

PhD in Computer Science and Engineering
Bologna, April 2016

Machine Learning

Marco Lippi

`marco.lippi3@unibo.it`



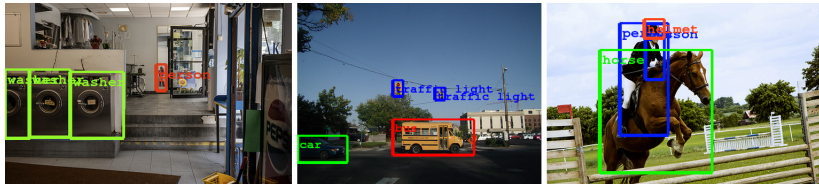
ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

Convolutional Neural Networks

Convolutional neural networks

Architecture inspired by biological processes, focused on **vision**:

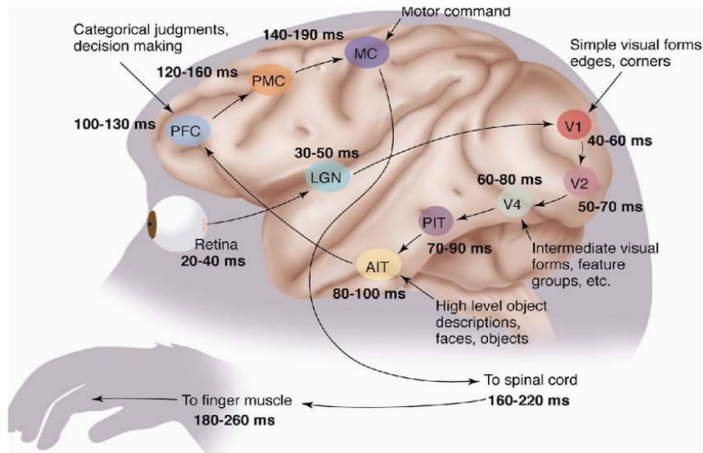
- neurons respond to overlapping regions in a **visual field**
- extremely **fast** computation (especially now with GPUs)
- pioneering models **back from the 80s-90s** !



[Figure from Google Research]

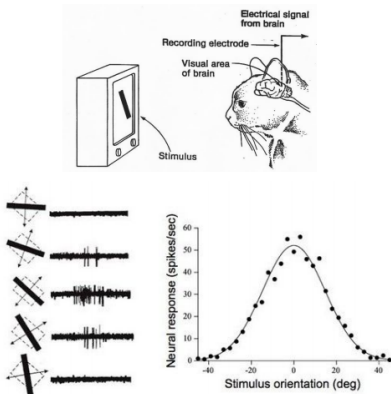
Quite a long history. . .

The human visual cortex is hierarchical



[Figure from nyu.edu, Simon Thorpe]

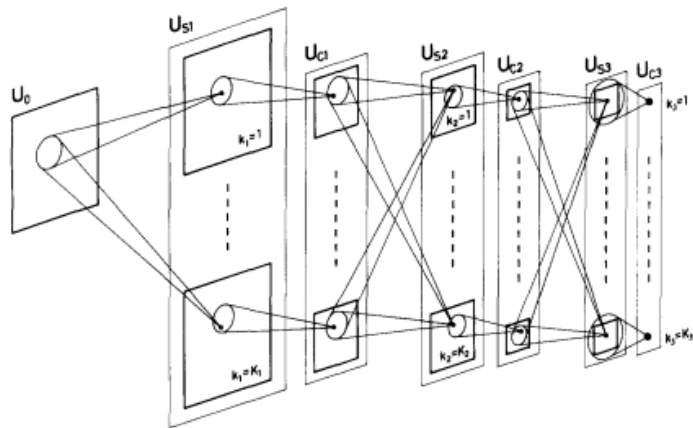
Hubel and Wiesel model (1962)



[Figure from nyu.edu]

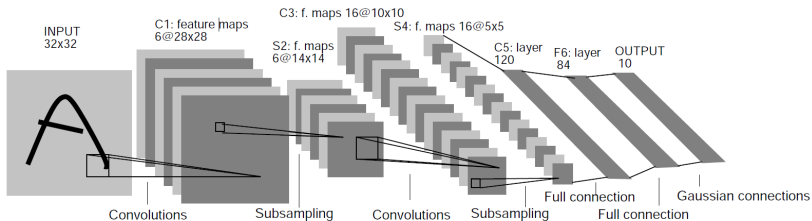
- **Simple cells** detect (edge-like) local features
- **Complex cells** receive input from simple cells and their receptive fields are spatially invariant

Cognitron and Neocognitron (Fukushima, 1974-1982)



[Figure from Fukushima, 1980]

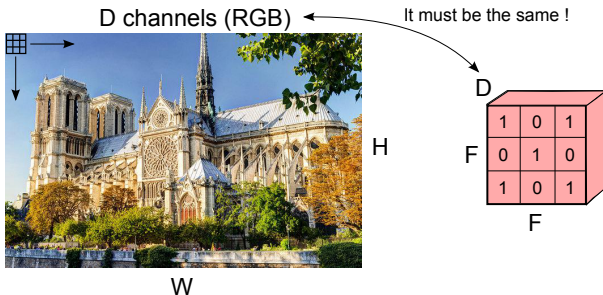
LeNet5 (LeCun et al., 1989)



[Figure from LeCun et al., 1989]

Receptive fields and convolutional filters

Convolutional filter: a matrix (tensor) of weights to be applied on the image to perform **convolutions**



Receptive fields and convolutional filters

Convolution between image patch and filter

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

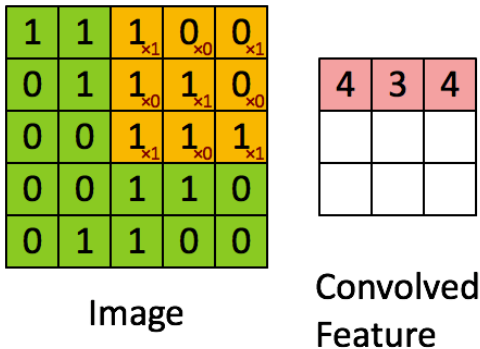
4	3	

Convolved
Feature

[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

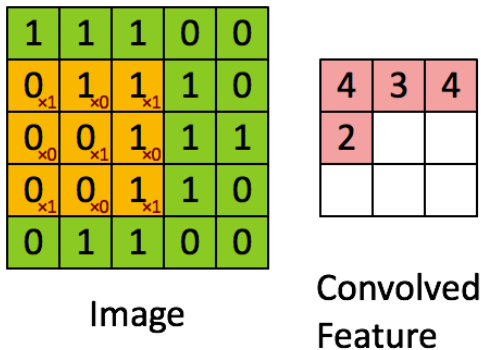
Convolution between image patch and filter



[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter



[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter

1	1	1	0	0
0	1 _{x1}	1 _{x0}	1 _{x1}	0
0	0 _{x0}	1 _{x1}	1 _{x0}	1
0	0 _{x1}	1 _{x0}	1 _{x1}	0
0	1	1	0	0

Image

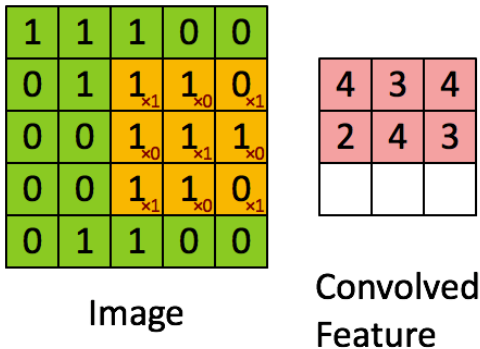
4	3	4
2	4	

Convolved
Feature

[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter



[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter

1	1	1	0	0
0	1	1	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0 _{x0}	0 _{x1}	1 _{x0}	1	0
0 _{x1}	1 _{x0}	1 _{x1}	0	0

Image

4	3	4
2	4	3
2		

Convolved
Feature

[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter

1	1	1	0	0
0	1	1	1	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0 _{x0}	1 _{x1}	1 _{x0}	0
0	1 _{x1}	1 _{x0}	0 _{x1}	0

Image

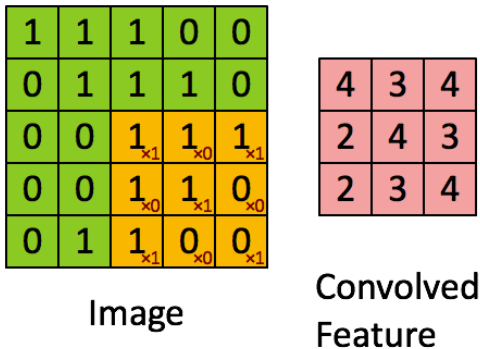
4	3	4
2	4	3
2	3	

Convolved
Feature

[Figure from <http://deeplearning.stanford.edu/>]

Receptive fields and convolutional filters

Convolution between image patch and filter



[Figure from <http://deeplearning.stanford.edu/>]

Stride

Hyper-parameter S indicating the “step” to be used when moving the filter on the image

- given a $W \times H$ image
- given a $F \times F$ filter
- $(W - F)/S$ and $(H - F)/S$ must be integers

Zero padding

Adding zeros along the border to allow convolutions on all pixels

- if $S = 1 \rightarrow$ zero padding with $(F - 1)/2$

Given a volume of size $W \times H \times D$

Choose hyper-parameters:

- # filters K
- filter dimension F
- stride S
- amount of zero padding P

Output is a volume of size $\hat{W} \times \hat{H} \times \hat{D}$ where:

- $\hat{W} = (W - F + 2P)/S + 1$
- $\hat{H} = (H - F + 2P)/S + 1$
- $\hat{D} = K$

Common settings: $K=2^m$, $F=3$, $S=1$, $P=1$ (or $F=5$, $S=1$, $P=2$)

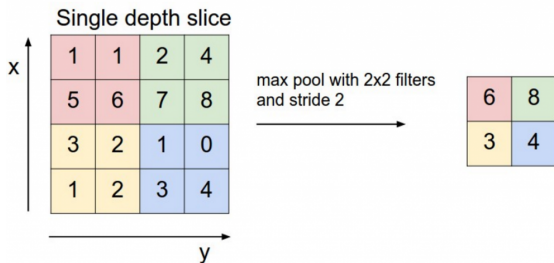
Typically placed **on top of a convolutional layer**:

- **drops to zero** negative inputs
- it is often added to further **augment non-linearity**
- operating on each activation map **independently**

Pooling Layer

An aggregation/subsampling function:

- operating on each activation map **independently**
- take the max/avg over a $M \times M$ filter with stride Z
- **no parameters to learn**
- just a **computation** above previous layer !



[Figure from <http://deeplearning.stanford.edu/>]

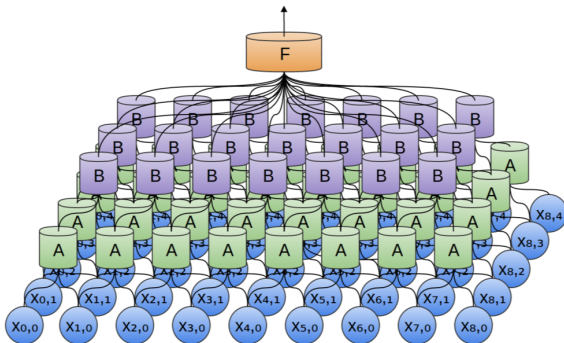
Recently introduced layer. . .

- operating on the output of max pooling
- subtracting mean and dividing by standard deviation of input
- this allows to obtain **brightness invariance**

Fully Connected Layer

Standard layer as in classic ANNs:

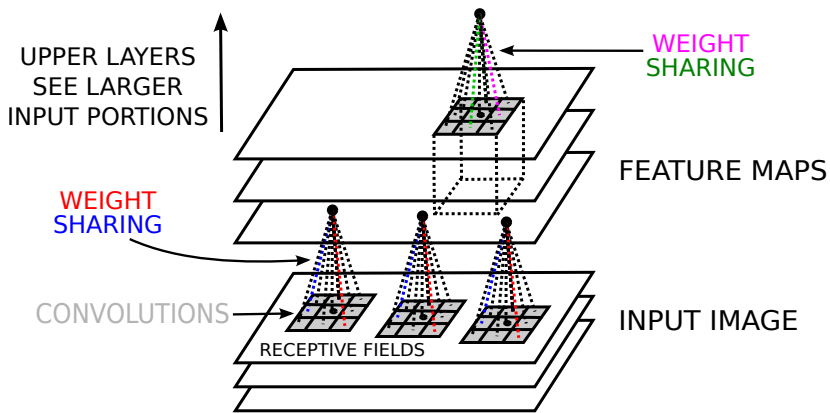
- **every** neuron connected to **every** neuron in previous layer
- typically implementing a **linear classifier**



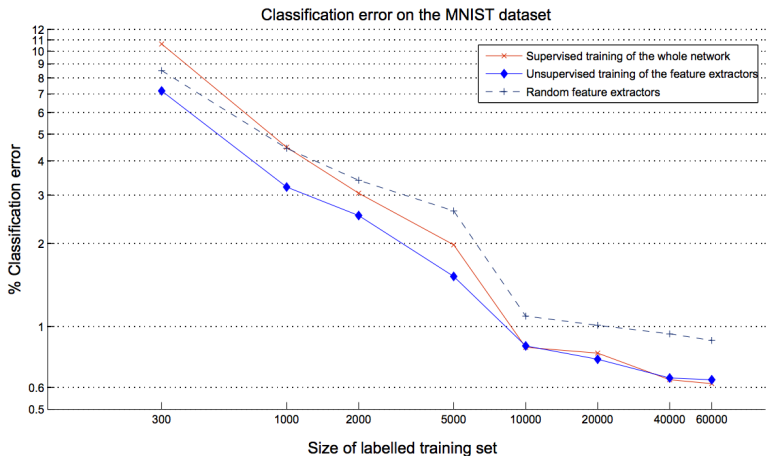
[Figure from colah.github.io]

Convolutional neural networks

One of the key advantages is to **share weights** between neurons !



A crucial advantage must be **in the structure** of a CNN !



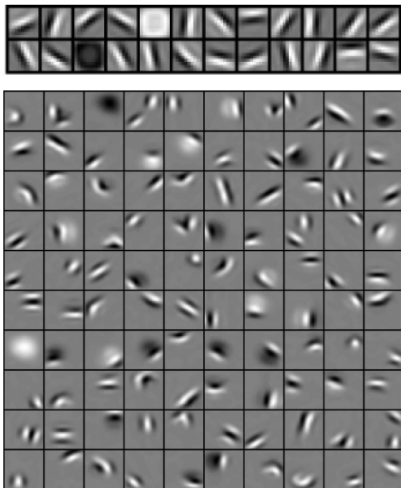
[Figure from Ranzato et al., 2007]

CNNs features have many nice properties:

- **compositionality** due to hierarchical structure
- **translation invariance** due to max pooling
- **scale invariance** via sub-sample processing

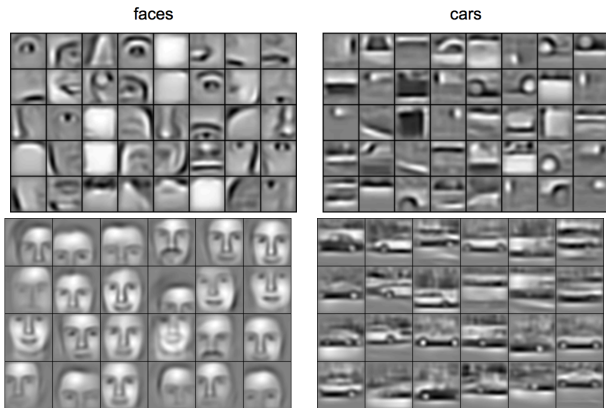
Feature extraction

Features obtained with various natural images



[Figures by Lee et al., 2009]

Features obtained with images belonging only to a given category



[Figures by Lee et al., 2009]

Features obtained with images belonging only to a set of categories

faces, cars, airplanes, motorbikes



[Figures by Lee et al., 2009]

The dream: build a computer vision system capable of **recognizing thousands of object categories**

- over 14 millions images
- over 20 thousands categories
- tagged via **crowdsourcing** !
- organized according to the **WordNet** noun hierarchy

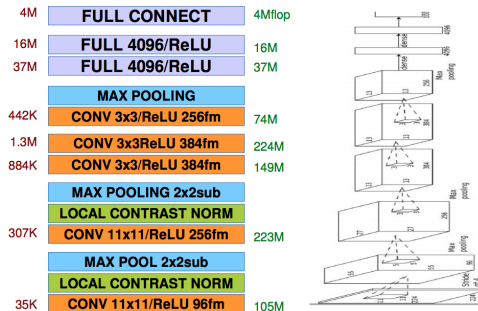


[Figure from vision.stanford.edu]

Annual competition:

ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)

Breakthrough in 2012 by AlexNet [Krizevsky et al., 2012]*



[Slide from YannLeCun]

*2012 paper with 4,377 citations. . .

Breakthrough in 2012 by AlexNet [Krizevsky et al., 2012]

- $\sim 60,000,000$ parameters
- $\sim 650,000$ neurons
- trained for **5-6 days** with 2 GPUs in parallel
- achieved a top-5 error rate of 18.2 % (second best 26.2 %)
reduced to 15.4 % with multiple models and pre-training
- achieved a top-1 error rate 40.7 %
reduced to 36.7 % with multiple models and pre-training

The ImageNet challenge

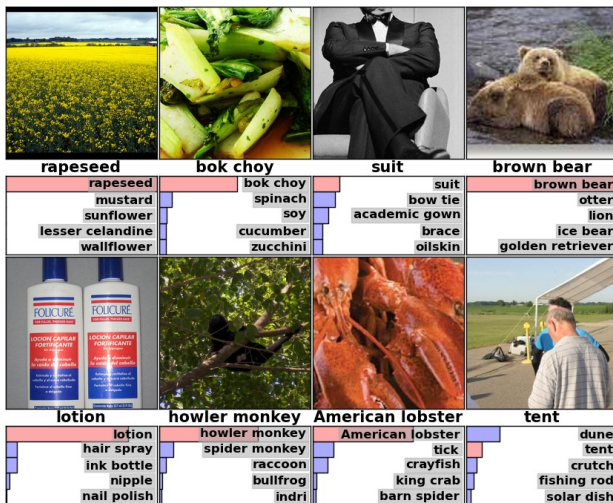
Predictions



[Figure from Krizhevsky et al., 2012]

The ImageNet challenge

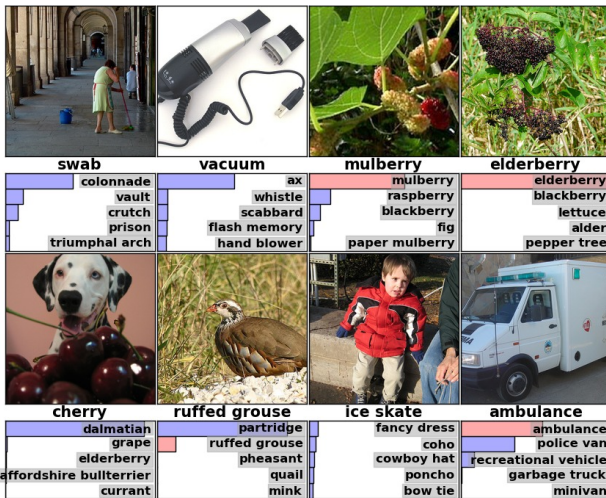
Predictions



[Figure from Krizhevsky et al., 2012]

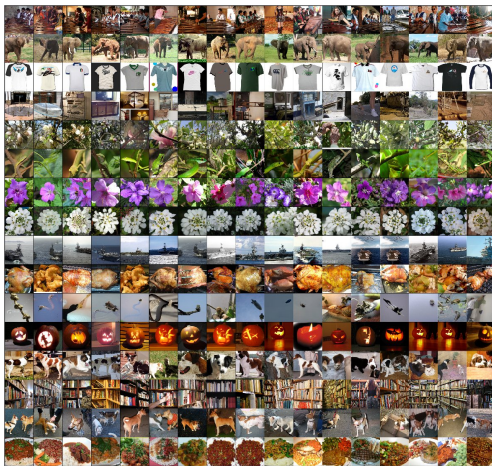
The ImageNet challenge

Predictions



[Figure from Krizhevsky et al., 2012]

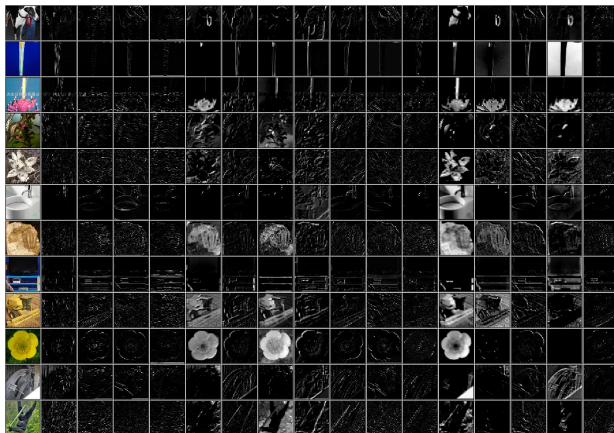
Retrieval



[Figure from Krizhevsky et al., 2012]

The ImageNet challenge

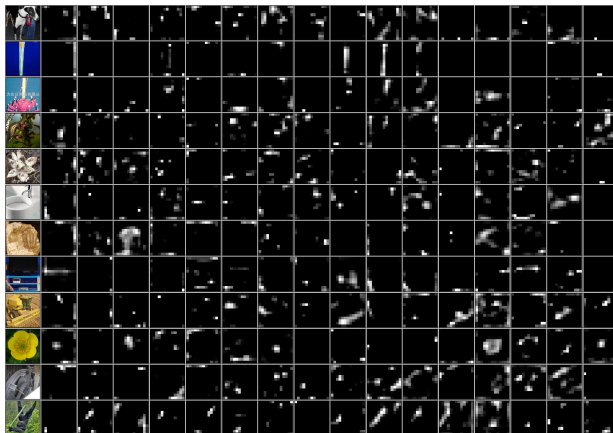
Feature maps (convolutional layer 1)



[Figure from Krizhevsky et al., 2012]

The ImageNet challenge

Feature maps (convolutional layer 1)



[Figure from Krizhevsky et al., 2012]

Features of the 2014 edition:

- 1.2 million images for training
- 50k images for validation
- 100k images for test
- 1,000 semantic categories
- Winner: **GoogLeNet** (22 layers !!!) → 6.67 % Top-5 error

In 2015 **ResNet** (Microsoft Research) → 3.6 % top-5 error !!!

- They employed a 152-layer net !!!
- They won all the tasks (localization, detection, segmentation)

A large number of **layers** implies a large number of **parameters**, thus a large number of **examples** needed to train the network.

What if one has only a small training set ?

- it has now become common to **pre-train** the network on a large dataset (e.g., ImageNet)
- **fine-tuning** is then performed on the (smaller) dataset related to a specific task

This is an instance of **transfer learning** !

[Torrey & Shavlik, 2009]

“Transfer learning is the improvement of learning in a new task through the transfer of knowledge from a related task that has already been learned.”

- **Complex but crucial** for any machine learning system
- One of the tasks that **humans** are very good at
- Need to map features/relations **across domains**

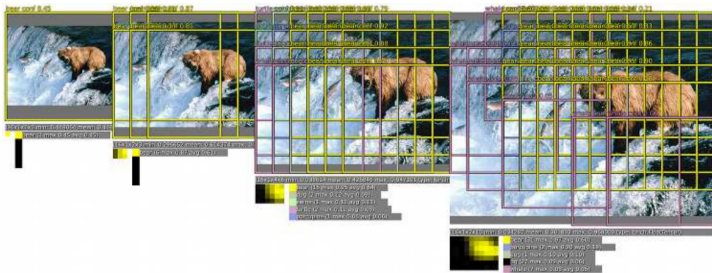
Two possible scenarios:

- 1 use a CNN trained on ImageNet as a **fixed feature extractor**
→ this means to remove the last fully connected layer
- 2 also **fine-tune** some/all the weights of the CNN with the smaller dataset → be careful of **overfitting** !

A fine-tuning of the whole network might be required if our own dataset is **much different** from the one used in pre-training, **and not too small...**

Quite straightforward adaptation of CNNs:

- sliding window over the image
- multi-scale resolution



[Figure by Sermanet et al., 2013]

Most of the tricks common to other kinds of deep networks

- dropout
- data augmentation (jittering, noise injection, ...)
- weight decay
- sparsity of hidden units
- carefully choose learning rate

In addition...

- visualize **feature maps**
- visualize **parameters** (filters)

Visual Turing Challenge



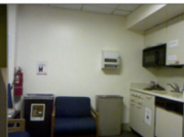
QA: (What is behind the table?, window)
Spatial relation like "behind" are dependent on the reference frame. Here the annotator uses observer-centric view.



QA: (what is beneath the candle holder, decorative plate)
Some annotators use variations on spatial relations that are similar, e.g. 'beneath' is closely related to 'below'.



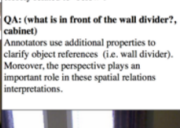
The annotators are using different names to call the same things. The names of the brown object near the bed include 'night stand', 'stool', and 'cabinet'.



Some objects, like the table on the left of image, are severely occluded or truncated. Yet, the annotators refer to them in the questions.



QA: (what is behind the table?, sofa)
Spatial relations exhibit different reference frames. Some annotations use observer-centric, others object-centric view
QA: (how many lights are on?, 6)
Moreover, some questions require detection of states 'light on or off'



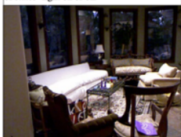
QA: (what is in front of the wall divider?, cabinet)
Annotators use additional properties to clarify object references (i.e. wall divider). Moreover, the perspective plays an important role in these spatial relations interpretations.



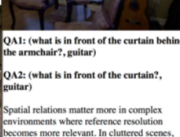
QA1: (How many doors are in the image?, 1)
QA2: (How many doors are in the image?, 5)
Different interpretation of 'door' results in different counts: 1 door at the end of the hall vs. 5 doors including lockers



QA: (How many drawers are there?, 8)
The annotators use their common-sense knowledge for amodal completion. Here the annotator infers the 8th drawer from the context



QA: what is at the back side of the sofas?
Annotators use wide range spatial relations, such as 'backside' which is object-centric.



QA1: (what is in front of the curtain behind the armchair?, guitar)
QA2: (what is in front of the curtain?, guitar)

Spatial relations matter more in complex environments where reference resolution becomes more relevant. In cluttered scenes, pragmatism starts playing a more important role



QA: (What is the object on the counter in the corner?, microwave)
References like 'corner' are difficult to resolve given current computer vision models. Yet such scene features are frequently used by humans.

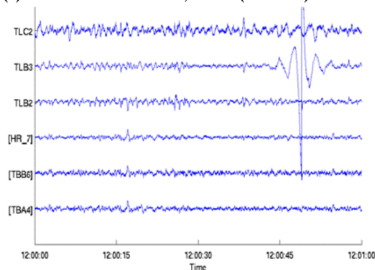


QA: (How many doors are open?, 1)
Notion of states of object (like open) is not well captured by current vision techniques. Annotators use such attributes frequently for disambiguation.

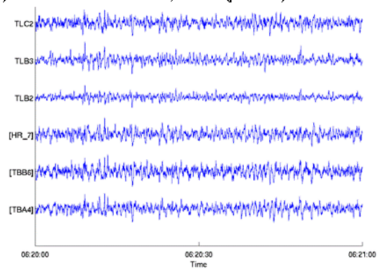
Prediction of epilepsy seizures from intra-cranial EEG

- Temporal CNNs [Mirowski et al., 2008]

(a) EEG on 06-Dec-2001, 12:00 (interictal)



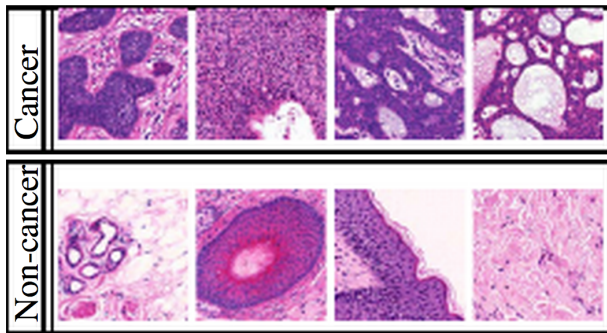
(c) EEG on 12-Dec-2001, 06:20 (preictal)



[Figure by Mirowski et al.]

Cancer diagnosis and classifications

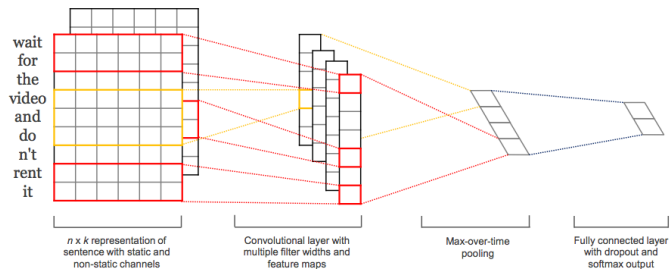
- Auto-Encoders [Fakoor et al., 2013]
- Convolutional Neural Networks [Cruz-Roa et al., 2013]



[Figure by Cruz-Roa et al.]

Sentence classification

- Sentiment analysis
- Question classification
- Subjectivity score
- ...



[Figure by Kim, 2014]

Word Embeddings, Language Models, etc.

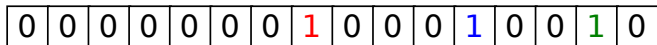
Modeling language

In Natural Language Processing/Understanding (NLP/NLU), one crucial element is to **represent sentences** for the desired task

- Classification
- Segmentation
- Tagging
- ...

A classic representation is that of Bag-of-Words (Bow)

The cat is walking in the garden



VECTOR LENGTH =
VOCABULARY DIMENSION

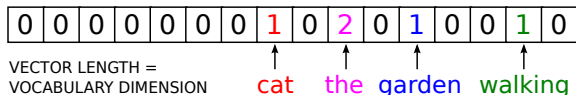
cat

garden

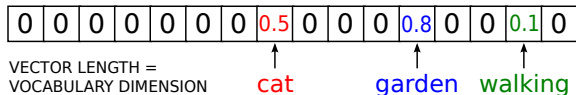
walking

Other BoW variants: consider **frequencies**

- frequency of a word **within a document**
 - Term Frequency (TF)



- frequency of a word **within a corpus**:
 - Inverse Document Frequency (IDF) \rightarrow $TF \times IDF$



Rare words in common are much more significant...
...But still, **it is not enough** !

A classic problem: not capturing **lexical** and **semantic** similarity !

```
The cat is walking in the garden
A dog was running towards the park
```

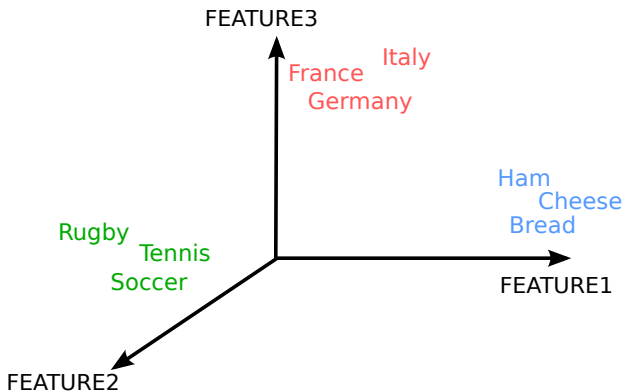
Almost **no similarity**: only article “the” in common...

A lot of approaches have tried to include lexical/semantic features

- Use of ontologies (e.g., WordNet)
- Analysis of co-occurrences
- Word disambiguation (e.g., bank: river/finance ?)

“You shall know a word by the company it keeps” (J.R. Firth, 1957)

- Use **context** to learn **word representations** (embeddings)
- A word will be represented by a **dense real-valued vector**



Learn a **probability distribution** over sequences of words

The probability of **observing a sequence of words** w_1, \dots, w_m is:

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1})$$

A classic approach employs **n-grams**:

$$P(w_m | w_{m-1}, w_{m-2}, \dots, w_{m-n+1}) = \frac{\text{count}(w_{m-n+1}, \dots, w_m)}{\text{count}(w_{m-n+1}, \dots, w_{m-1})}$$

Just counting words. . .

The idea of neural embeddings dates back to over one decade:

- **predict next word** given current (and previous ones)

$$f(w_t, \dots, w_{t-n+1}) = P(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-n+1})$$

Function f is decomposed in two parts:

- a mapping $C : |V| \rightarrow \mathbb{R}^d$
- an ANN: $g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$

C is the matrix of **word embeddings** (or word vectors)

Neural Network Language Model

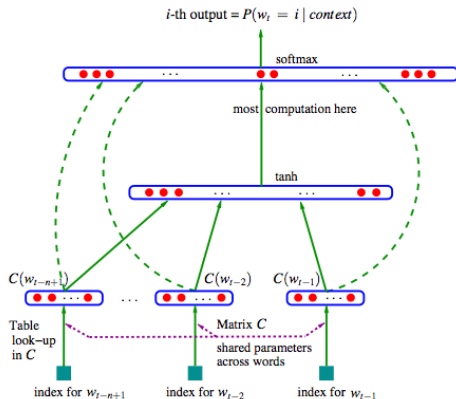


Figure 1: Neural architecture: $f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$ where g is the neural network and $C(i)$ is the i -th word feature vector.

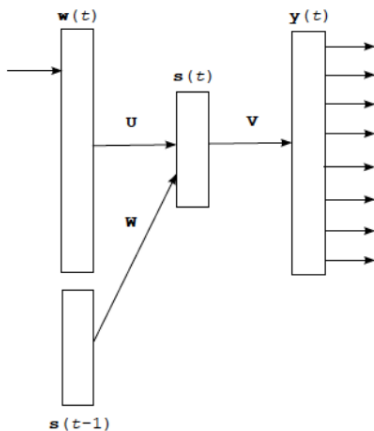
[Figure by Bengio et al., 2003]

Additional considerations...

- Output layer has $|V|$ size \rightarrow quite expensive
- Reducing to $\log |V|$ via **hierarchical softmax**
- Partitioning the output spaces into a hierarchical structure
- Using a **bit vector encoding** of words (e.g., binary Huffman tree)
- Exploiting WordNet hierarchy

Recurrent Neural Network Language Model

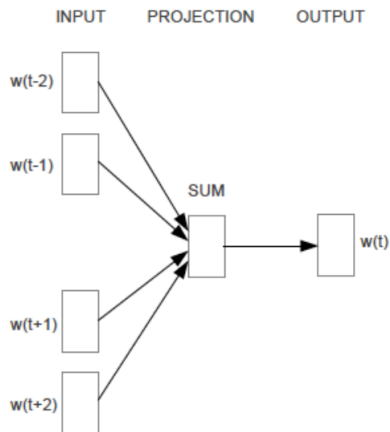
Next lecture !



[Figure by Mikolov et al., 2010]

Continuous bag-of-words

Predict word given past and future context



[Figure by Mikolov et al., 2013]

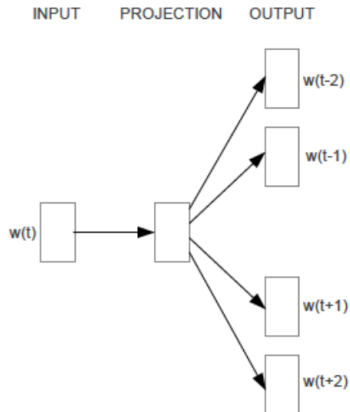
Continuous bag-of-words

Objective function: maximize the **average log probability**

$$\frac{1}{T} \sum_{t=k}^{T-k} \log p(w_t | w_{t-k}, \dots, w_{t+k})$$

Skip-gram

Predict contextual words of a given word



[Figure by Mikolov et al., 2013]

Skip-gram

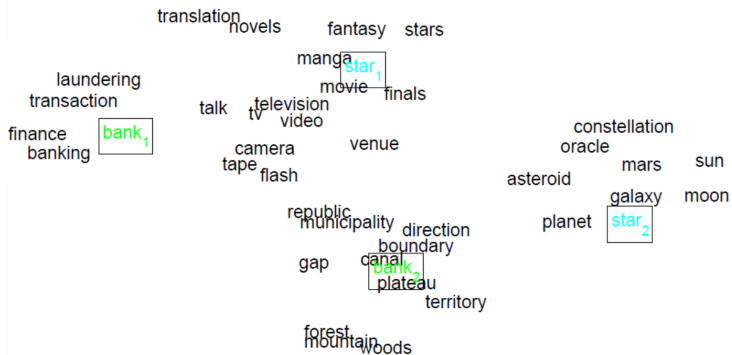
Objective function: maximize the **average log probability**

$$\frac{1}{T} \sum_{i=1}^T \sum_{-c \leq h \leq c, h \neq 0} \log p(w_{t+h} | w_t)$$

Trade-off to choose size of context c :

- if larger, the model is more accurate. . .
- . . . but it is computationally more expensive

Word vectors (word2vec)

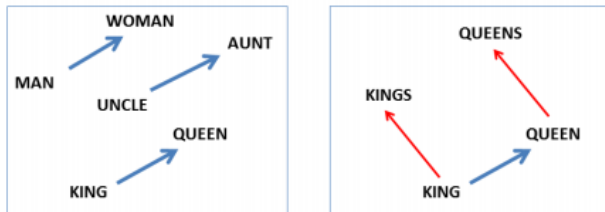


The power of word embeddings

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
454	1973	6909	11724	29869	87025
AUSTRIA	GOD	AMIGA	GREENISH	NAILED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLUISH	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/S
ITALY	SATAN	IPOD	PURPLISH	POPPED	BAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATS
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	KBIT/S
NORWAY	VISHNU	HD	GRAYISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARVATI	GEFORCE	SILVERY	SLASHED	GBIT/S
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

[Table by R. Collobert et al., 2011]

The power of word embeddings



[Figure by T. Mikolov]

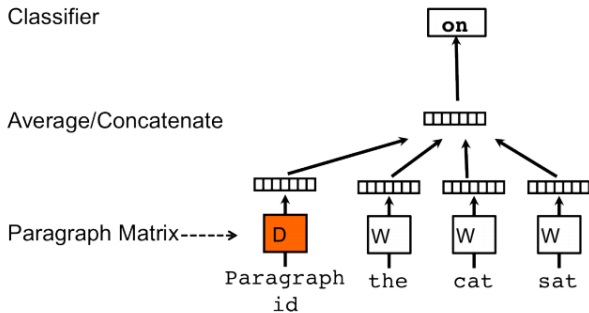
Transfer learning !

- named entity recognition
- part-of-speech tagging
- parsing
- semantic role labeling
- machine translation

Sentence, paragraph, document vectors (doc2vec)

Generalization of word embeddings to any text

- W is the word embedding matrix
- D is the paragraph/document embedding matrix



[Figure by Le & Mikolov, 2014]

Pre-trained word vectors:

- **GloVe** (Global Vectors for word representation) @ Stanford
<http://nlp.stanford.edu/projects/glove/>
Different versions (Wikipedia, Twitter, Common Crawl)
- **word2vec** @ Google
<https://code.google.com/archive/p/word2vec/>
Trained on GoogleNews
Naming version trained on FreeBase
- ...