Feb 9th, 2017 Convolutional Neural Networks

Motivation: MLPs and AEs have too many parameters. (= weights)
eg.
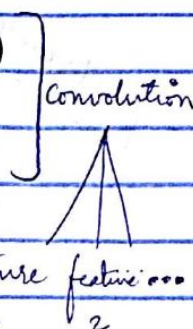* a tiny pic of $30 \times 30$ resolution: for a neural network, it is 900 inputs.
  first layer will have (80,000 weights)!!
* $224 \times 224$ images of faces: 65,000 pixels. ∴ 65,000 inputs for MLP.
  AE is not designed for this kind of task. 1st layer = (600,000 weights)

we combined 3 ideas:
① Receptive fields (neighborhoods) in the vision system (of cats)
② weight sharing (for a bunch of inputs)
   (eg. share 9 weights for 230,000 numbers)
③ down-sampling (idea from A.E.)
   = pooling

convolution

feature   feature ...
   1        2

So far, feature extraction was manually engineered. (SIFT)
                                    (manually handrafted features)
   ⎧ (Network shrud itself extract features. Not rely on anybody else
Breakthrough.                                              outside )

## What is Convolution?

physical system: an input acts on a linear physical
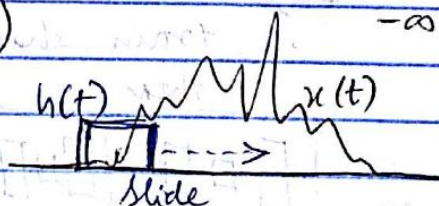              system to produce an output.
$x(t)$ = input                $h(t) \stackrel{?}{=}$ unit impulse response of system
$y(t)$ = output

$$y(t) = x(t) * h(t) = \int_{-\infty}^{+\infty} x(\tau) \, h(t-\tau) \, d\tau$$

$$y(n) = \sum_{-\infty}^{\infty} x(k) \, h(n-k)$$

$\tau$ = dummy variable

$h(t)$ ⟍⟋⟍⟋ $x(t)$
       slide

$$h = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

$X = \begin{bmatrix} & 200 & \\ & \boxplus & \end{bmatrix} 200$

$\sum$
| 10 | 10 | 10 |
| 50 | 50 | 50 |
| 100 | 100 | 100 |
$*$
| -1 | 0 | 1 |
| -1 | 0 | 1 |
| -1 | 0 | 1 |
$= \sum$
| 10 | 0 | 10 |
| -50 | 0 | 50 |
| -100 | 0 | 100 |
$= 0$

vertical edge detector                no vertical edges
              (has horizontal edge)          ∴ 0

## Weights as filters?

$x_1 \xrightarrow{w_1}$
$x_2 \xrightarrow{w_2}$  ◯ $\to y = w_1 x_1 + w_2 x_2 + w_3 x_3.$
$x_3 \xrightarrow{w_3}$

Draft idea:

input image → convolution (features) feature maps → MLP → y

(too many weights)

if: many images.

way more ip than MLP can handle.

(image comes in) → many neurons → convolution → pooling → down sampled feature maps → MLP → y.

image → convolⁿ + pooling → convolⁿ pooling → convolⁿ pooling → fully connected layer → y

Pooling (= down sampling)

| 1 | 2 | 2 | 1 |
|---|---|---|---|
| 3 | 4 | 3 | 2 |
| 5 | 6 | 5 | 7 |
| 10 | 8 | 9 | 8 |

2×2

| 4 | 3 |
|---|---|
| 10 | 9 |

max pooling
(maximum in each 2×2 block)

Convolution = getting features.

[ Pooling = Reduce dimensionality
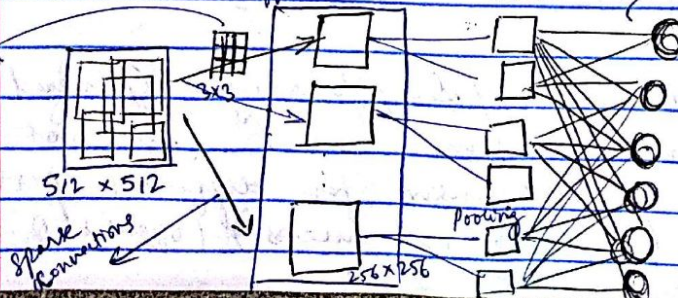  MLP layer = classification

MLP: Can never handle raw data.
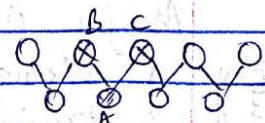∴ it needs features. Shallow structure

CNN • Sampling makes it invariant to translation. ∴ many samples with diff. convolution mask.

dense connections.

512 × 512

3×3

sparse connections

256×256

pooling

B   C
⊗   ⊗
   A

if: change Ⓐ, only Ⓑ & Ⓒ are changed. = weight sharing.
Local Sparsity

earlier: neural network was manually derived of but, the filter was automatically organised in neural network — this filter was in D — of part 2
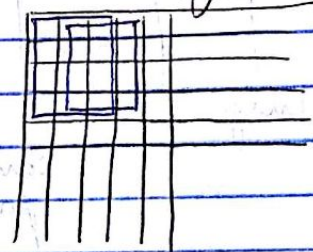
Feed forward Feature extraction
* convolve input with learned filters.
* non-linearity (logistic)
* spatial pooling
* normalization          (usually you convolve ⟶ Normaliz ⟶ then pool

(weight sharing = filter sharing)

Convolution
* dependencies are local.
* translation invariance - immensely important
* few parameters (due to weight sharing)
* stride (= overlap)

Non-linearity
$\hookrightarrow$ sigmoid $f(x) = 1/(1+e^{-x})$

$\hookrightarrow$ Rectified linear unit
       (ReLU)
   $f(x) = max(0, X)$
   This simplifies Backprop.

(has an overlap of 1 column. (= stride = 1)
if required, we can have a stride of
2 or more. If it is 1, then highly
probably computationally expensive)

Common Models
- AlexNet      (2012) winner on imageNet
- ZF Net       (2013)
- GoogleNet    (2014)
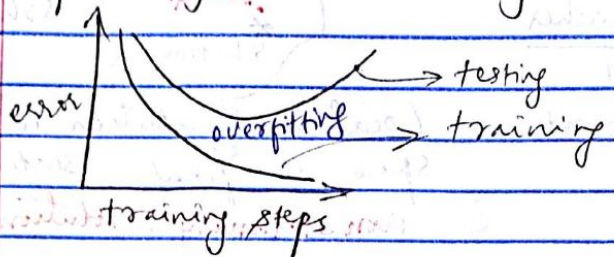- VGG          (2014 runner up)
- ResNet       (2015) winner on imageNet.

Major issues with CNN/Deep learning in General
* design requires some care    (but not THAT bad)
* training is still expensive for difficult problems.
* generalization (inspite of all success) (opposite of overfitting)
    biggest problem

difficulty ⇓

# Regularization

Any change to increase generalization capability (but not necessarily reducing training error). It is about the error of testing. Not training.

error ↗ overfitting → testing
→ training
training steps

## Types of Regularization:

① Data Augmentation
- translation, scaling, rotation
(increase the # of training data)

② Early Stopping

error

Testing
Training

→ stop here. don't train further

③ Dropout

- randomly remove neurons
- neurons become more insensitive to other neurons. (avoids overfitting)
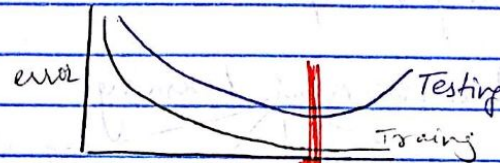(learn local dependencies. They become more independent)
- equivalent of averaging several models.
(like training multiple N.N.)    (ensemble).

④ Weight Penalties ($L2$ & $L1$ regularization)

Total error $E = error + \|w\|^2$ (size of weights)

$E = error + \|w\|$    reduce the error, but do it using very small weights.

(already minimized)    force it to choose small weights.

to reduce $E$, the only option is to reduce $\|w\|$.