l,k} = \sum_j w_{l,kj} * h_{l-1,j} + b_{l,k}$$

$$h_{l,k} = \sigma \left(\sum_j w_{l,kj} * h_{l-1,j} + b_{l,k} \right)$$

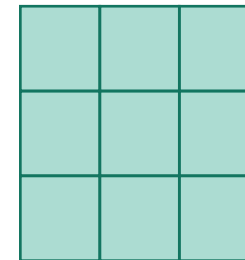
Convolution

Image: x

544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

$$w * x$$

Convolution kernel: w



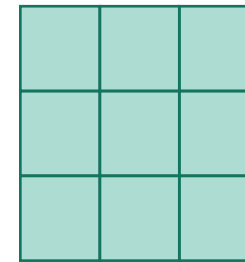
Convolution

Image: x

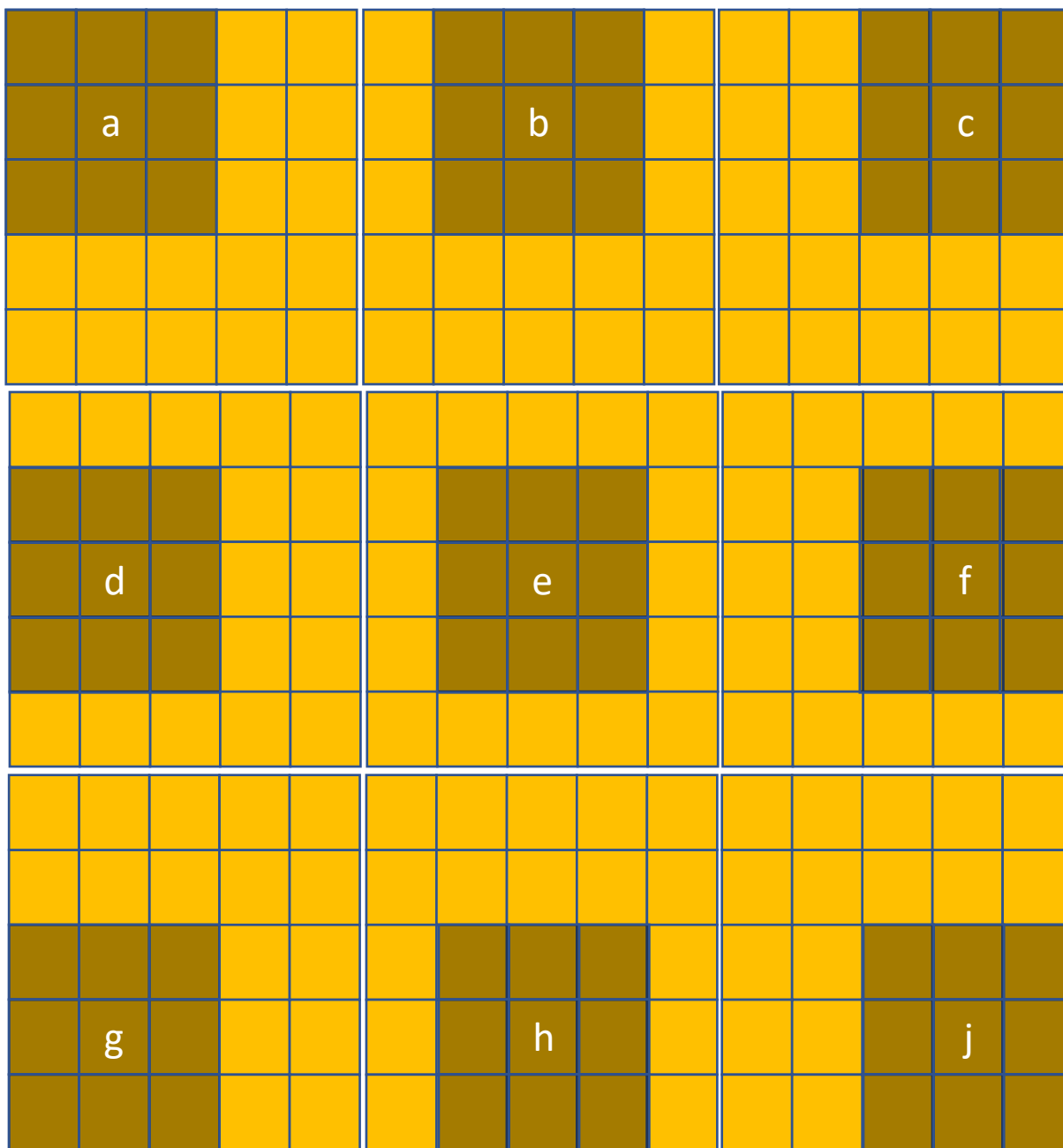
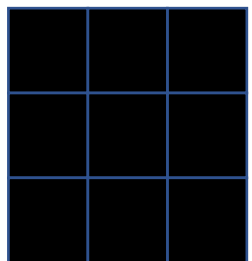
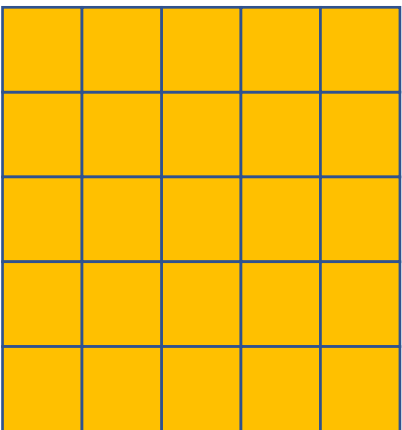
544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

$w * x$

Convolution kernel: w



$$a_{ij} = \sum_p \sum_q x_{(i-p)(j-q)} w_{(p)(q)}$$



a	b	c
d	e	f
g	h	j

Larger hidden layers

Image: $x \in \mathbb{R}^{M_0 \times N_0}$

544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

d_1



Fully connected: $h_1 \in \mathbb{R}^{d_1}$

M_1



Convolutional: $h_1 \in \mathbb{R}^{d_1 \times (N_1 \times M_1)}$

Channel size

Channel size is linked with kernel size and the type of convolution

544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

$$a_{ij} = \sum_p \sum_q x_{(i-p)(j-q)} w_{(p)(q)}$$

- When the kernel is placed in the image – no problem

Channel size

Channel size is linked with kernel size and the type of convolution

	544	552	570	585	600	607	608	581	558	577
	549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582	
529	514	486	476	483	509	552	584	604	586	
506	499	468	421	459	547	588	596	598	603	
567	561	519	484	510	557	586	612	603	565	
579	594	581	563	557	553	572	587	575	575	
590	601	594	586	580	563	559	587	602	585	
596	602	602	595	586	585	592	577	545	557	
593	614	589	568	588	625	610	546	519	557	

$$a_{ij} = \sum_p \sum_q x_{(i-p)(j-q)} w_{(p)(q)}$$

- When the kernel is placed in the image – no problem
- When it is placed on the boundary – it is not well defined
- Out-of-boundary values are not defined
- Two options:
 1. Valid convolution: only evaluate convolution when all the elements are defined
 2. Padding (Same): pad the boundaries so that result of the convolution will have the same size

Valid convolution

	544	552	570	585	600	607	608	581	558	577
	549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582	
529	514	486	476	483	509	552	584	604	586	
506	499	468	421	459	547	588	596	598	603	
567	561	519	484	510	557	586	612	603	565	
579	594	581	563	557	553	572	587	575	575	
590	601	594	586	580	563	559	587	602	585	
596	602	602	595	586	585	592	577	545	557	
593	614	589	568	588	625	610	546	519	557	

- If the kernel is centered, i.e. $w_{(0)(0)}$ is the center of the kernel, then convolution can only be evaluated within the red area
- You loose a pixel at each end of the picture
- If the kernel center is the top left corner, then the green area is the valid area
- You loose two pixels at the bottom and right of the image

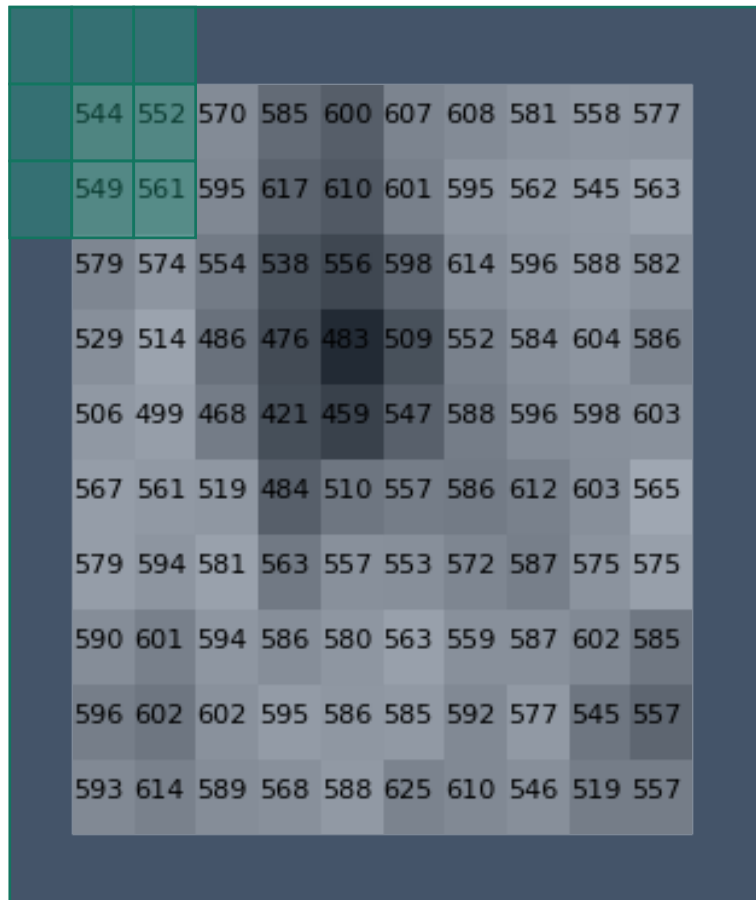
$$h_{l-1} \in \mathbb{R}^{d_{l-1} \times (M_{l-1} \times N_{l-1})}$$

$$w_{l,..} \in \mathbb{R}^{k_1 \times k_2}$$

$$h_l \in \mathbb{R}^{d_l \times (M_l \times N_l)}$$

$$M_l = M_{l-1} - k_1 + 1 \quad N_l = N_{l-1} - k_2 + 1$$

Same padding



Alternatively you can pad the image on the boundaries so that channels will have the same size across layers

$$h_{l-1} \in \mathbb{R}^{d_{l-1} \times (M_{l-1} \times N_{l-1})}$$

$$h_{l-1} \in \mathbb{R}^{d_{l-1} \times ((M_{l-1} + k_1 - 1) \times (N_{l-1} + k_2 - 2))}$$

$$h_l \in \mathbb{R}^{d_l \times (M_l \times N_l)}$$

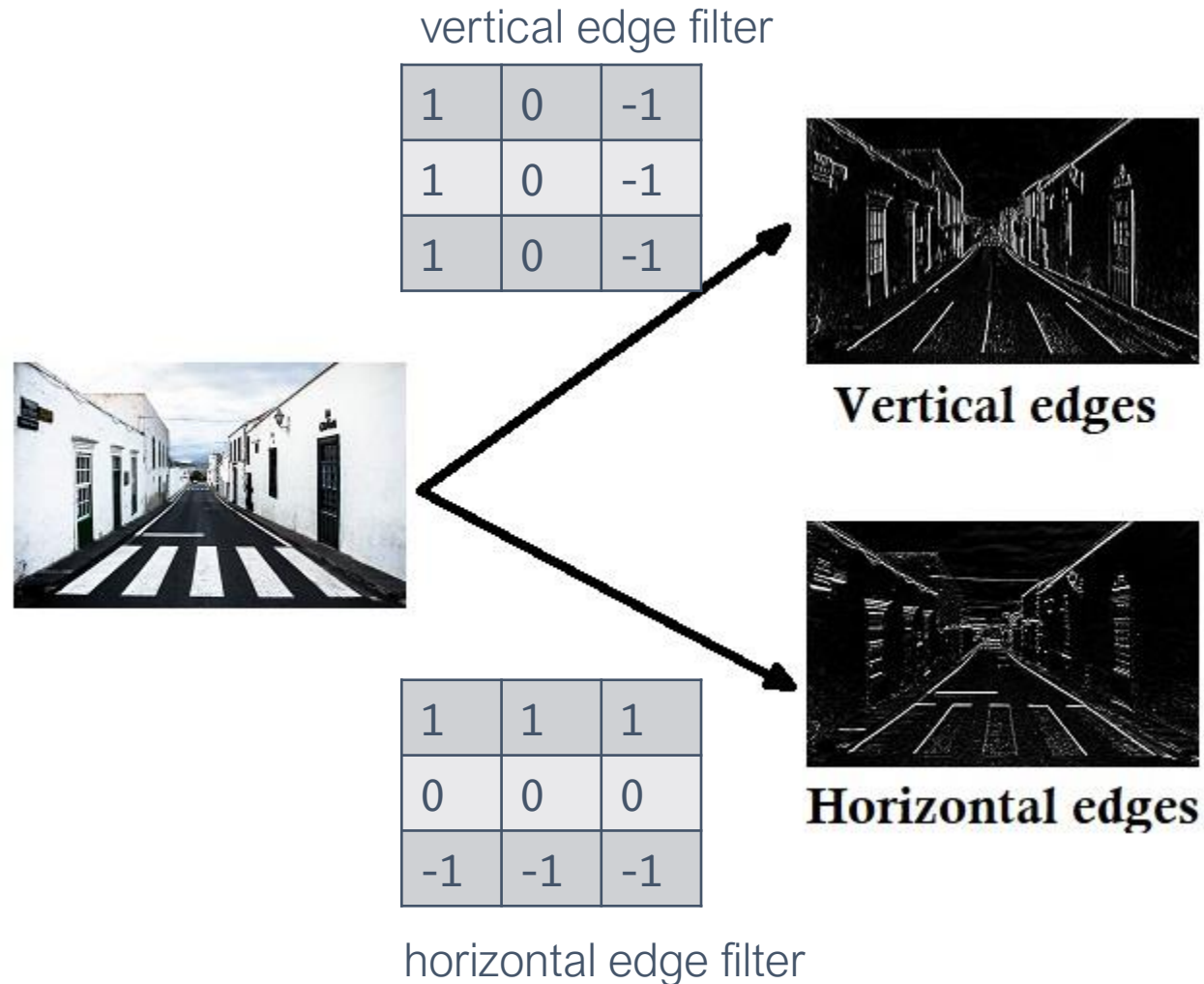
$$M_l = M_{l-1} \quad N_l = N_{l-1}$$

Where you pad depends on where the center of the kernel is.

Commonly you would use centered kernels – padding around the image as shown on the left

The value you pad is a parameter, 0 is used often but you can use symmetric padding for certain applications

Convolution operation and filter (kernel)



Pooling

544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

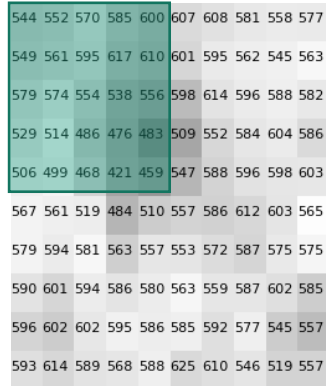
- Pool information in a neighborhood
- Represents the region with one number >> summarizes information
- Applied to each channel separately
- **Max-pooling** – maximum of the activation values
- **Min-pooling** – minimum of the activation values
- Both are non-linear operations, like median filtering
- **Averaging pooling** – linear operator
- Max-pooling is the most commonly used version

Max pooling

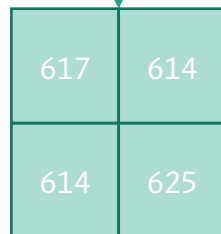
544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557

- Represents the entire region with the neuron that achieves the highest activation
- Leads to partial local translation invariance
- 617 in the highlighted area can be in any of the neurons, the pooled value will not change
- Does not lead to complete translation invariance
- Often applied with strides equal to the size of the kernel

Dimensionality reduction



544	552	570	585	600	607	608	581	558	577
549	561	595	617	610	601	595	562	545	563
579	574	554	538	556	598	614	596	588	582
529	514	486	476	483	509	552	584	604	586
506	499	468	421	459	547	588	596	598	603
567	561	519	484	510	557	586	612	603	565
579	594	581	563	557	553	572	587	575	575
590	601	594	586	580	563	559	587	602	585
596	602	602	595	586	585	592	577	545	557
593	614	589	568	588	625	610	546	519	557



617	614
614	625

- Leads to a substantial dimensionality reduction
- Even when the pooling kernel is of size 2x2, it can halve the image!
- As the size of the pooling kernel increase, the reduction increases as well
- Non-linear dimensionality reduction
- Only the most prominent activation is transmitted to the next layer
- More advanced pooling mechanisms exist
CapsuleNets [Sabour, Frosst and Hinton 2017]

Essential blocks lead to powerful algorithms

Convolutional layers and pooling are the essential blocks

They have been used to create complicated networks

First one

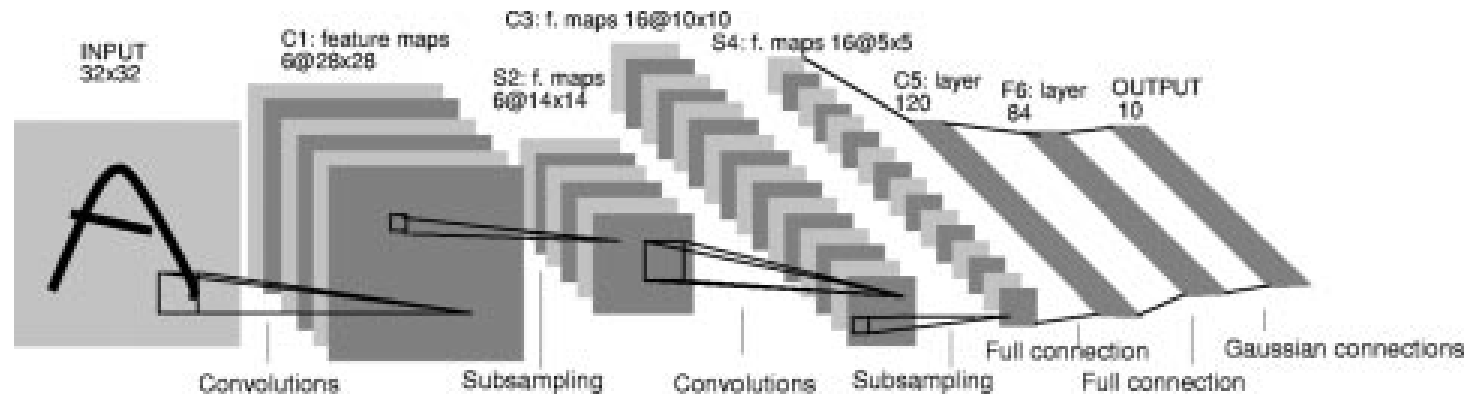


Fig. 2. Architecture of LeNet-5, a convolutional NN, here used for digits recognition. Each plane is a feature map, i.e., a set of units whose weights are constrained to be identical.

[Lecun, Bottou, Bengio and Haffner; Gradient-based learning applied to document recognition; 1998]

Then silence for a long time

Why silence

Models had too many parameters

They overfit for small datasets

We did not have very large datasets

Even for large sets, we did not have enough computation power to train the models until...



General purpose Graphical Processing Units (GPUs)
Allowed parallel processing

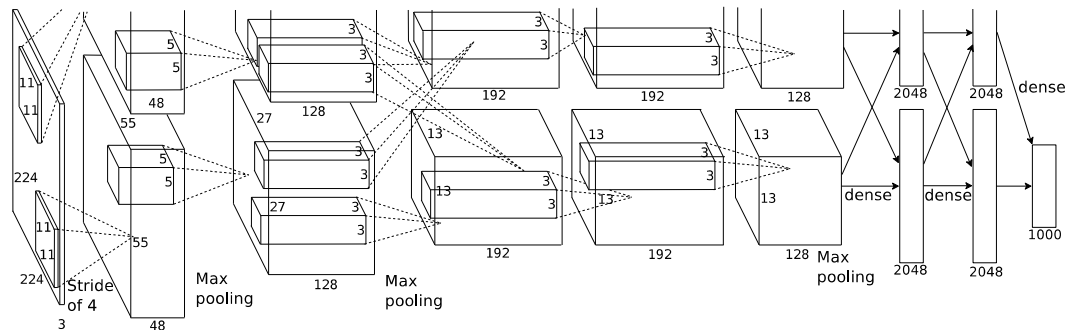
Then in 2012

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca



Network

Krizhevsky et al. almost halved the error rate in the ImageNet challenge

A simple CNN



First layer filters 11x11

A word about the ImageNet challenge

ImageNet: A Large-Scale Hierarchical Image Database

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li and Li Fei-Fei
Dept. of Computer Science, Princeton University, USA
{jiadeng, wdong, rsocher, jial, li, feifeili}@cs.princeton.edu

[CVPR 2009]

[International Journal of Computer Vision](#)

December 2015, Volume 115, [Issue 3](#), pp 211–252 | [Cite as](#)

ImageNet Large Scale Visual Recognition Challenge

Authors

[Authors and affiliations](#)

Olga Russakovsky , Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, Li Fei-Fei

Large scale object recognition challenge started in 2010

1.2 million training images

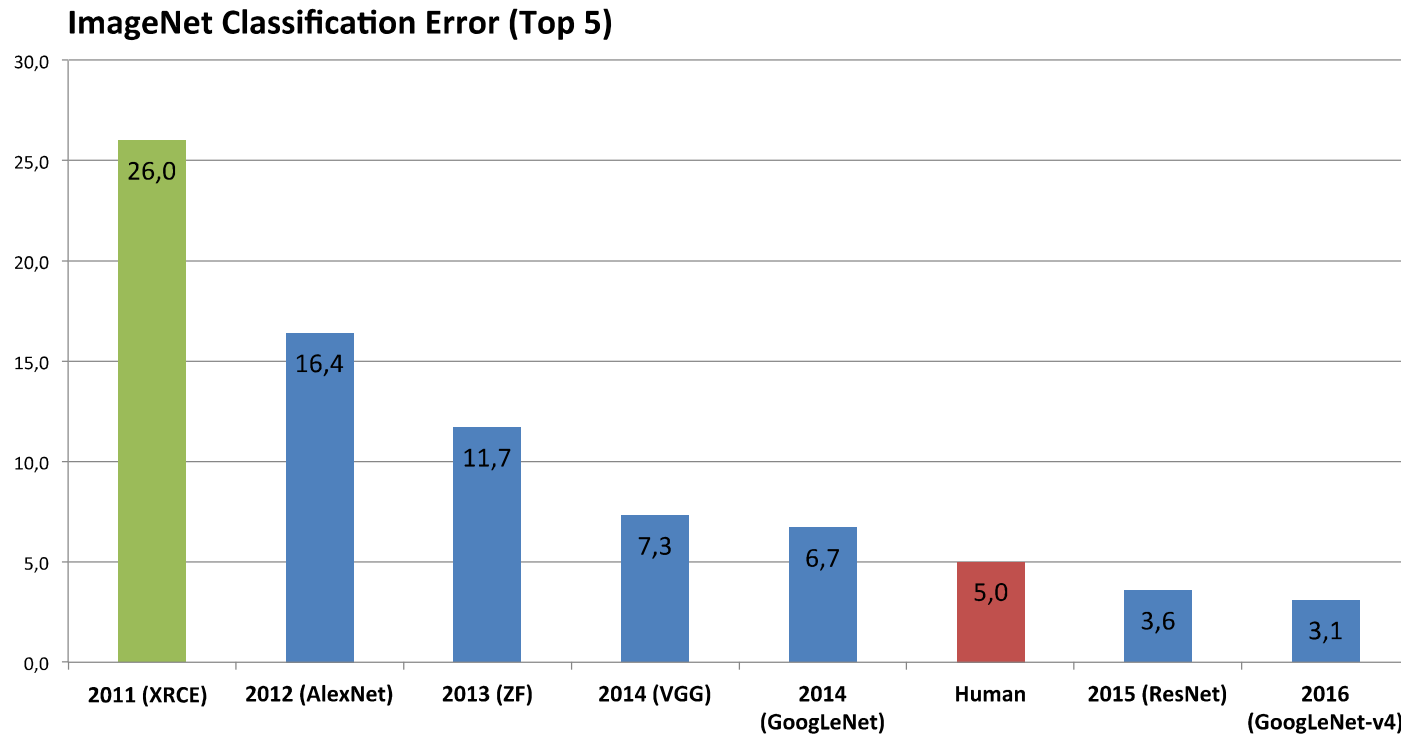
1000 object categories

200k validation and test images

2017 – 3 challenges:

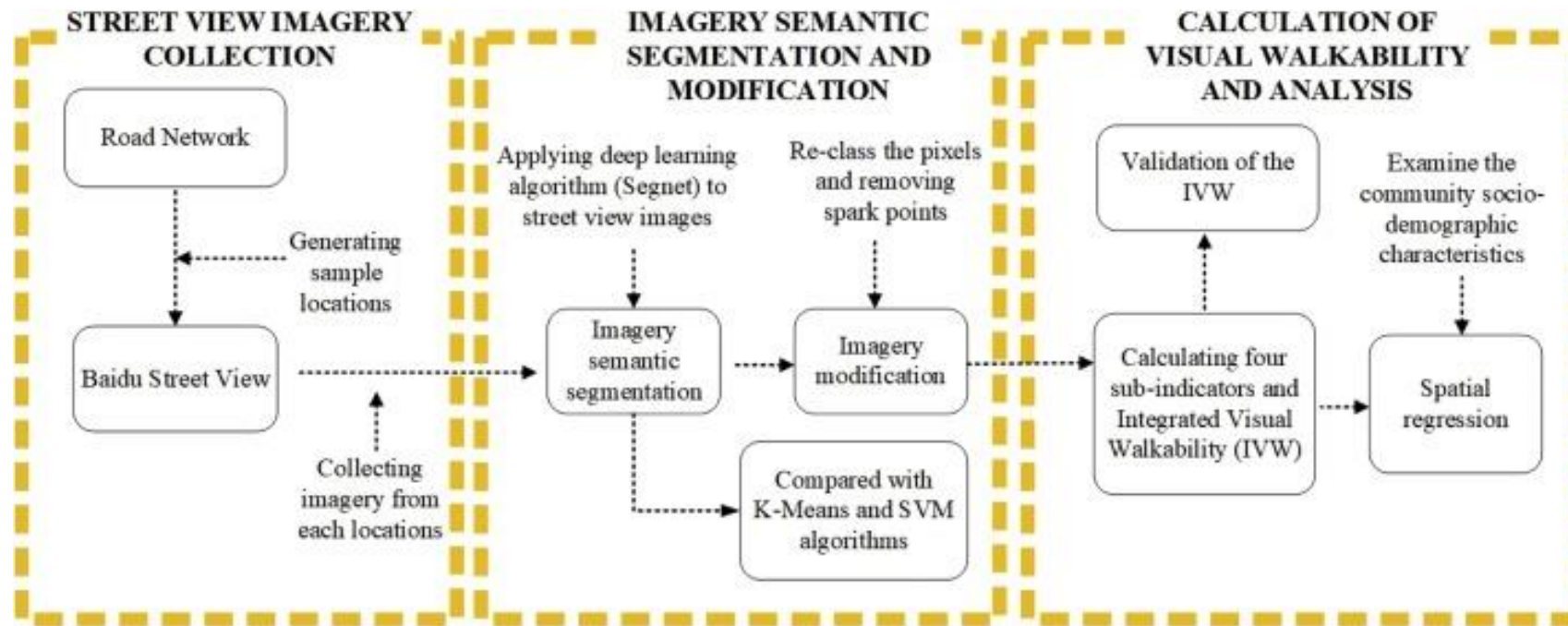
- Object localization
- Object detection
- Object detection from video

Historical evolution of the ImageNet Challenge



CNNs for Geospatial Research

Example 1: Visual walkability



[Zhou et al., 2019](#)

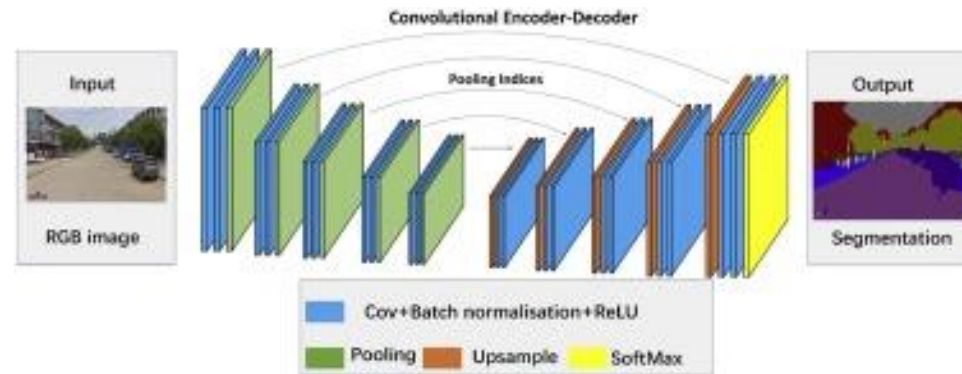
Visual walkability

- These indicators are derived from street images and used to build the IVW index.

Indicators	Definition	Formula	Explanation
Psychological Greenery	Extent to which the visibility of street vegetation can influence pedestrian psychological feelings	$G_i = \frac{(\sum_1^6 T_n)}{6 * Sum}$	T_n is the number of tree pixels; Sum is the total pixel number;
Visual Crowdedness	Extent to which the visibility of obstacles can influence pedestrian experiences	$C_i = \frac{(\sum_1^6 C_n)}{6 * Sum}$	C_n is the number of obstacles pixels; Sum is the total pixel number;
Outdoor Enclosure	How the room-like outdoor space is (the ratio of vertical objects to horizontal features)	$S_i = \frac{\sum_1^6 B_n + \sum_1^6 T_n}{\sum_1^6 P_n + \sum_1^6 R_n + \sum_1^6 F_n}$	B_n is the number of building pixels; T_n is the number of tree pixels; P_n refers to the number of pavement pixels; R_n refers to the number of road pixels; F_n refers to the number of fence pixels
Visual Pavement	Psychological impacts of the proportion of road and sidewalk on pedestrian experience	$D_i = \frac{\sum_1^6 P_n + \sum_1^6 R_n}{\sum_1^6 R_n}$	R_n refers to the number of road pixels; P_n refers to the number of pavement pixels; F_n refers to the number of fence pixels;
Integrated Visual Walkability (IVW)	Integrated Visual Walkability index	$IVW = (G\text{-level} + C\text{-level} + S\text{-level} + D\text{-level}) * 5$	

Zhou et al., 2019

Visual walkability



(a). Structure of SegNet



(a). Segmentation sample of SegNet



(b). Segmentation sample after reclassified



(c). Street view image and corresponding labelled training image

Example 2: measuring urban inequality

Measurements of inequalities: based on census costly and infrequent

High spatial and temporal resolution – challenging even for data-rich cities and countries

Opportunities to leverage emerging large-scale datasets: street-level images, satellite data, etc.

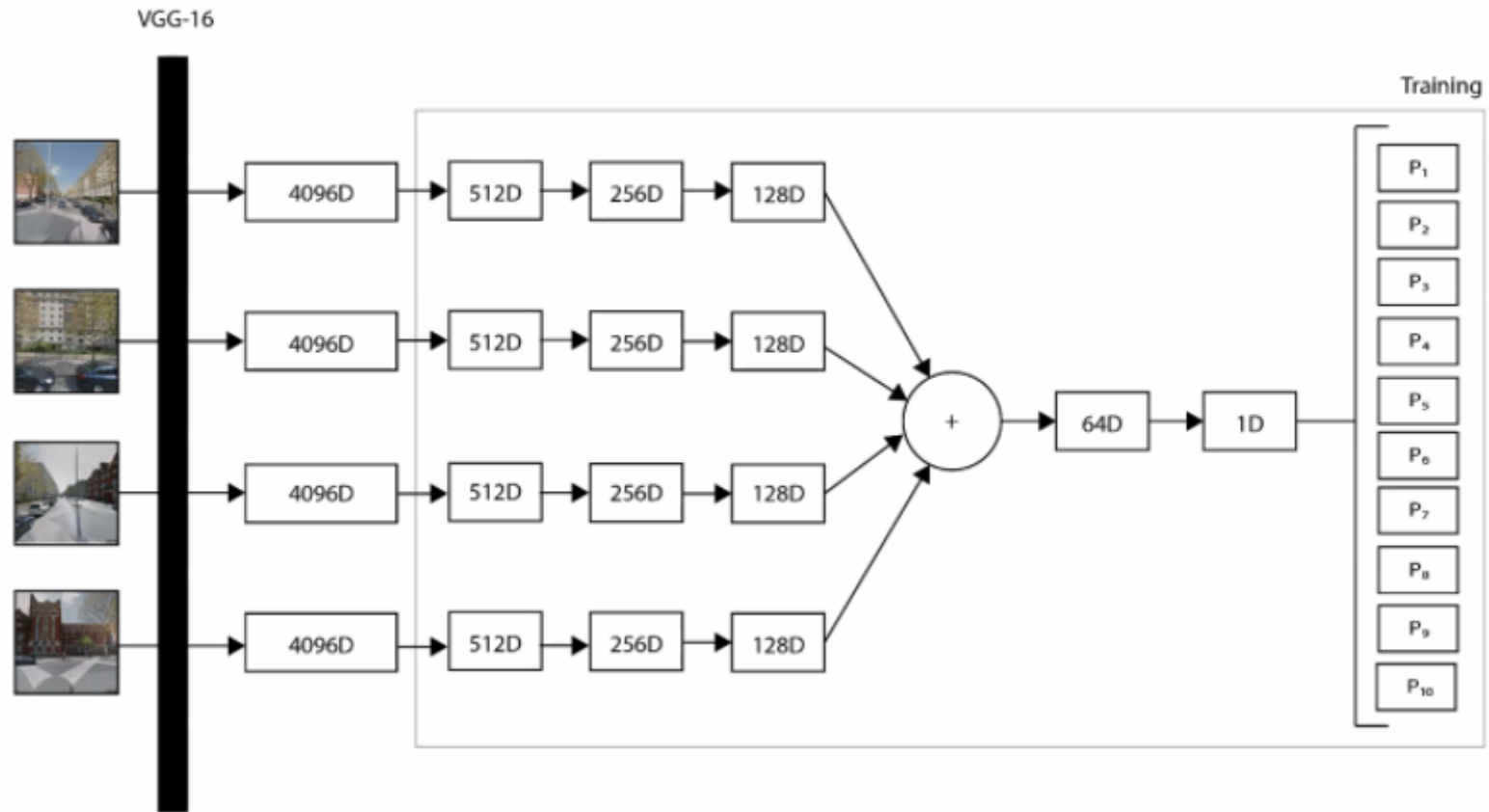
Example 2: measuring urban inequality

Case study: using street images to predict the decile of index of multiple deprivation (IMD)

~525,000 images from ~180,000 postcodes in London

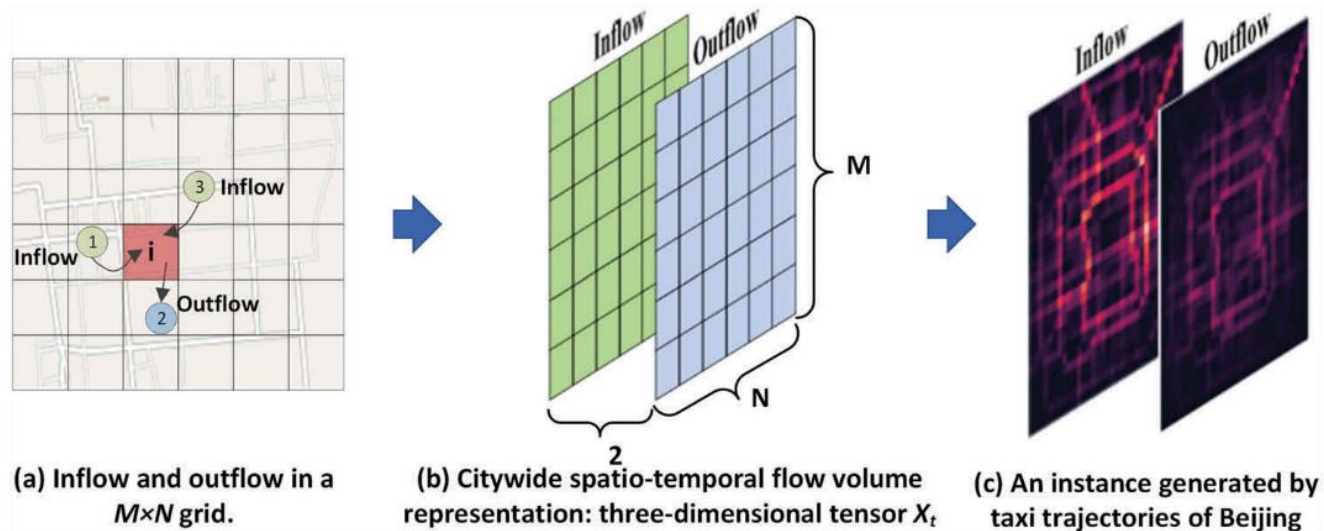


Example 2: measuring urban inequality



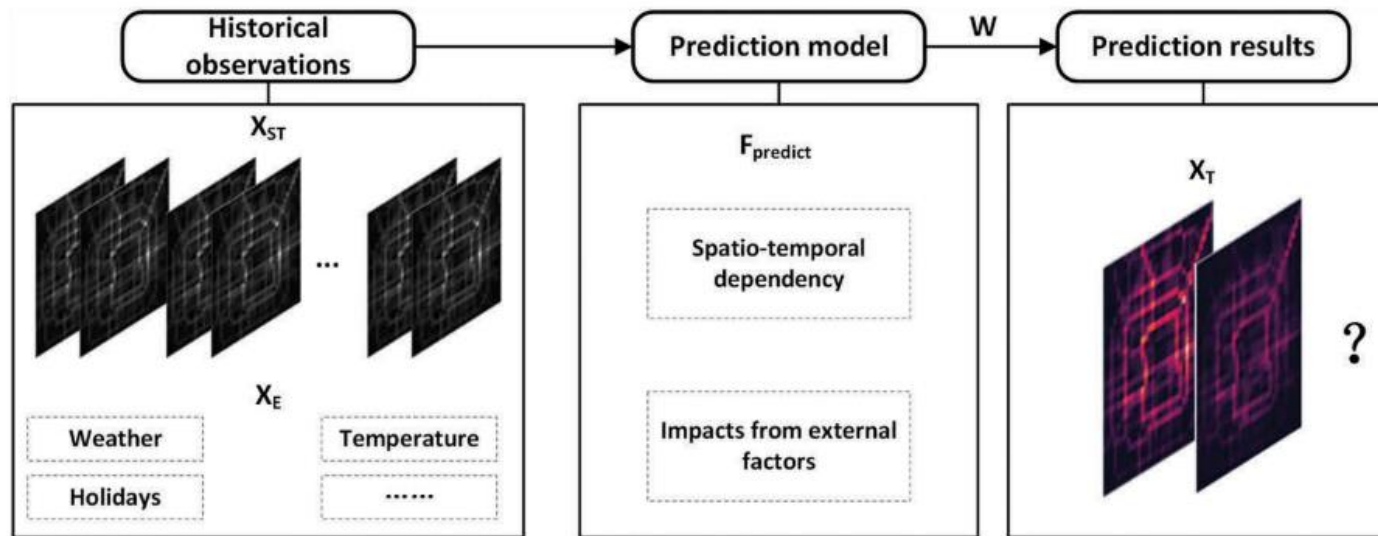
Example 3: spatio-temporal prediction of traffic flows in urban cells

The idea is to transform traffic flow into image-shape data and then apply CNN to predict flows.



Example 3: spatio-temporal prediction of traffic flows in urban cells

The idea is to transform traffic flow into image-shape data and then apply CNN to predict flows.

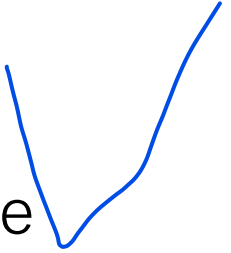


Summary – deep learning for images

DL provides powerful CNN-based tools to extract patterns and knowledge from images for varied purposes.

Most models are open-source and easy to reuse. Most models have similar structure (incl. Conv, Pooling, FC layers).

For most applications, you don't need to design a CNN from scratch. It is advisable to test classical models and adjust the models if needed.



Potential Problems with CNN

Garbage in, garbage out

If labels are not reliable, the model trained would be problematic.

“A flock of birds flying in the air”



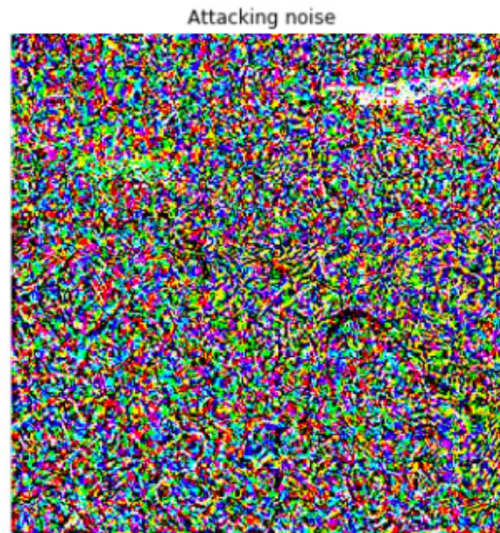
“A group of flowers in a field”



Adversarial Attacks



"A car"



"A toaster"

- You can fool a neural network classifier by adding noise to data.
- The noise might not be even visible to a human eye.

Adversarial Attacks



- Stop sign is no longer recognised by a neural network after a few stickers are placed on it.

Deep learning for text data

Keyword Detection

- Keyword extraction identifies the most important tokens or n-grams within a piece of text
- Input: a piece of text
- Output: list of most important tokens (words), optionally grouped by topics



Sentiment Analysis

- Sentiment Analysis assigns a sentiment score to each word in a piece of text

I love data science, and our teachers are awesome	+4 (Strongly positive)
Beer is disgusting, why do people even like it?	-1 (Weakly negative)
It's so great that my train is late every single day	+1 (Weakly positive)

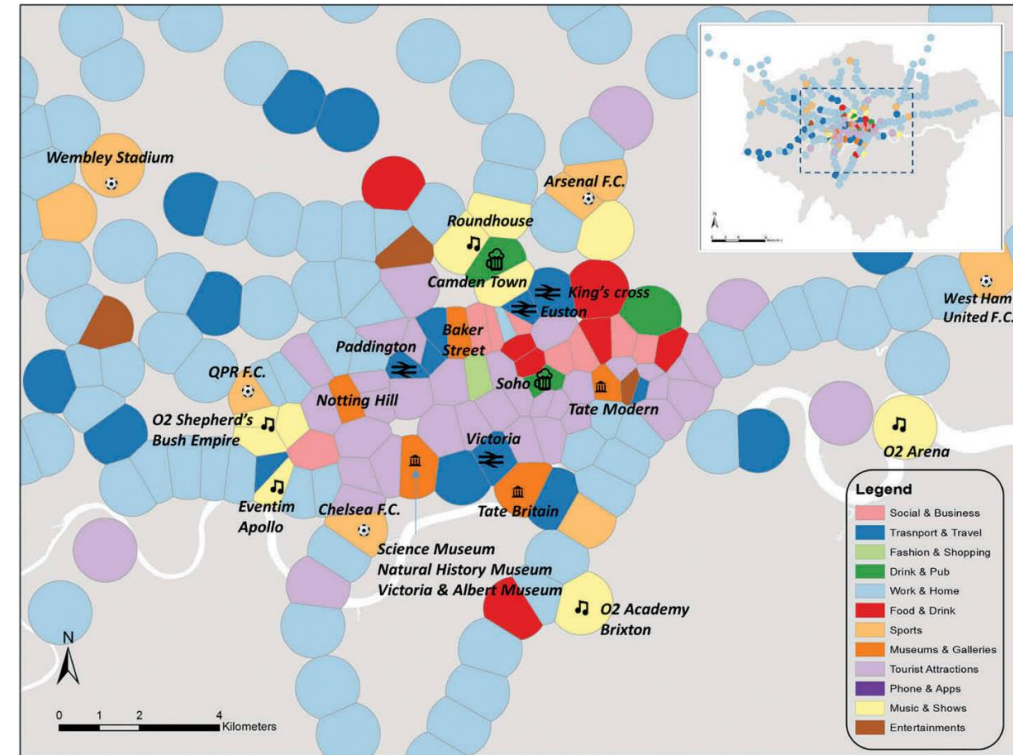
- Input: a piece of text
 - Output: a piece of text with a sentiment score assigned to each word
-

Topic Extraction



- Topic modelling extracts topics from a collection of documents.
- Input: a collection of documents
- Output: a list of topics with words that belong to them

Example:
Dominant topics on Tweets around the stations



Lai, J., Cheng, T., and Lansley, G. (2017) Improved targeted outdoor advertising based on geotagged social media data. *Annals of GIS*, 1–14. <https://doi.org/10.1080/19475683.2017.1382571>

Summary

- Unstructured data are common in the big data era, including images, text, etc.
- Deep learning provides powerful tools for utilising unstructured datasets. Most tools are open-source and ready to use.
- For applications, you don't need to fully understand the mathematics of these DL models.
- We cover the foundation and classical models of CNN here. More details can be found in the Coursera module “Convolutional Neural Networks”.

<https://www.coursera.org/learn/convolutional-neural-networks>

Workshop

In this week's workshop you will learn how to train your own deep learning models in Python for a range of tasks:

- Image classification
 - Face recognition
 - Semantic image segmentation
-