

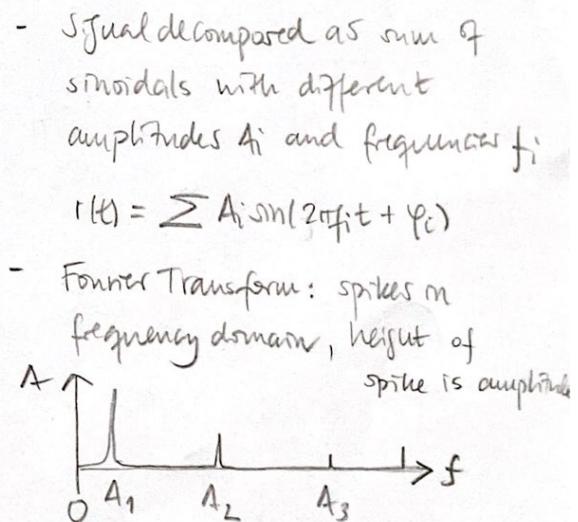
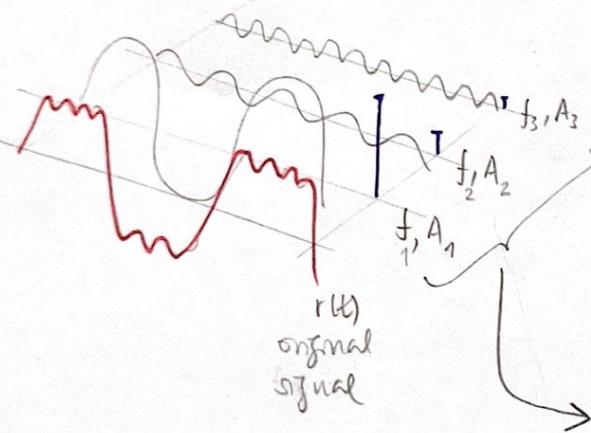
① INTRODUCTION TO COMPUTER VISION

IMAGE REPRESENTATION & CLASSIFICATION (see TXT & notebooks)

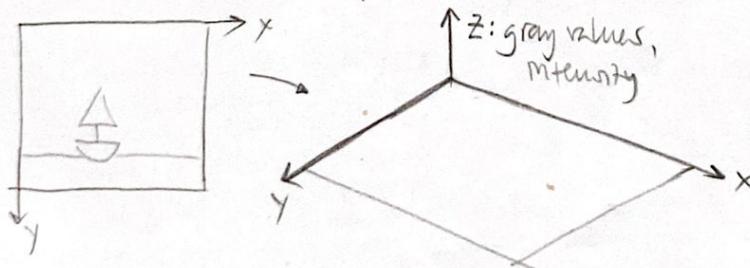
- Pixelmaps
- Thresholding
- Color maps / color transformations (HSV)
- Image features (e.g. brightness)
- Classification example (& standardization of images): Day and Night image classification
- Image processing & Computer Vision Pipeline

CONVOLUTIONAL FILTERS & EDGE DETECTION

Fourier Transforms



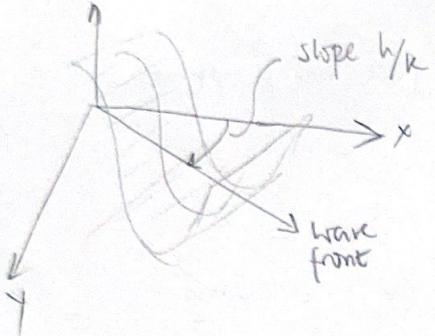
- Image: it's also a varying function but not in time, in 2D space



- an image can be understood as a 2D surface
 - its decomposition elements would have this form:
- $a_i \sin(h \cdot x + k \cdot y)$
- amplitude
- frequency (of repetition of the wave in x/y dir)
- spatial pixel coords

$$z = \sum a_i \sin(h_i x + k_i y)$$

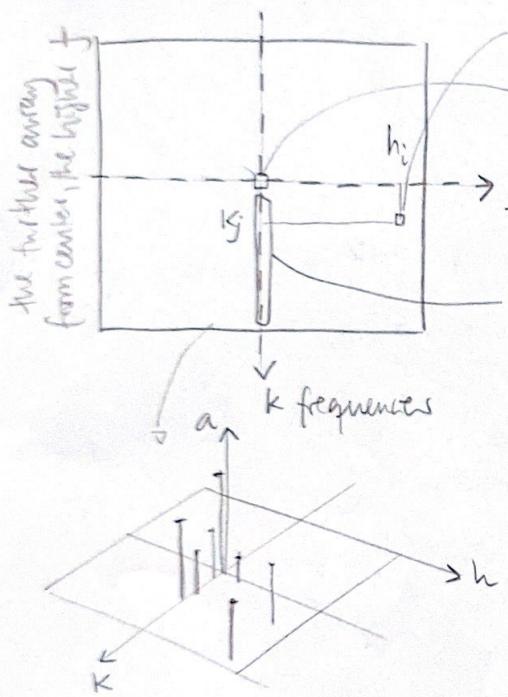
$h_i, k_i \neq 0 \Rightarrow$ wave i travels in angle with slope h_i/k_i (it's front)



- Therefore, an image is decomposed as a sum of 2D sine waves.

the amplitude²
is the intensity

- Interpretation of the FT: the brighter the pixels, the bigger their amplitude



amplitude for frequency pair (h_i, k_j)

$(h_i, k_j) = (0, 0)$: usually bright: flat plane with no ripples - its intensity/brightness is the average value of the gray values on the image

waves with $h=0$ frequencies; usually there's many white pixels here (and also in $k=0$), because human environments have many vertical ($h=0$) horizontal ($k=0$) lines

→ the image is the sum of all these 2D sine waves, with this amplitudes a and frequencies h & k

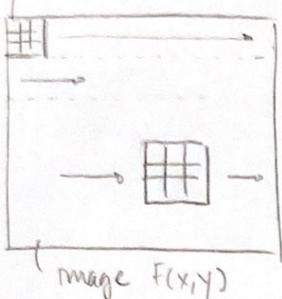
- Practical tricks

- If FT image brighter in center \Rightarrow image mostly solid colors, few texture (low f)
- If FT image horizontal line in center \Rightarrow vertical structures in image, e.g. human beings!
- " - vertical line - " \Rightarrow horizontal structures in image

NOTE: in the FT image, each pixel is not linked to the same pixel in the original image... in FT images pixels are frequency pairs and their value is the amplitude associated to the frequency!

Convolution

convolutional kernel K



$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

$$0 + (-1) + 0 + (-1) + 4 + (-1) + 0 + (-1) + 0 = 0$$

weights

Convolution kernel
for edge detection

$$K * F(x_i, y_i)$$

convolution
operator: kernel
is applied on all
pixels of the image

$$\begin{bmatrix} K \\ F(x_i, y_i) \end{bmatrix}$$

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 120 & 140 & 120 \\ 225 & 220 & 205 \\ 255 & 250 & 230 \end{bmatrix}$$

(x_i, y_i)
is center
pixel

the sum
of all elements
must be 0,
because this
filter is computing
differences between
neighbor pixels!

$$0 \cdot 120 + (-1) \cdot 140 + 0 \cdot 120 \\ + (-1) \cdot 225 + 4 \cdot 220 + (-1) \cdot 205 \\ + 0 \cdot 255 + (-1) \cdot 250 + 0 \cdot 230 = \underline{\underline{60}}$$

- Kernel cannot be overlaid perfectly on edges & corners - what to do?
 - extend: border pixel values extended outwards
 - padding: 0 values (black) set outwards
 - cropping: borders are removed \rightarrow output image slightly smaller
- Weights: the higher the relative value, the higher the importance
 - \hookrightarrow their spatial disposition defines the feature we want to remove/enhance.

Sobel Filters

Approximation of image gradient

$$\text{gray scale image, } F \quad \xrightarrow{\quad \quad} S_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \quad \xrightarrow{\quad \quad} \boxed{11}$$

intensity
image derivative in X direction,
emphasis of VERTICAL lines

$$\xrightarrow{\quad \quad} S_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \quad \xrightarrow{\quad \quad} \boxed{-}$$

image intensity derivative in Y dir,
emphasis of HORIZONTAL lines

The strongest edges are found with the magnitude:

$$|F_x'| = |S_x \cdot F| = \sqrt{\text{Sobel}_x^2}$$

$$|F_y'| = |S_y \cdot F| = \sqrt{\text{Sobel}_y^2}$$

$$|F_{xy}'| = \sqrt{\text{Sobel}_x^2 + \text{Sobel}_y^2} \quad \rightarrow \text{gradient magnitude in both directions}$$

Direction of gradient: $\arctan(\text{Sobel}_y / \text{Sobel}_x)$ \rightarrow to find lines with an angle upwards...

Lowpass Filters \rightarrow usual procedure:

- ① apply lowpass filter to remove noise
- ② apply highpass filter to enhance features

They remove high frequency components; usually some sort of averaging is done with pixel values

$$\frac{1}{9} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \rightarrow \frac{1}{9} \cdot (1+1+1+\dots+1) = \underline{1} \quad \text{sum must be 1!}$$

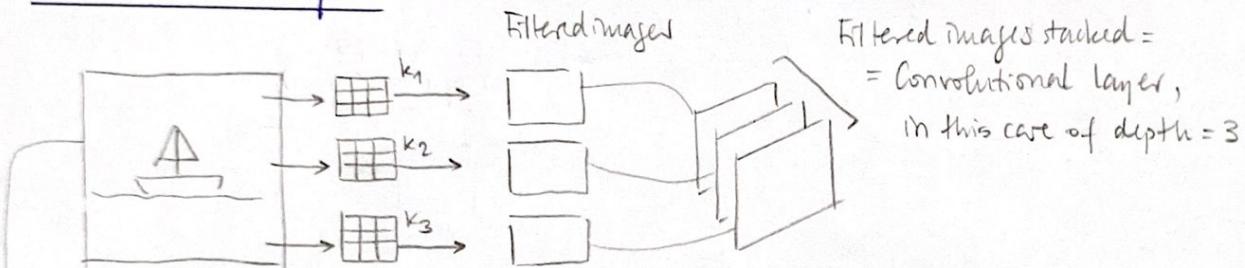
$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \rightarrow \Sigma = 1!$$

since it's so common, it has its own function in OpenCV
it's basically a weighted sum

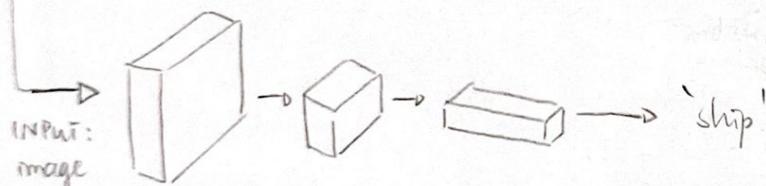
Gaussian Blur: high-freq.

removed BUT edges are preserved!

Convolutional Layers



- A convolutional neural network (CNN) can have several convolutional layers, each with a depth; some other layer types are also possible



- Weights of convolutional layer are learned during training, but basically they converge to be
 - lowpass filters, e.g. to compress image by filtering out irrelevant info
 - highpass filters, e.g. to enhance features like edges

Canny Edge Detector

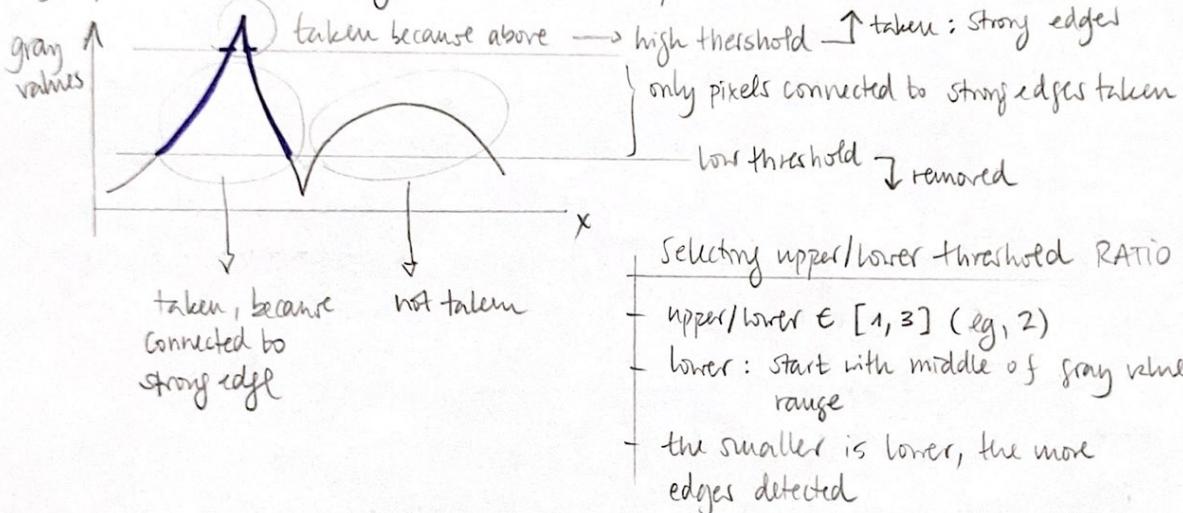
we face up many issues here

- Edge detection process: grayscale \rightarrow lowpass \rightarrow highpass \rightarrow binary threshold

solution

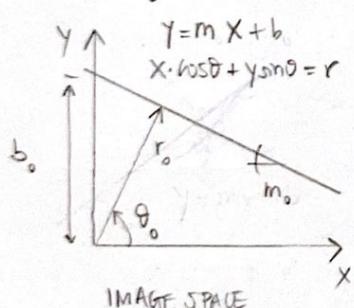
- CANNY EDGE DETECTOR: Steps (4)

- 1) Filter out noise with Gaussian blur
- 2) Sobel X & Y to compute strength & direction of edge
- 3) Non-maximum suppression: strongest edges selected and thinned
- 4) Hysteresis thresholding to isolate best edges

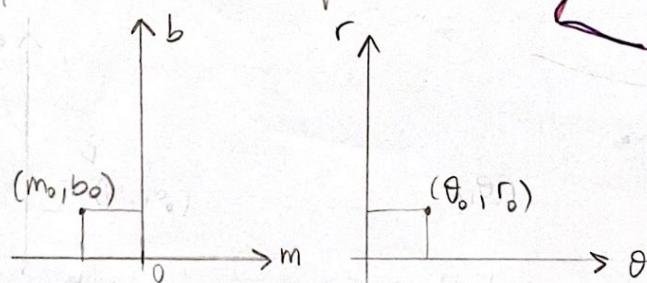


Hough Transform \rightarrow Line Detection as base of Shape Detection

- Most shapes are a collection of line segments
- Hough transform helps detect lines on image



- we represent a line parametrically
- two options:
 (m, b) cartesian
 (r, θ) polar



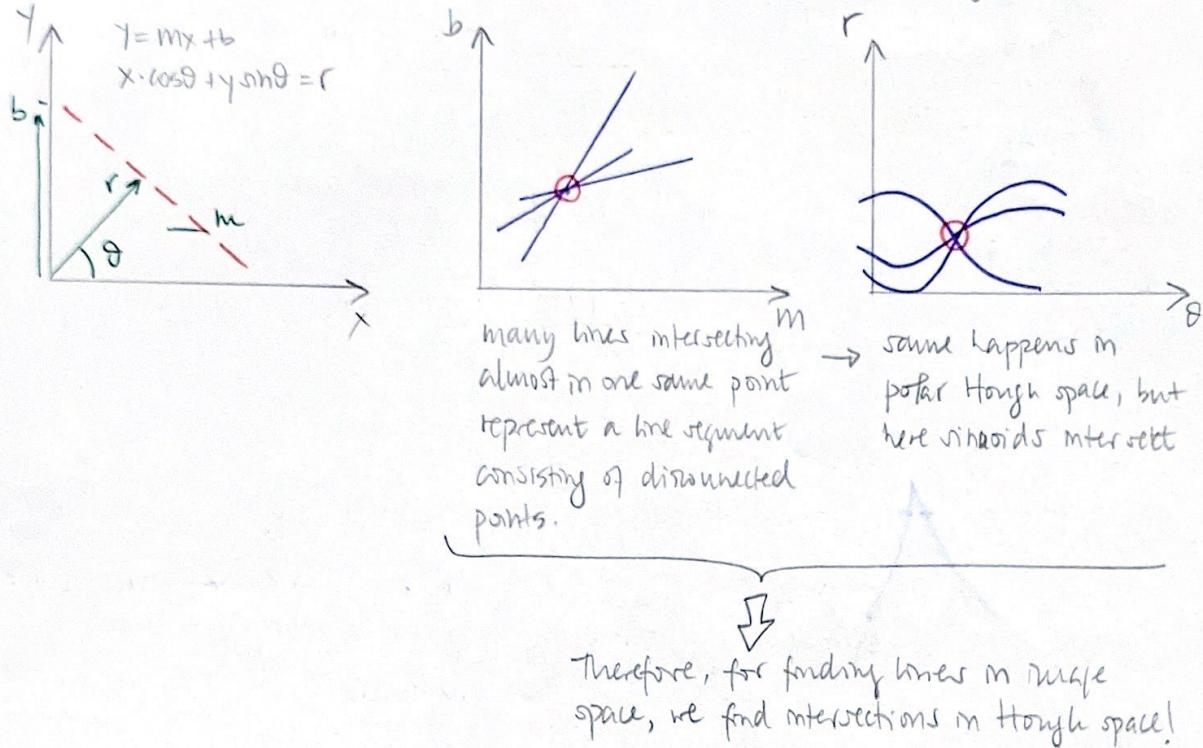
- the parameters of a line are represented
- a line becomes a point!



- a better representation is with polar coordinates, because otherwise slope of vertical lines $m \rightarrow \infty$

original binary edges
 ~detected line segments

- Intersection points in Hough space are discontinuous lines in image space!



- How does the Hough transform work?

- We apply an edge detector to image, eg Canny
- Each point (x_i, y_i) of the binarized image is transformed as (θ_i, r_i)
- The Hough space is discretized in accumulator cells that contain how many (θ_i, r_i) points lie in them
 - cells with resulting high values are related to intersections \Rightarrow lines!

- The same approach can be applied to other analytical shapes, eg. circles:

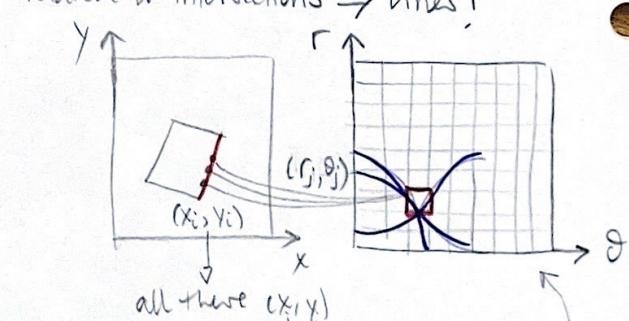
$$(x-a)^2 + (y-b)^2 = R^2$$

↳ 3 params: $a, b, R \Rightarrow$ 3D Hough space with a 3D accumulator

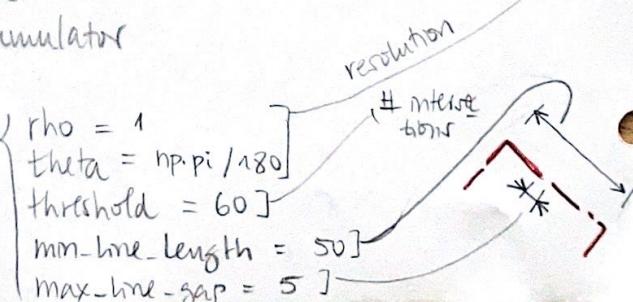
- Parameters for OpenCV:

`cv2.HoughLinesP(cannyImage, ...)`

→ Note: in `cv2.HoughCircles()`, the gray blurred image is passed, not the canny image!

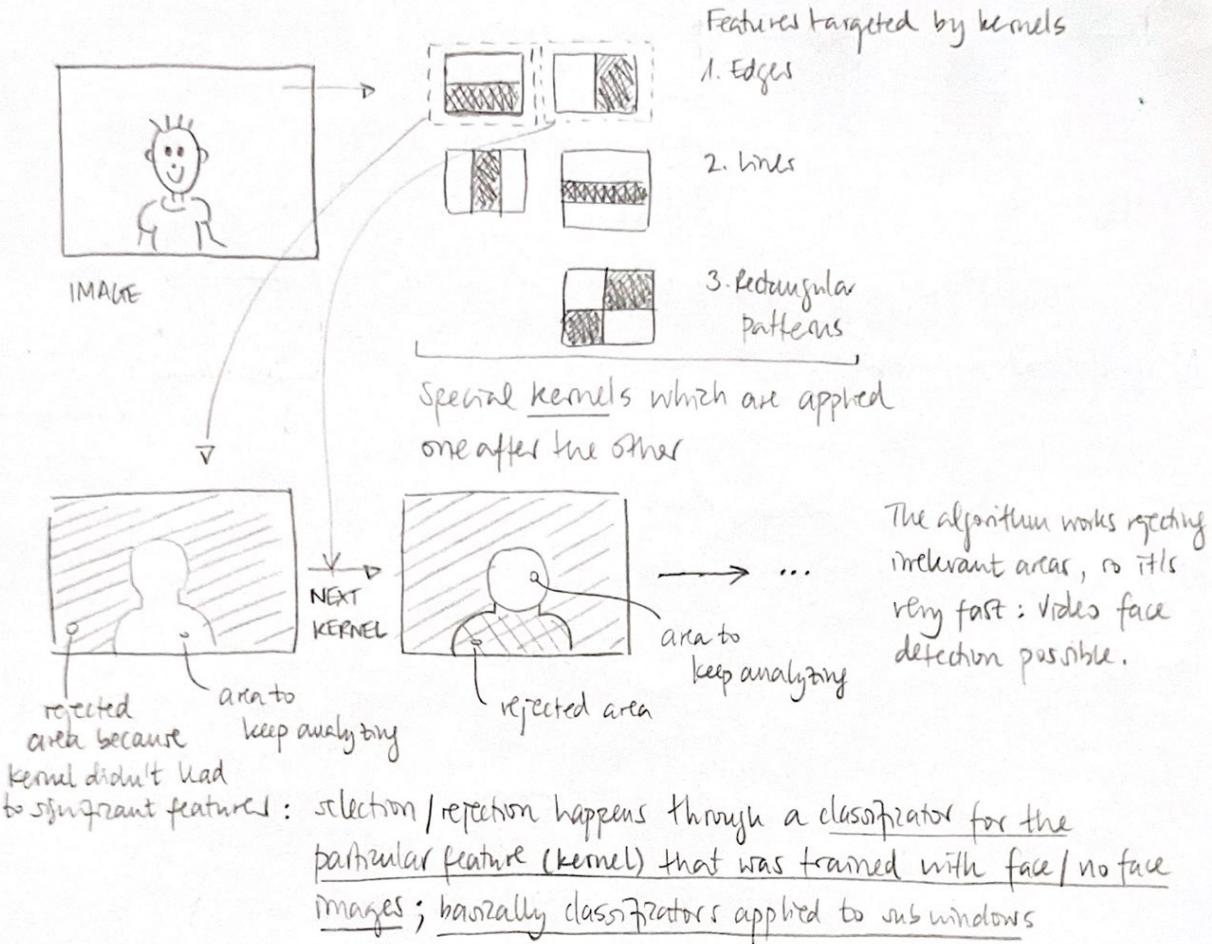


pairs produce a similar (r, θ) !



$\text{rho} = 1$
$\text{theta} = \text{np.pi}/180$
$\text{threshold} = 60$
$\text{min_line_length} = 50$
$\text{max_line_gap} = 5$

Haar Cascades: Applied to Face Detection (LaVilla-Jones Paper/Algorithm)



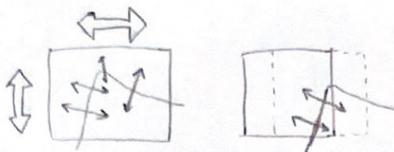
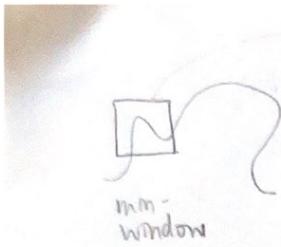
- Haar cascades can be used for segmenting areas of interest for further processing.
- We can define Haar cascade classifiers for different objects, not only faces, but features and kernels must be defined accordingly, and the classifiers trained! (these are trained with yes/no labelled images beforehand)

TYPES OF FEATURES AND IMAGE SEGMENTATION

Corner detection: Harris

- Typical features can be
- Corner detection can work as follows:
 - We take a mini-window and observe the image gradient magnitude and direction
 - Then we shift the window in all directions
 - If the magnitude and direction change considerably → CORNER!

- 1. Edges: high gradients or intensity changes
 - 2. Blobs: regions or areas with a color/intensity / texture
 - 3. Corners: intersections of edges
- this are the most salient!
because they can match unique parts of the image
⇒ one of the most important!



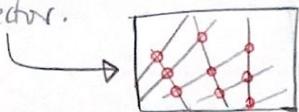
magnitude of the gradient: $\rho = \sqrt{S_x^2 + S_y^2}$
Sobel

we shift window
(we add a small
displacement
vector)

direction of the gradient

$$\varphi = \arctan\left(\frac{S_y}{S_x}\right)$$

- This is implemented in the **Harris** corner detector.



Morphological operations

Often performed on binary images to improve their properties

* Dilation → cv2.dilate (image, kernel, iterations)

* Erosion → cv2.erode (image, ", ")

[None, default] usually 1.

* Opening: erosion + dilation

remove
noise and
shrink region

enlarge
object again

remove salt noise

cv2.morphologyEx(image, cv2.MORPH_OPEN, kernel)

* Closing: dilate + erode

close small
gaps and
enlarge region

shrink
obj again

close holes = remove pepper noise

cv2.morphologyEx(..., cv2.MORPH_CLOSE, ...)

Contours (Motivation: Image Segmentation)

- Contours: complete closed boundaries that contain objects

- Continuous curves that follow the edges along a boundary
- Computed after edges or a binary threshold
- Provide info on contained object: area, shape center, perimeter, bounding box

shape
center
perimeter
bounding box

contour features

there are much more...
see TXT

there can
be used to
analyze
Shape of
objects

`cv2.threshold()` → `cv2.findContours()` → `cv2.drawContours()`

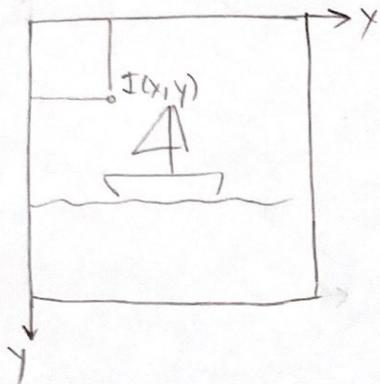
try to use white background, then apply inverse threshold

retrieve contour tree or hierarchy
+ list of contours

each contour is a list of coordinates in the image — I understand that, joining them, we draw the contours!

so we don't have subpixel precision?

Image Moments



$$M_{pq} = \iint_{-d_1-d_2}^{d_1+d_2} x^p y^q I(x, y) dx dy \quad (\text{continuous})$$

$$[M_{ij} = \sum_x \sum_y x^i y^j I(x, y)] \quad (\text{discrete, image})$$

sum of grayvalues: M_{00}

for BINARY images: area = M_{00}

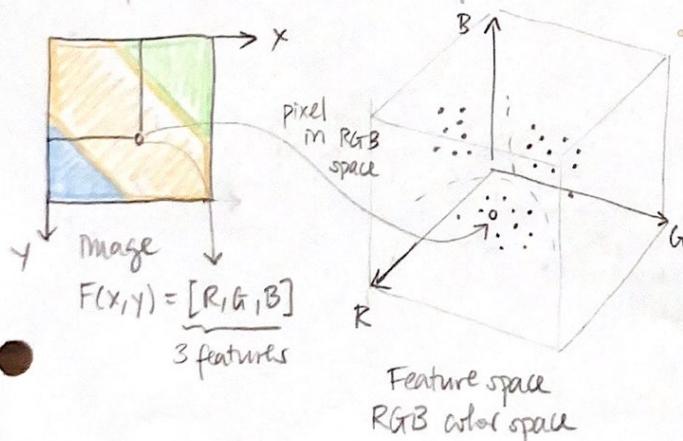
$M = cv2.moments(\text{contour})$

$M_{00} = M['m00'] \rightarrow$ dictionary with all moments

$$\cdot \text{Centroid: } (\bar{x}, \bar{y}) = \left(\frac{M_{10}}{M_{00}}, \frac{M_{01}}{M_{00}} \right)$$

K-means clustering (Motivation: Image Segmentation)

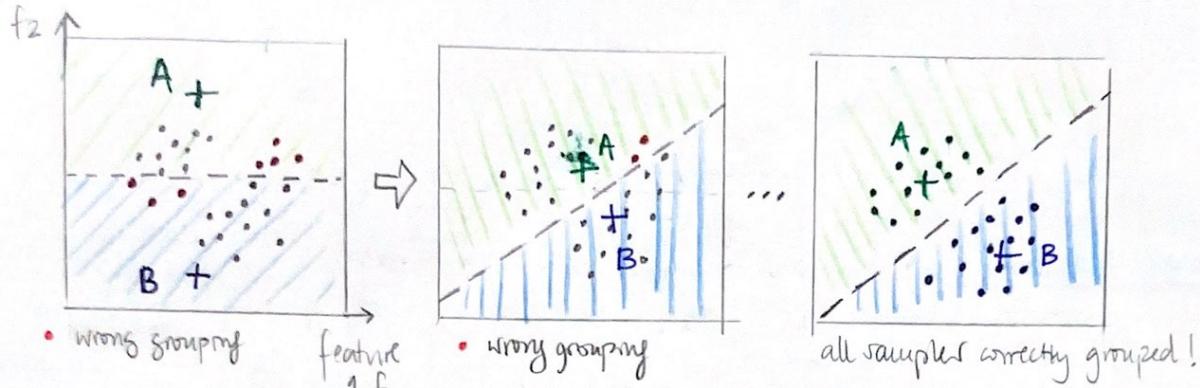
- It separates an image into segments by clustering data points that have similar features/trait.
- Unsupervised learning: no need of labelling
- Input variable: K groups to cluster data



Example:

- Color segmentation of a flag
- $K=3$
- (x, y) pixels set in RGB color space
- $K=3$ clusters identified
- Assumption: there are 3 natural or expected color groups!

- clustering (2D example):



- 1) $K=2$ centerpoints A, B are randomly set in feature space
- 2) Assign each data point to closest centerpoint A/B

- 3) Move centerpoints A & B to center of mass of clustered points
- 4) Repeat 2): Assign each point to closest centerpoint A/B and then repeat 3)

- The algorithm is expected to converge if there's an underlying structure.
- STOP criteria:
 - num iterations
 - if centerpoints move less than a threshold

! The notebook with the example is very interesting because it has non-straightforward reshaping & masking operations that are banz ...

FEATURE VECTORS (4 OBJECT IDENTIFICATION)

→ Feature vectors are used to detect objects, going beyond object parts such as corner

- Corners are distinct image features that can be used as matching patterns
 - BUT: we want to find whole objects, not just corners
- ↳ SOLUTION: we implement feature vectors that contain several (eg corner) feature data

- Example

Gradient direction feature vector

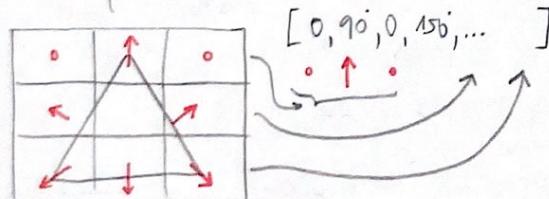
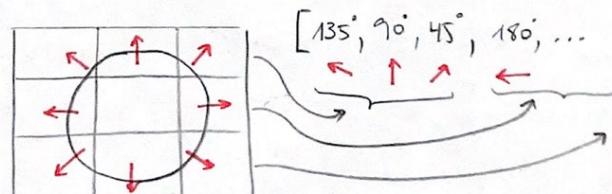


image divided in cells and average gradient direction computed



for applications {
• Self-driving cars
• AR

- ADDITIONALLY: we want to detect objects in REALTIME

Rublee, Rabaud, Konolige, Bradski - 2011, Willow Garage
 ORB = Oriented FAST, and rotated BRIEF

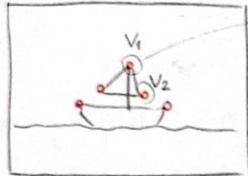
Binary robust independent elementary features

Features from accelerated segment test

→ State-of-the-Art algorithm feature detector (2011):

- used for real-time object recognition or 3D reconstruction
- alternative to SIFT & SURF: scale-invariant feature transform → because there are patented!
- robust wrt illumination changes, noise & image transformations!

1) FAST: feature detection



Keypoints: small distinctive regions, eg, corners

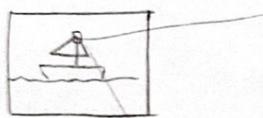
2) BRIEF: feature vector creation

$$V_1 = [0101110\dots] \quad \text{keypoints encoded as bitstrings = Binary feature vectors}$$

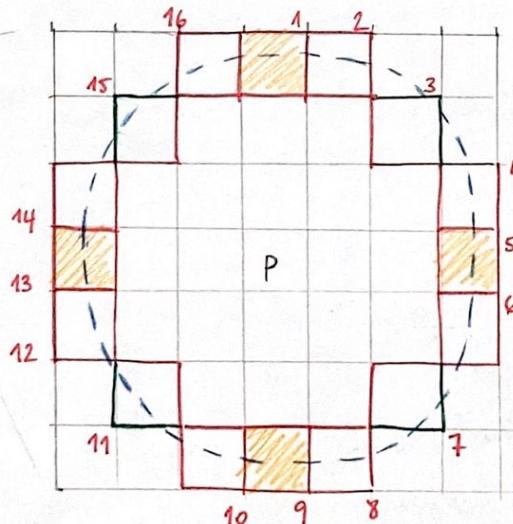
$$V_2 = [10010011\dots]$$

Sequence depicts properties of keypoint area: patterns of intensity

① FAST: Feature or keypoint detection (Features from Accelerated Segment Test)



We stop in a pixel P



- I_p : intensity of pixel p
- h: threshold intensity
- 16 pixels in a circle around p are classified to be:
 - BRIGHTER than p if $I > I_p + h$
 - DARKER if $I < I_p - h$
 - SIMILAR otherwise

→ Pixel P is a keypoint iff > 8 CONNECTED pixels in this set of 16 are either BRIGHTER or DARKER

→ Faster version with same result: consider set of 4 equidistant pixels, (4X faster)
 e.g., 1, 5, 9, 13

→ p is keypoint iff 2 consecutive pixels are either BRIGHTER or DARKER

- keypoints found by FAST account for the location (not orientation) of object

DEFINING EDGES

(Binary Robust Independent Elementary Features)

② BRIEF: Feature vector creation mem. efficient & fast to manipulate

- keypoints are encoded as binary vectors or bit-strings of 128 to 512 bits or descriptors

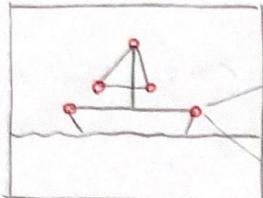
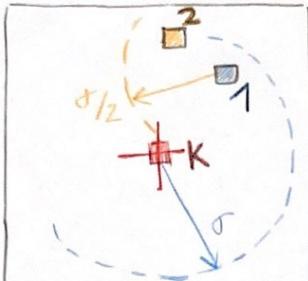


Image with keypoints
+ BLURRED to remove
high freq. noise
(Gaussian kernel)

for each
keypoint K



patch: neighborhood of
w x h pixels around
keypoint K

- Two random pixels are selected around K

- 1) First pixel within the Gaussian spread with σ around K
- 2) Second pixel within the Gaussian spread $\sigma/2$ around pixel 1

- Then, intensities of 2 pixels are compared:

$I_1 > I_2 \rightarrow 1$, & otherwise
(1 brighter than 2)

?

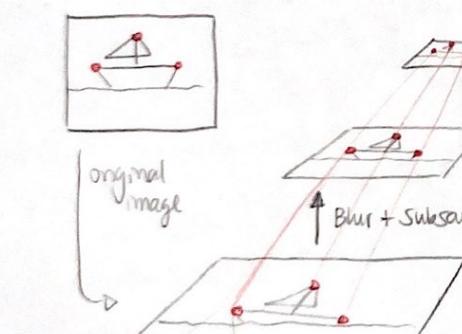
BUT random pixels are always the same or each time different?

- The process is repeated N times for each key-point ($N = \text{length of binary string}$):

$$V(K) = [\underbrace{100110\dots}_{N}]$$

It contains intensity/brightness comparisons of N (128..512) random pixel pairs around keypoint K

SCALE & ROTATION INVARIANCE : how is it achieved?



: (higher levels)

Level 2: $\frac{W}{4} \times \frac{h}{4}$, $\frac{1}{4}$ resolution

Level 1: $\frac{W}{2} \times \frac{h}{2}$, $\frac{1}{2}$ resolution

Level 0: $w \times h$

(low levels)

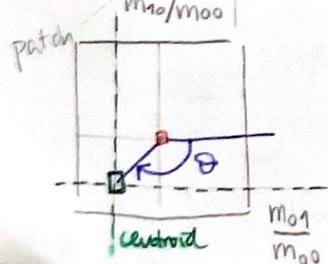
- Image pyramids are computed

- For each level image, keypoints are located

→ this way, we are scale invariant!

?

Basically, keypoints of original level are scaled to higher levels without recomputing them, right?



- For each keypoint of each level, its centroid (center of mass of intensities) is computed in a patch of constant size (number of pixels)

We take angle θ of centroid

constant for all levels, thus, keypoints in higher levels are bigger!

— After computing θ , a modified version of BRIEF is used

↓
— RBRIEF = Rotation aware BRIEF : it introduces rotation invariance

- It starts like BRIEF: $N (= 256, \text{eg})$ random pixel pairs are created

↓
• Then, the set of pixel pairs is rotated θ !

↓
• This way, we always compute the bit-string with respect to the axis given by the centroid!

— The set of keypoints and feature vectors is known as ORB descriptor; all

Feature matching

?) Why is it called training if there's no iterative optimization?

BRIEFS stacked together.

corr.

— Keypoints of training image are saved

— Keypoints of query image are computed (that's what I understand)

— Training & query keypoints are compared: their similarity is measured with a metric

- Hamming distance is commonly used = # positions at which symbols are different

— Best keypoint matches are returned) → (?) Is the relative position of the matched keypoints analyzed?

— If # matches above/threshold ↳ object found!

↳ usually a % of all keypoints; it's a user parameter (e.g. 35%)

Applications

• Video (because it's very fast)

• Face detection

• Concrete object detection, no matter background

— BUT: ORB doesn't work well for general object detection, in other words, to use ORB, objects need to

— In the Notebook (jupyter), we see that ORB is:

- scale invariant
- rotation invariant
- illumination invariant
- noise invariant

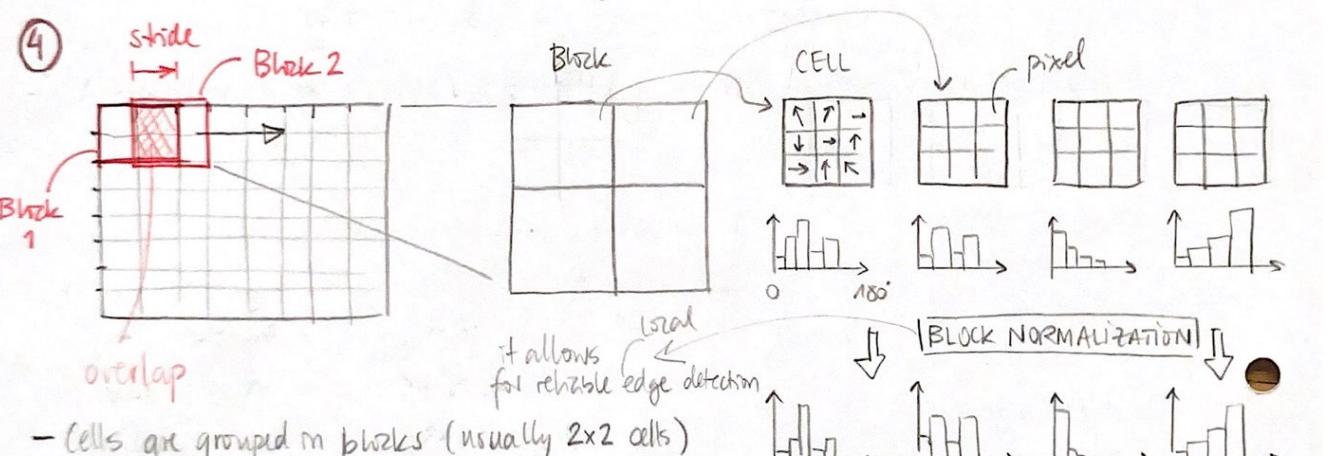
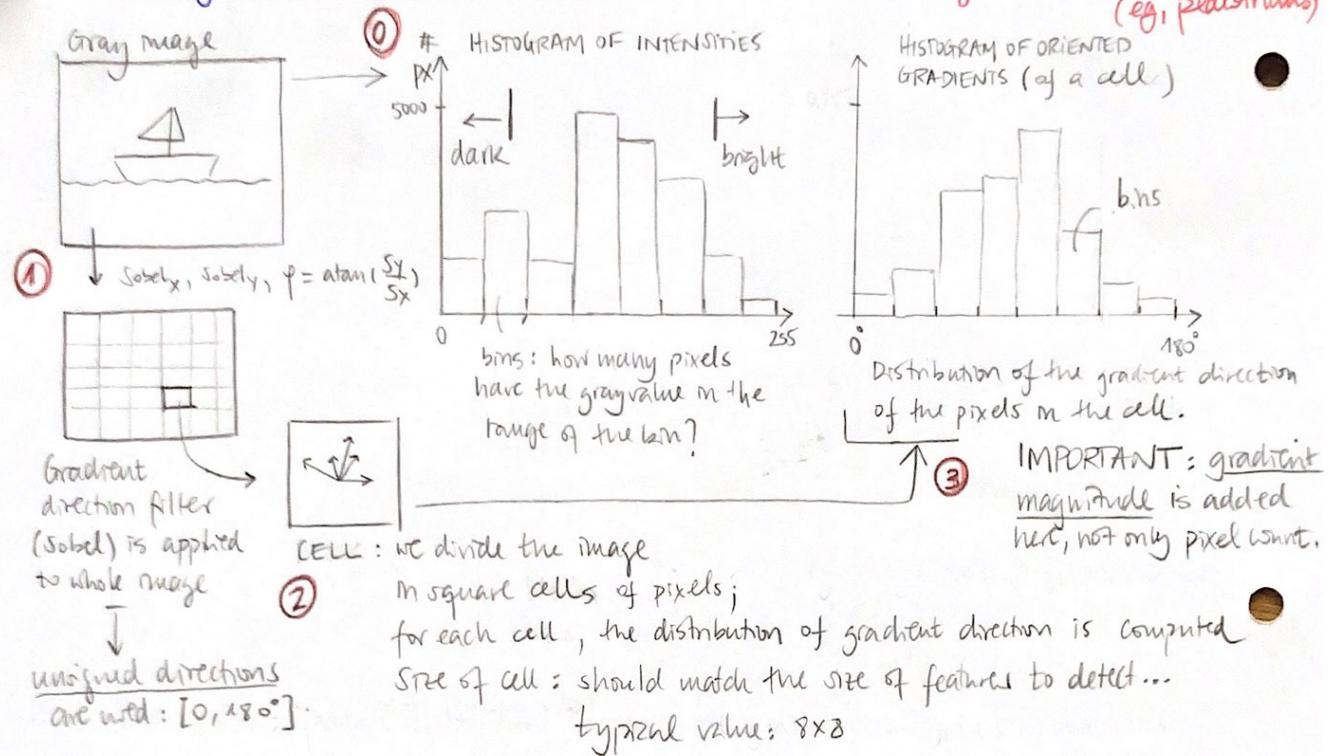
?) If BRIEF is done with random pixel positions each time, how is that going to work?

have clear distinctive and consistent internal features

for other cars, we can use HOG with a classifier.

?) The real useful application would be to determine the scaling & pose of the template in the image — is there a function for that?

HOG = Histogram of Oriented Gradients (Dalal & Triggs 2005) → Object detection of generic object classes (e.g., pedestrians)



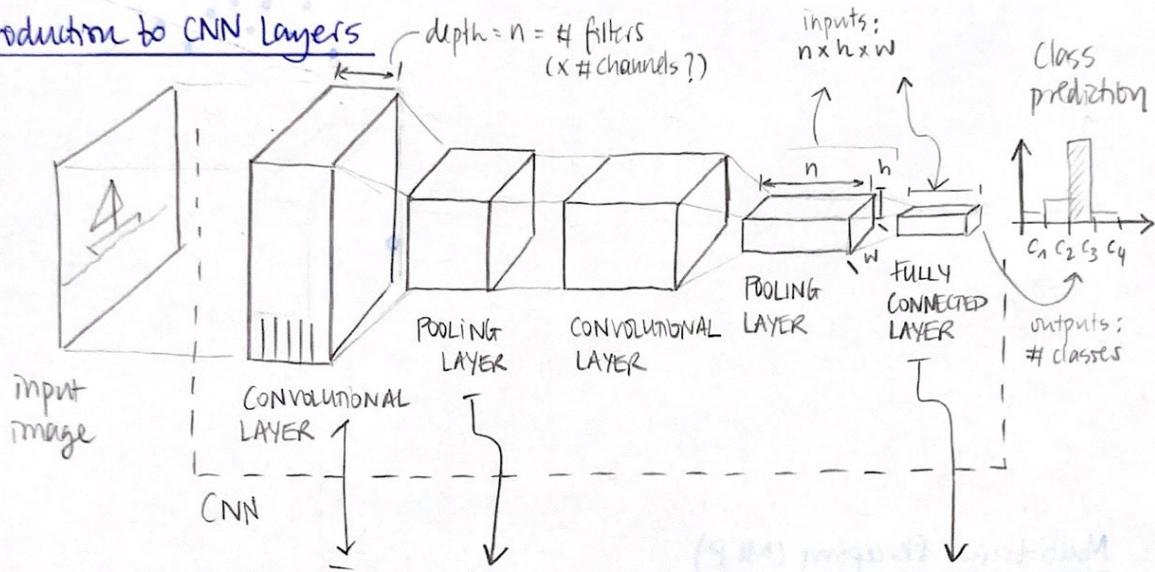
- Cells are grouped in blocks (usually 2x2 cells) that usually overlap
- The distance between blocks is called **stride**; it's recommended having a stride that generates overlap (typically: 1 cell)
- The HOGs of a block are normalized together
- Thus, assuming blocks overlap, most cells will be normalized wrt several blocks → this increases accuracy
- **HOG descriptor**: feature vector that contains all normalized HOG histograms of all the cells in the image → block-cells
- ↓ A classifier, e.g. a Support Vector Machine (SVM) is trained with many HOG descriptors of the object class, and many HOG descriptors of other object classes e.g. pedestrian
- **Detection** is done with a sliding window over image trying to look for image part similar to HOG descriptor

- Note: the # of histograms in the Hog descriptor \neq # cells, but
 $= \# \text{ blocks} \times \text{cells/block}$

Therefore, if there's block overlap (as recommended), the histogram of a cell will come up in the vector several times with a different normalization.

CNN LAYERS AND FEATURE VISUALIZATION → continues on page 28

Introduction to CNN Layers



In a **CONVOLUTIONAL LAYER** we have filters that are applied to the image

- filter weights are learned during training with backpropagation
- filters highlight important features (edges, color, etc.)
- If we have n filters, we'll produce n convoluted maps and the convolutional layer will have depth = n .

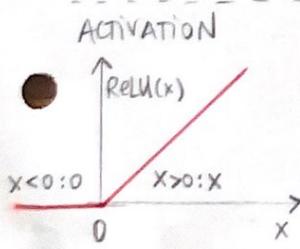
In a **POOLING LAYER** the size of the filtered images is reduced by taking relevant pixels. For instance, if we do maxpooling, the max pixel value of a cell of pixels is taken; example

20	30	80	90
140	45	40	50
200	100	6	7
15	16	50	250

⇒

140	90
200	250

In a **FULLY CONNECTED** layer we take as input all the pixels from last layer and map them to class nodes in a fully connected manner



- Images are normalized to $[0, 1]$ region
- Negative pixel values are set to 0 with ReLU activation

This:

- reduces dimensionality
- increases robustness wrt small grayvalue changes
- we increase field of view to later filters

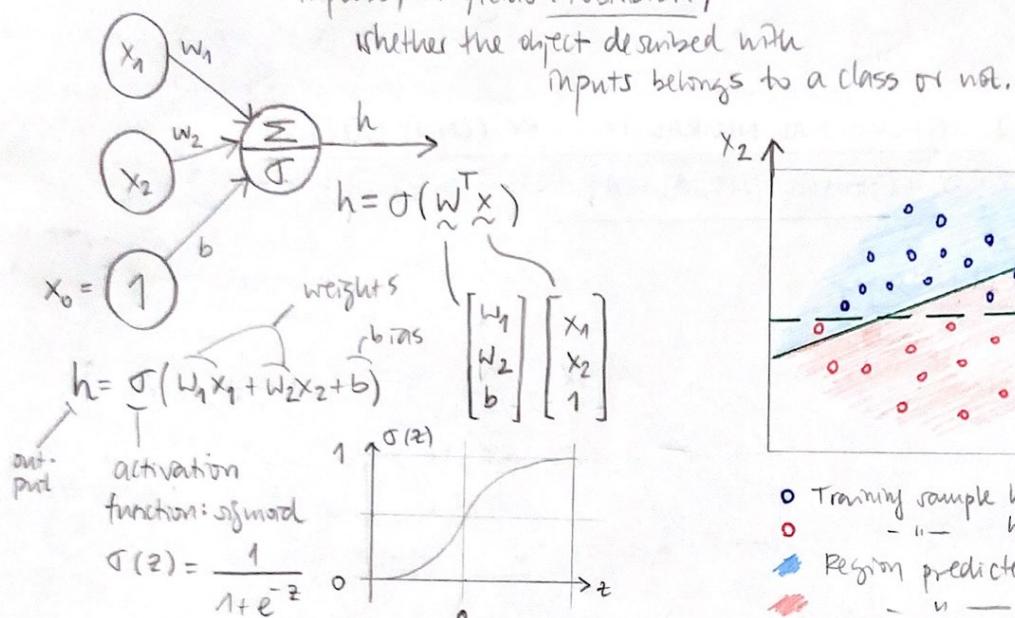
NEURAL NETWORKS

Optional excursions (from several sources)

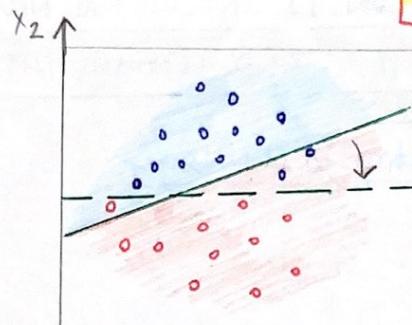
but, in general, following Udacity's guideline

Perception: binary classifier: given n

inputs, it yields PROBABILITY



Neural networks with non-linear activations work like universal function approximators.



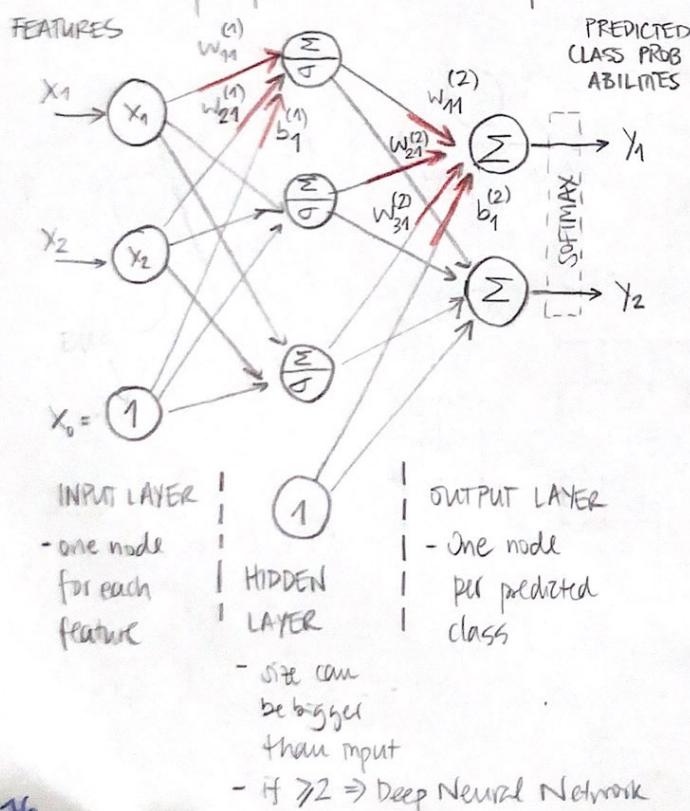
- Training sample $h=1$ (ground truth)
- " " $h=0$ (- " -)
- Region predicted as $h=1$
- " " $h=0$

Boundary trained, composed of the weights

During training the weights are modified so that the boundary better splits the training samples in their categories

Multi-Layer Perceptron (MLP)

Several perceptrons stacked one after the other



- All layers are fully connected
- All layers except output have bias b ($= x_0$)
- Several activation functions are possible:
 - Heaviside step: $\sigma_h(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$

• ReLU: $\sigma_r(z) = \max(0, z)$

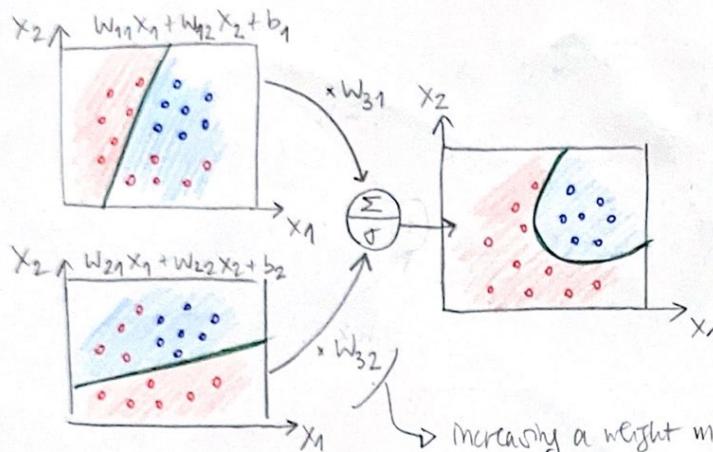
It's non-linear, but the most effective & used ...

- For classification, softmax is applied with inputs: normalized exponential function

$$\sigma_i(\vec{z}) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}, \quad i = 1, \dots, K, \quad \vec{z} = [z_1, \dots, z_K]$$

classes

- While a single perceptron draws a linear boundary between 2 regions, a MLP is a linear combination of perceptrons that achieves non-linear boundaries



The weights and biases define non-linear boundaries in feature space!

Each region is a class!

- the # features defines the dimension
- the # classes defines the regions

Increasing a weight means giving more importance to the output of that node, as it happens with neurons in the brain: as Hebb suggested: when a neuron triggers another, their connection grows stronger.

Feedforward: the process of obtaining a prediction after applying the network to the input

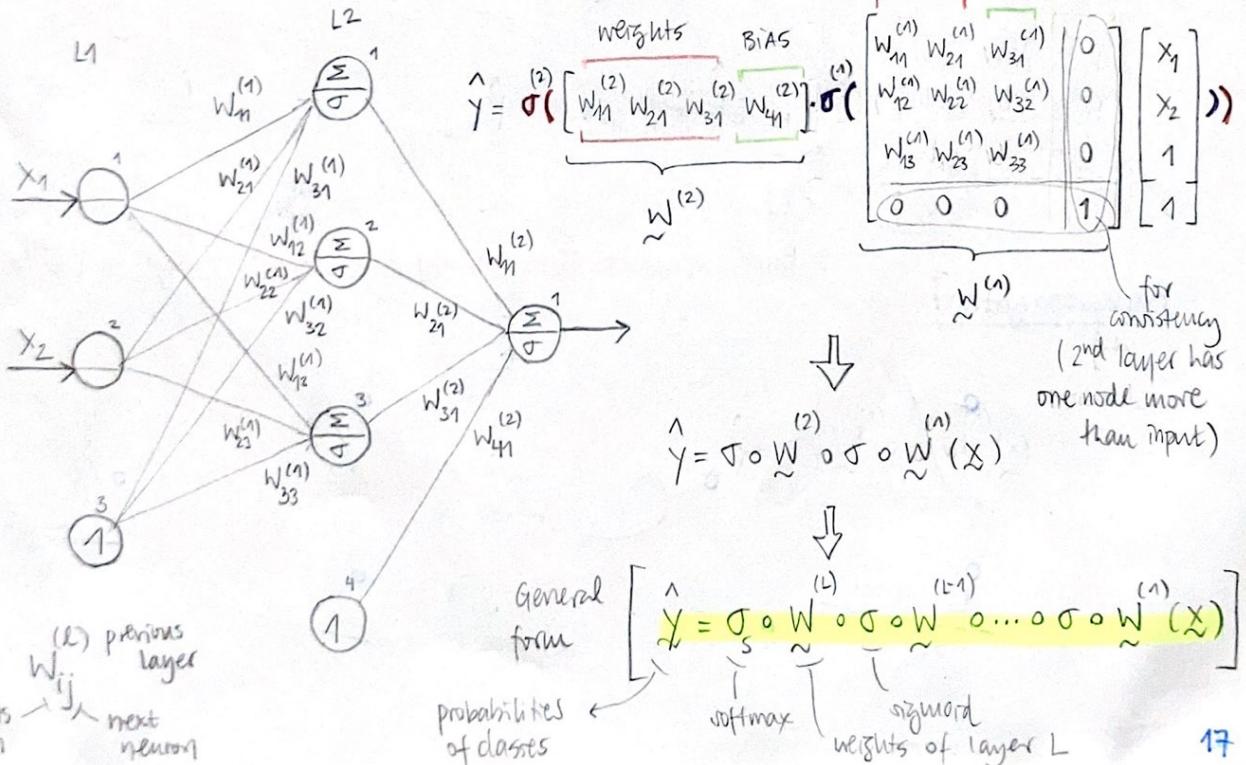
Notation: Function composition: $g \circ f(x) \equiv g(f(x))$

$$\begin{array}{ccc} \xrightarrow{x} & \boxed{\text{MLP}} & \xrightarrow{\hat{y}} \\ \text{feedforward} & & \\ x \rightarrow \hat{y} & & \end{array}$$

$\hat{x} = [x_1 \ x_2 \ x_3]^T$ $\hat{y} = [y_1 \ y_2]^T$

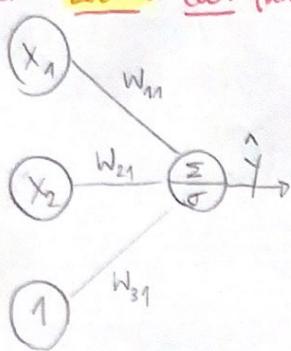
Inputs: features (eg 3) # classes: probabilities (eg, 2 classes)

more neurons & layers result in more regions & better performance



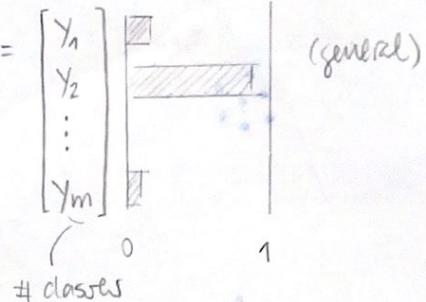
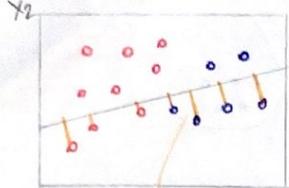
Error function: A measures the error between the ground truth and the prediction

aka Loss or Cost function



$$y = \sigma(\tilde{w} \cdot \tilde{x}) \text{ (Perception)}$$

$$\hat{y} = \sigma \circ w^{(1)} \circ \dots \circ w^{(n)}(x) = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$



classes

ground truth
predicted

The error is a metric that measures the distance between \hat{y} and y

- One metric would be the Mean Squared Error (MSE)

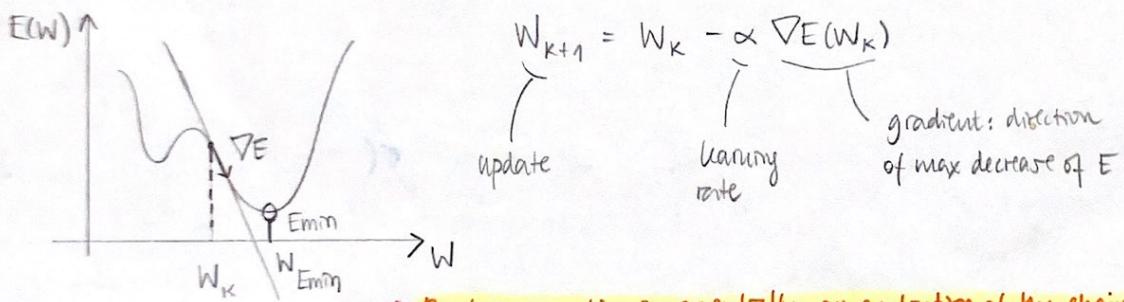
$$E_{\text{MSE}}(\tilde{w}) = \frac{1}{m} \sum_{j=1}^m (y_j - \hat{y}_j)^2 \quad (\text{that's actually related to the distance between } y \text{ and } \hat{y})$$

- Since \hat{y} measures probabilities of classes, it's better to use the Cross-Entropy:

$$E_{\text{CE}}(\tilde{w}) = -\frac{1}{m} \sum_{j=1}^m y_j \ln(\hat{y}_j) + (1-y_j) \ln(1-\hat{y}_j)$$

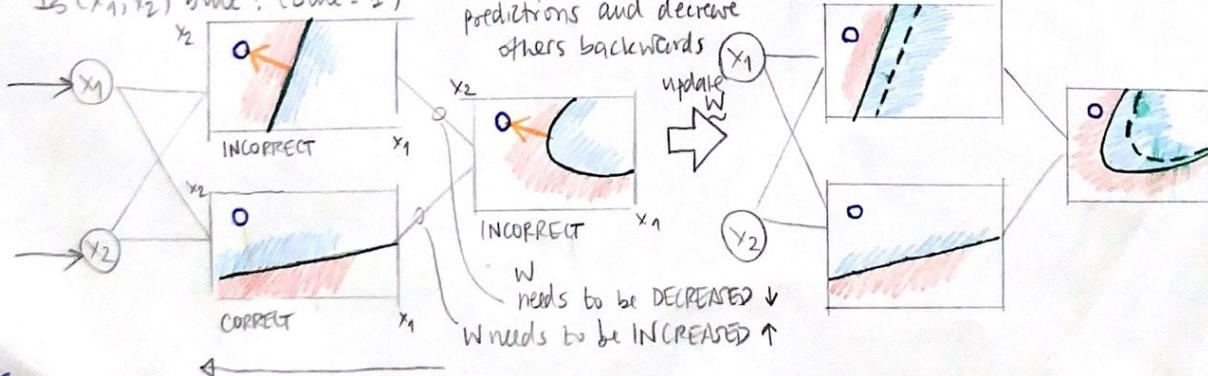
| There are more error functions

Gradient Descent: we modify \tilde{w} in the direction in which $E(w)$ decreases

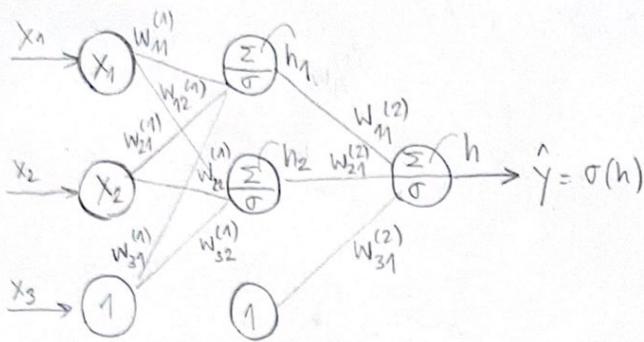


Backpropagation is essentially an application of the chain rule to find

Backpropagation: Intuition: we start in the output and check whether the prediction corresponds with ground truth; we increase the weights that contribute correct predictions and decrease others backwards



The gradient for gradient descent



Chain rule



$$y = f(x) \quad z = g(y) = g(f(x)) \\ = g \circ f(x)$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

- $\hat{y} = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$

↳ usually they are vectorized

- $\nabla E(\tilde{W}) = [\dots, \frac{\partial E}{\partial W_{ij}^{(l)}}, \dots]$ The gradient of E is also vectorized

- Weights are updated like this: $W_{ij}^{(l)} = W_{ij}^{(l)} - \alpha \frac{\partial E}{\partial W_{ij}^{(l)}}$ (Gradient descend)

$\alpha \frac{\partial E}{\partial W_{ij}^{(l)}}$

which are these?

- Computation of $\frac{\partial E}{\partial W_{ij}^{(l)}}$ (with the example):

$$h_1 = W_{11}^{(1)} x_1 + W_{21}^{(1)} x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{32}^{(1)}$$

$$h = W_m^{(2)} \sigma(h_1) + W_{21}^{(2)} \sigma(h_2) + W_{31}^{(2)}$$

$$\hat{y} = \sigma(h) = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

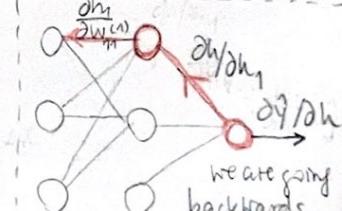
$$E(\tilde{W}) = -\frac{1}{m} \sum_{k=1}^m y_k \ln(\hat{y}_k) + (1-y_k) \ln(1-\hat{y}_k) = E(W_m^{(1)}, W_{21}^{(1)}, W_{31}^{(1)}, W_{12}^{(1)}, W_{22}^{(1)}, W_{32}^{(1)}, \dots)$$

$$\nabla E(\tilde{W}) = \left[\frac{\partial E}{\partial W_m^{(1)}}, \frac{\partial E}{\partial W_{21}^{(1)}}, \frac{\partial E}{\partial W_{31}^{(1)}}, \dots \right]$$

$$\frac{\partial E}{\partial W_m^{(1)}} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} \cdot \frac{\partial h}{\partial h_1} \cdot \frac{\partial h_1}{\partial W_m^{(1)}}$$

all terms are here
all terms here:

BACKPROPAGATION VISUALLY



we are going backwards...

$$\frac{\partial E}{\partial \hat{y}} = \hat{y} - y \quad , \quad \frac{\partial h}{\partial h_1} = W_m^{(2)} \frac{\partial \sigma}{\partial h_1} = W_m^{(2)} \underbrace{\sigma(h_1)(1-\sigma(h_1))}_{\sigma'(h_1)}$$

$$\frac{\partial \hat{y}}{\partial h} = \sigma(h)(1-\sigma(h)) = \hat{y}(1-\hat{y}) \quad , \quad \frac{\partial h_1}{\partial W_m^{(1)}} = x_1$$

$\frac{\partial}{\partial h_1} \sigma'$
sigmoid derivative!

So, basically, each $\frac{\partial E}{\partial w_{ij}^{(l)}}$ is computed following the chain rule and the vector $\nabla E(\tilde{w})$ is assembled

Training happens by updating the weights with the gradient:

$$\tilde{w}' = \tilde{w} - \alpha \nabla E(\tilde{w})$$

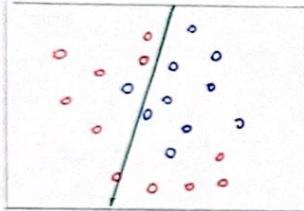
In short, backpropagation consists in applying gradient descent to the $E(\tilde{w})$ function, and for that, using autodifferentiation and the chain rule.

TRAINING: Test split & complexity

- In the dataset, ALWAYS define a training split and a test split.
- Use the test split to evaluate how the model works.
- Prefer simple models over complex:

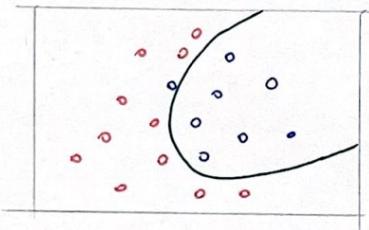


TRAINING: Overfitting & Underfitting : we check our fitting with the test split!!!



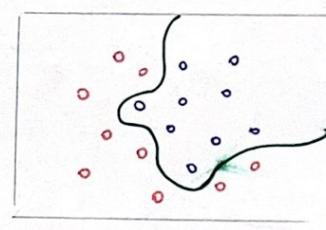
UNDERFITTING

- Too simple model
- Error by bias because we have a high bias
- Usually simplistic architectures



CORRECT

- Incorrect prediction is treated as noise
- Model generalizes correctly
- Slightly complex architectures



OVERTFITTING

- Too complex model
- Error by variance, because high variance is learned
- Model is TOO SPECIFIC and has problems to generalize
- Noise is fitted
- Very (too) complex architectures

Methods to prevent overfitting:

1. Early stop
2. Regularization
3. Dropout
4. Batch Normalization

usually we start with a complex architecture and prevent it to overfit.

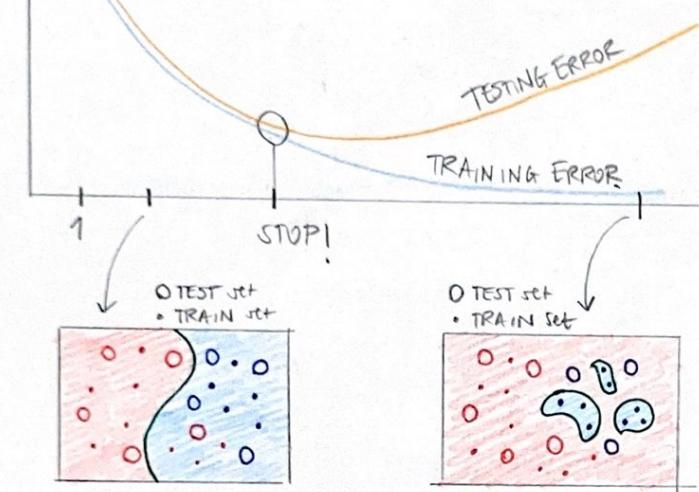
TRAINING: Early Stopping (It prevents overfitting)

— Divide the dataset in randomly distributed TRAINING & TESTING sets

— Train with the TRAINING set + compute E_{train}

— Test " " TESTING set + compute E_{test}

Error
UNDERFITTING OVERFITTING

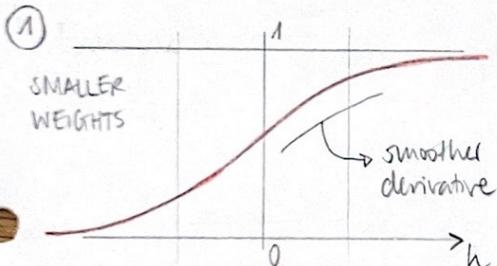


↑
not each by network
for computing
weights

→ Stop test when
 E_{test} starts to
increase while E_{train}
decreases

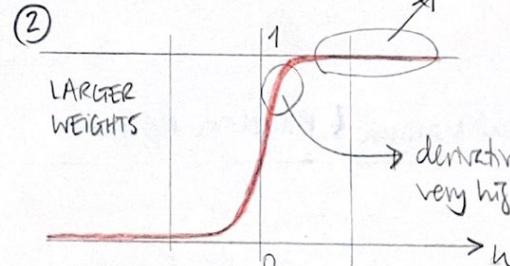
Epochs: each epoch the
whole training samples
have been processed

TRAINING: Regularization (It prevents overfitting)



$$\hat{y} = \sigma(\tilde{W}^T \tilde{x})$$

$$\hat{y} = \sigma(x_1 + x_2), W = [1 \ 1 \ 0]^T$$



$$\hat{y} = \sigma(k \tilde{W}^T \tilde{x}), k > 1$$

$$\hat{y} = \sigma(10x_1 + 10x_2), W = [1 \ 1 \ 0]^T, k=10$$

— Larger weights (②) shrink the sigmoid function and it becomes more sensitive and unstable

- Derivative is almost \neq very soon
- and it suddenly increases a lot close to the origin

— Too large weights stretch the output of σ to the border values (1 or 0), so the model is TOO CERTAIN = WE ARE OVERFITTING!

→ **SOLUTION: REGULARIZATION**: Penalize large weights by adding them to the error

→ 2 major methods:

$$\begin{cases} L1 : \sum |W_{ij}| \\ L2 : \sum |W_{ij}|^2 \end{cases}$$

L1 REGULARIZATION

$$E(\tilde{W}) = -\frac{1}{m} \sum_{j=1}^m y_j \ln(\hat{y}_j) + (1-y_j) \ln(1-\hat{y}_j) + (\lambda) \sum |w_{ij}|$$

L1 Regularization term

$\lambda \uparrow$ if we want to penalize more

- We tend to end up with sparse vectors:
 - small weights become 0
 - bigger weights become 1
- Good for FEATURE SELECTION

L2 REGULARIZATION

$$E(\tilde{W}) = -\frac{1}{m} \sum_{j=1}^m y_j \ln(\hat{y}_j) + (1-y_j) \ln(1-\hat{y}_j) + (\lambda) \cdot \sum (w_{ij})^2$$

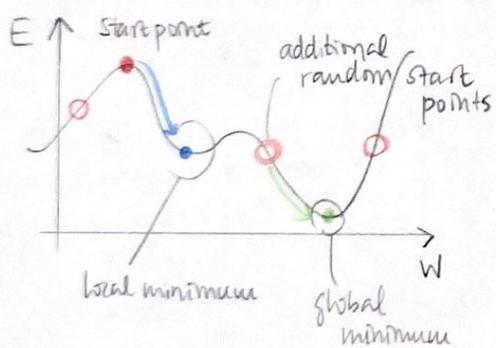
L2 regularization term

- We tend to have homogeneously small weights in the vectors
- It's the most used, because usually we want to train models with homogeneously SMALL WEIGHTS

TRAINING: Dropout (It prevents overfitting)

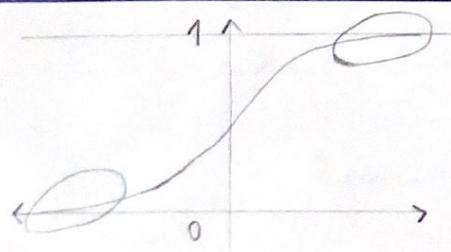
- sometimes some weights are considerably larger than others and they end up dominating the training
- ↓ Solution: we randomly deactivate nodes during training & let the rest train = DROPOUT.
- The probability of any node to deactivate is fed as parameter, e.g. 0.2.

TRAINING: Local Minima & Random Restart



- Gradient descent has the flaw or risk of leading to a local minimum.
- SOLUTION: Start from several random initial sets = RANDOM RESTART
 - We cannot guarantee we achieve the global min
 - BUT: the local min will be better than before!

TRAINING: Other activation functions



- The problem of the vanishing gradient: the sigmoid function has a very small slope very soon
- Since the gradient is packed of those derivatives, we're going to have VERY SMALL steps!

$$\nabla E(\tilde{w}) = [\dots, \frac{\partial E}{\partial w_{ij}^{(l)}}, \dots]$$

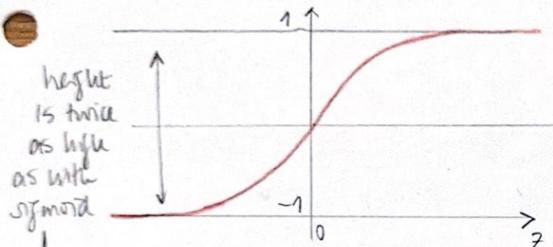
$$\frac{\partial E}{\partial w_{ij}^{(l)}} = \underbrace{\frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} \cdot \frac{\partial h}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ij}^{(l)}}}$$

all are derivatives of a sigmoid \Rightarrow small

product of small values yields a TINY value!

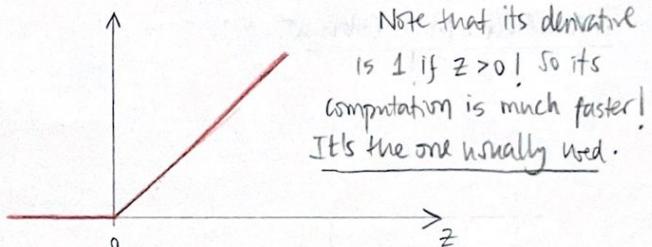
\rightarrow The best way to solve this issue are other activation functions with larger slopes:

HYPERBOLIC TANGENT



$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

RELU: RECTIFIED LINEAR UNIT



$$\text{ReLU}(z) = \max(0, z) = \begin{cases} 0, & z < 0 \\ z, & z \geq 0 \end{cases}$$

\Rightarrow CONCLUSION:

- Usually the $\text{ReLU}(z)$ function is used as activation function, not the sigmoid $\sigma(z)$
- The last (output) nodes have a softmax or in the case of a unique class:
 - the sigmoid, if we want probability
 - ReLU if we want some regression

TRAINING: Gradient descent vs STOCHASTIC vs BATCH gradient descents

? That's basically the same? YES, if it's meant the same!

\downarrow

- In regular gradient descent we enter all the samples to the network, compute the error & the gradient and make an update step in the weights. That's called an Epoch.

BUT: this approach has some drawbacks

- We need to process all the data to make a small step
- Loading the whole data every step to memory might be impossible in practice!

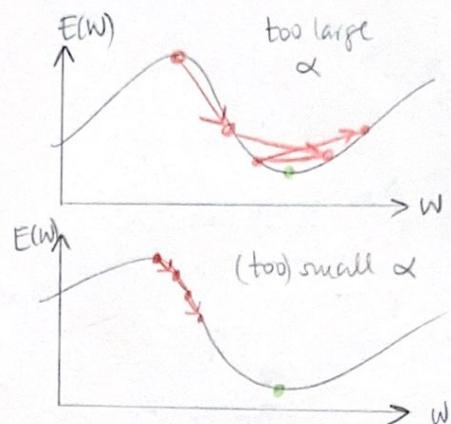
\downarrow

- SOLUTION: STOCHASTIC GRADIENT DESCENT: data is processed in batches or subsets
 - Each batch approximates the error function, but it's more memory efficient to do that
 - A mini-step after each batch is taken, which is better than one step after one Epoch

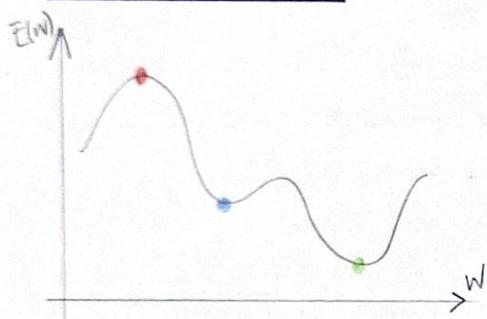
- The mini-steps with batches are less accurate, but more effective at the end

TRAINING: Learning Rate (α)

- Too large learning rate:
 - big steps ; good at beginning,
 - BUT we might miss the minimum
- Too small learning rate:
 - we make too small steps \Rightarrow too slow
 - BUT : usually, tend to use small learning rates
- GOOD strategy: DECREASE learning rate as we approach minimum



TRAINING: Momentum



- In order to avoid falling in local minima, the step is extended with the previous steps:
 - Local minima will be left behind
 - Target minima will be found
 - Since the step size is bigger, convergence around target minimum will take a bit longer

$$\left[\begin{array}{l} z^{k+1} = \beta z^k + \nabla E(w^k) \\ w^{k+1} = w^k - \alpha z^{k+1} \end{array} \right] \rightarrow w^{k+1} = w^k - \alpha (\underbrace{\nabla E(w^k)}_1 + \underbrace{\beta \nabla E(w^{k-1})}_{\beta} + \underbrace{\beta^2 \nabla E(w^{k-2})}_{\beta^2} + \dots)$$

\dots

$\beta \in [0, 1]$

no momentum 0.99 boosts search

that's how it's done in practice:
another step variable is accumulated

that's the effect: previous steps are summed weighting them with β^j , $j = \#$ step back in past;
older steps have less influence

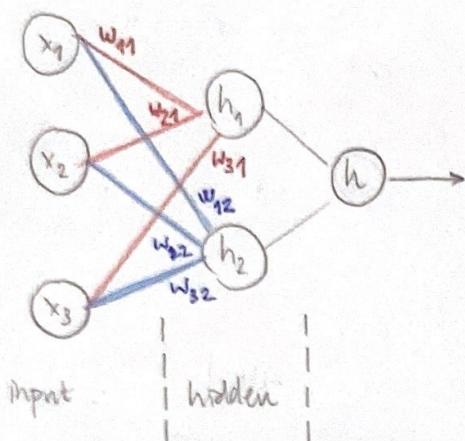
PYTORCH

see notes on Udacity slide related to the Notebooks!!!
I have very simple examples in these handwritten

Notes on how to use it.

hidden units

notes → ADVANCED STUFF
IN JUPYTER!!!



the feed forward propagation between 2 layers can be expressed with matrices

$$\vec{h} = [h_1 \ h_2] = [x_1 \ x_2 \dots x_n] \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ \vdots & \vdots \\ w_{n1} & w_{n2} \end{bmatrix}$$

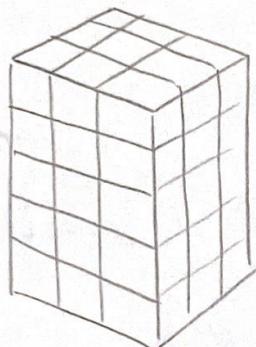
$$\hat{y} = f_2(f_1(\vec{x} \vec{W}_1) \vec{W}_2)$$

Tensors: similar to N-D arrays in numpy

1	1	11
2	2	12
3	3	13
4	4	14
5	5	15

[5]
column vector

[5, 3]
rows x columns



[5, 3, 3]
rows columns depth

- Datasets represented with tensors
- Samples are also tensors, which are passed usually in batches, e.g.:

(64, 1, 28, 28)
batch size channels rows x columns

Networks (import nn, F, etc.)

class Network(nn.Module):

def __init__(self):

super().__init__()

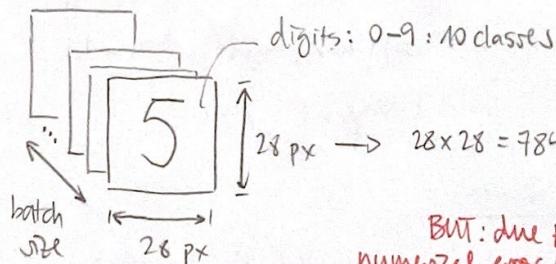
self.fc1 = nn.Linear(784, 128)
self.fc2 = nn.Linear(128, 64)
self.fc3 = nn.Linear(64, 10)

full connection
is like linear
transformation!

def forward(self, x):

x = self.fc1(x)
x = F.relu(x)
x = self.fc2(x)
x = F.relu(x)
x = self.fc3(x)
x = F.softmax(x, dim=1)
return x

data batch
x swept through
all layers



ReLU/Log-softmax is also often used due to num. error with softmax! BUT: due to numerical error, to avoid maxima with softmax, ReLU is also often used for training. This is called logits and it's softmaxed in mfunc25

fully connected
ReLU ReLU

Training & Simple Inference

```

model = Network()
criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(model.parameters(), lr=0.01)
epochs = 3
step = 0

for e in range(epochs):
    running_loss = 0
    for images, labels in iter(trainloader):
        step += 1
        prepare batch
        clear grads
        1. forward
        2. compute loss
        3. compute grad
        4. grad descent
        output
        | images, size = (images.size[0], images.size[2]*images.size[3])
        | optimizer.zero_grad()
        | output = model.forward(images)
        | loss = criterion(output, labels)
        | loss.backward()
        | optimizer.step()
        | running_loss += loss.item()
        | if (step % 40 == 0):
            print("Epoch {} | Loss {}".format(e, running_loss / 40))
            running_loss = 0
    validation() comes here!!!
    images, labels = next(iter(testloader))
    img = images[0].unsqueeze(0, 1, images.size[2]*images.size[0][3])
    with torch.no_grad():
        logits = model.forward(img)
        ps = F.softmax(logits, dim=1)
    → plot(img, ps) ...
  
```

we're passing 1 image!

Validation

- Model is evaluated during training to check that its accuracy is increasing
 \downarrow (\Rightarrow no overfitting) while loss decreases.
- Implemented in a `validation()` function that uses the test split / subset.

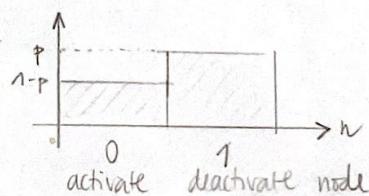
Further Topics in Jupyter Notebooks & Utility Guide TXT.

- General networks with arbitrary num of hidden layers created
- Use log probabilities (log-softmax) as output activation/transformation:
 - larger numbers \Rightarrow more stable
 - multiplication of probs becomes addition of logits \Rightarrow more efficient
- Dropout is added in forward(); it's activated with model.train()
deactivated with model.eval() \rightarrow in evaluation()
 \hookrightarrow used only during training!
- Save and load networks (state & architecture params): it's saved as a dict!
 - weights & biases of each layer
 - {
 - input size
 - output size
 - hidden layer sizes
- Loading any dataset (from folders of images to pre-processed tensor datasets)
 - Data augmentation: introduce variation to generalize the network
 - Normalization: scale gray values to $[-1, 1]$
- Transfer learning: use pre-trained big networks to leverage DL applications
 - Last layers (classifier) implemented and trained
 - GPU device: use (CUDA) if available: models, images & labels (or any tensor) must be transferred to device: `model.to(device)`

QUESTIONS

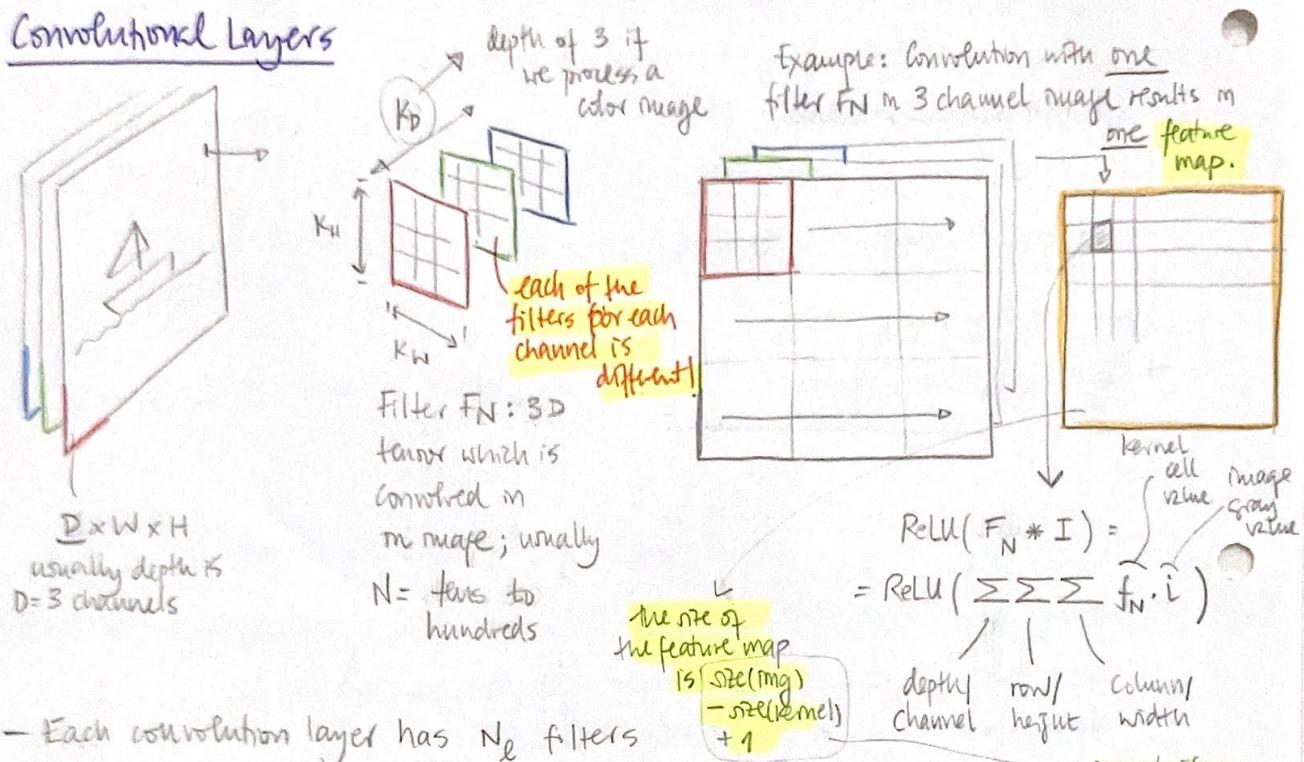
- Sometimes the code produces errors on my computer; for example:
 - transforms.RandomRotation()
 - Is the order of the transforms relevant?
 - Dropout: [Why is it done after `relu(forward(x))`?
Does here `.dropout()` mean it's just activated for next passes?]
- look at Jupyter Notebooks & fc-model.py & helper.py!

→ No; the `dropout()` function randomly zeros some of the elements following a Bernoulli distribution



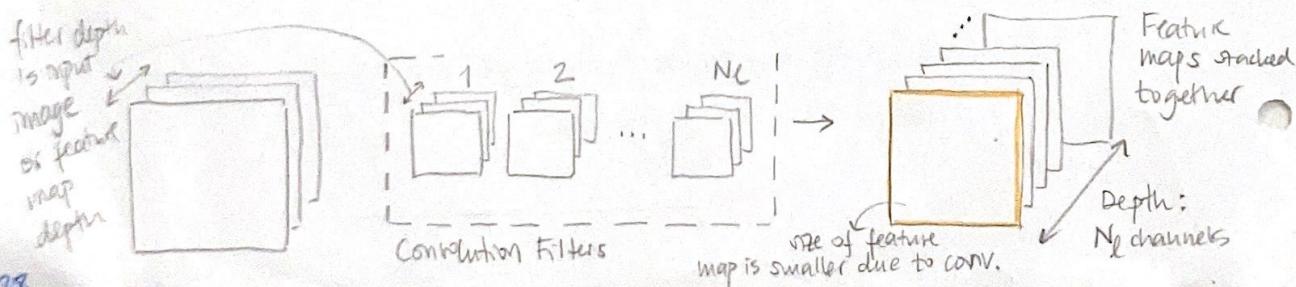
CNN LAYERS AND FEATURE VISUALIZATION (CONT.) (from page 15)

Convolutional Layers



- Each convolution layer has N_f filters
- Each of the filters is a 3D tensor and produces one output
 - the depth of the filter is the number of outputs from previous layer
 - each output is called **FEATURE MAP / ACTIVATION MAP**, and is like a channel
 - if the input image is the original preprocessed image, the filters of this first convolution layer are going to have expectedly a depth of 3 : R, G, B.
- All outputs or **FEATURE MAPS** are stacked to form a 3D matrix of N_f channels

- Therefore:
 - each 3D filter produces a feature map, which is like a channel
 - Convolutional layers can be concatenated
 - each feature map highlights patterns
 - successive convolutional layers discover patterns within patterns



- CNN vs MLP (dense or fully connected layers)

- In a fully connected or dense network all previous and successor nodes are connected, whereas in CNNs nodes are only LOCALLY connected
- Both have weights & biases, but in the case of CNNs these are in the form of FILTERS: the CNN learns which patterns to detect from the data.
- Training & inference works the same way.

- Notes:

- Images are usually normalized either to $[0, 1]$ or even $[-1, 1]$

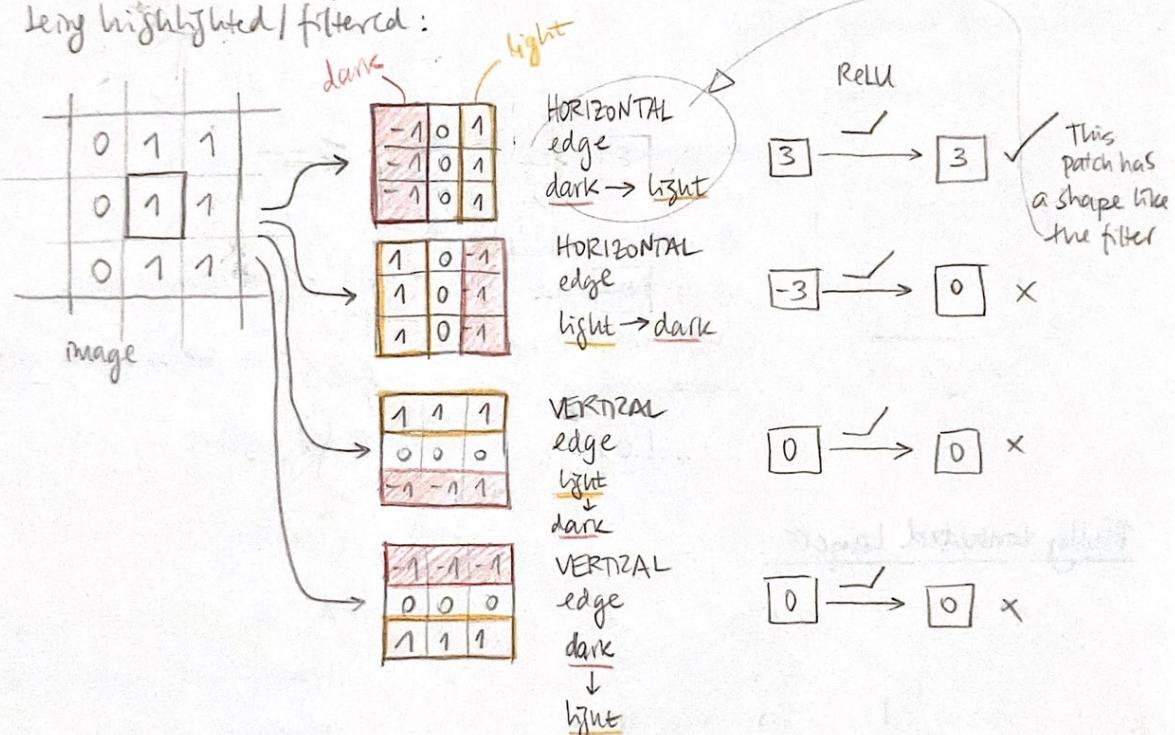
$\div 255$

$[0, 1] \rightarrow (0.5)$ mean
 (0.5) std

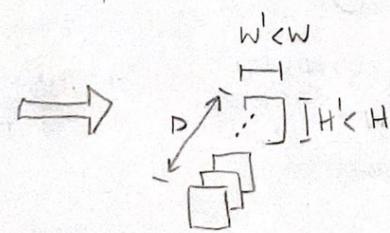
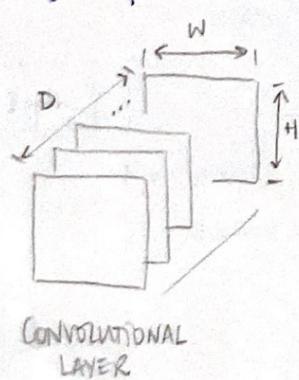
ToTensor() does that

Normalize()
does that

- filter kernel values give an idea of what's being highlighted / filtered:



Pooling layers: usually come after a convolutional layer and they reduce dimensionality



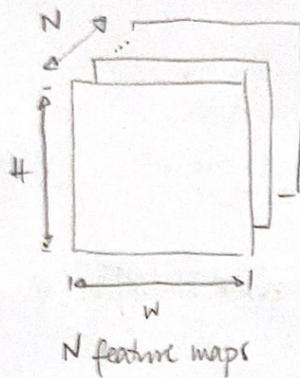
POOLING LAYER:

selection or computation of pixel values
reducing size of feature maps — BUT same depth = # feature maps

• Convolutional layers have usually many channels / filters, which results in many features
 \Rightarrow many parameters \Rightarrow
 \Rightarrow OVERFITTING
 \Rightarrow dimensionality reduction is the solution!

- There are several pooling layer types (see documentation of PyTorch). 2 most known:

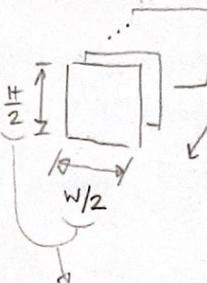
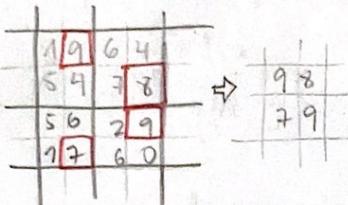
1) Max Pooling



Parameters

- window size, eg 2×2
- stride, eg 2

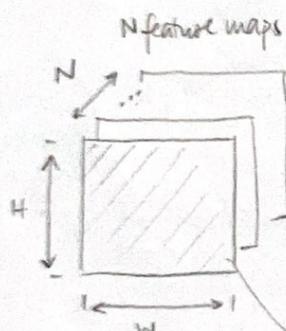
for each map, separately,
max pixel value within
a window selected



N : number of output
feature maps
is the same!

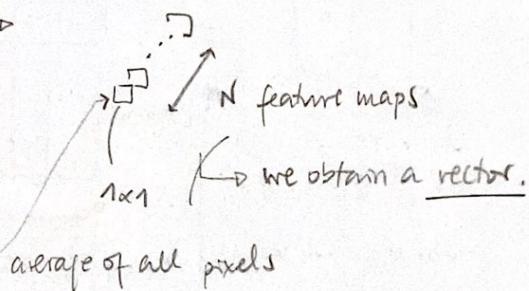
in the particular
case of window size 2×2
& stride 2

2) Global Average Pooling

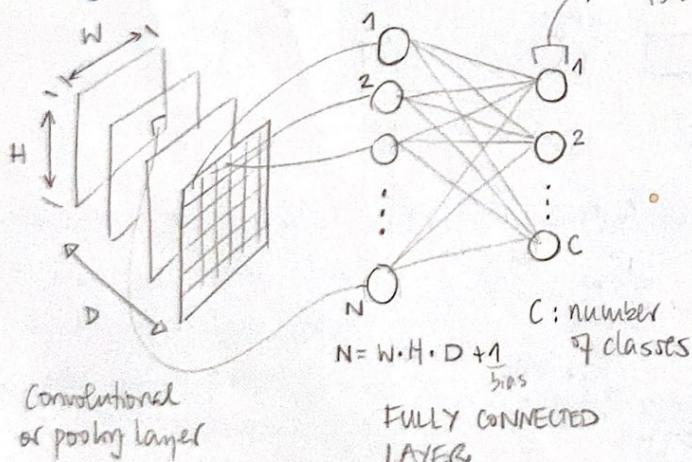


No params defined

Average pixel/node
value taken for
each feature map
independently



Fully Connected Layers



the output
is a vector of classes!

- A fully connected layer maps
linearly all pixels of the last
layer to C nodes, being usually
 C the number of classes we
want to predict.

$$h = \underbrace{W \cdot x}_{C \times 1 \text{ classes}} + \underbrace{b}_{C \times (W \cdot H \cdot D + 1) \text{ weights bias}}$$

feature maps vector

- Several fully connected layers can appear at the end one after the other
 - Intermediate f.c. layers usually have ReLU activation
 - LAST f.c. layers ends up usually with a SOFTMAX or a LOGMAX
 - ↳ this converts the output to be class probabilities.

- DROPOUT

- Dropout consists in randomly deactivating nodes, which reduces overfitting
- This is implemented in practice adding a layer which zeros outputs
- The probability p of deactivating a node is our param, which follows Bernoulli dist.

Network Structure : Which layers should we use when?

- To prevent overfitting

1. Early stop: track loss & accuracy and stop if loss \downarrow but accuracy \downarrow
2. Regularization (I understand that should be used in torch.optim? How?)
3. Dropout
 - Usually at the end of the network
 - Example: between fully connected layers \Rightarrow we prevent nodes from overdominating
4. Batch normalization

↳ L2 regularization implemented
through param [weight_decay]
in optim function

- Usually at the beginning of the network
- Just as with dropout, a layer is added
 - { nn.BatchNorm1d (after linear)
 - nn.BatchNorm2d (after conv2d)
- What's this?

For each batch, when applied, the mean of the outputs is subtracted and the std. deviation divided

- Which is the effect?
 - The values of the hidden weights don't shift around, they're more stable
 - Each layer learns more independently from others
 - We can use larger learning rates

- Which layers?

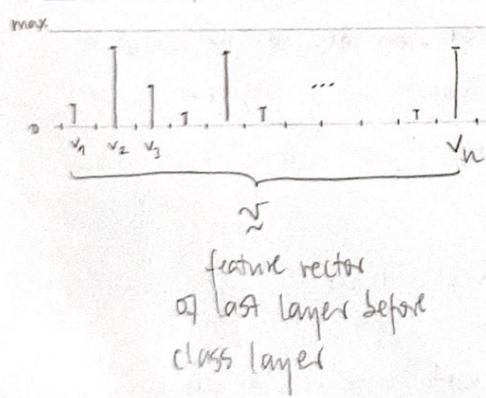
- In general: more layers help see more complex structures
BUT for simple structures, few layers are enough!
- LOOK AT KNOWN NETWORKS: and inspire on them! People already spent time figuring out how to arrange CNNs. 31

- Always try changes and see what happens
 - Number of layers
 - Size of conv. kernels
 - Loss & optimization functions
 - Hyperparameters: learning rate, momentum
 - Batch size

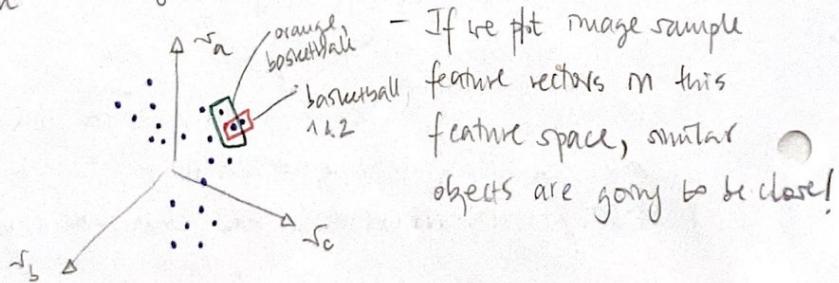
Feature Visualization

- Sometimes the network learns from the image incorrect ones that seem to work but don't generalize
 - Example: if we show wolfs in snow, the model will learn to detect wolfs when there's snow and a dog in the image!
- Therefore, we need to visualize the feature maps learned, to check the correct ones are learned.
- First layer
 - Usually highpass filters are learnt, which detect any kind of inclined lines or corners
- Successive layers
 - More complex shapes appear
 - As we go deeper, object parts, objects, start to appear; from layer 4/5, animals & people might appear
 - Field of view is increased in the filters (padding does that)
 - Important note: deep layers have as input the feature/activation map from previous layers, not images! Therefore, sometimes it's hard to understand
- Last layers / linear / fully connected layers: the output activation maps are vectors!!
- See images of conv. kernels in different layers (Matt Zeiler)

Last Feature Vector, PCA, t-SNE



- The feature vector before the class layer encodes objects with a characteristic object signature



- If we plot image samples feature vectors in this feature space, similar objects are going to be close!

→ Applications of that feature space

* Similar objects can be found using k -nearest neighbors

- as distance metric, L1 or L2 can be applied on samples

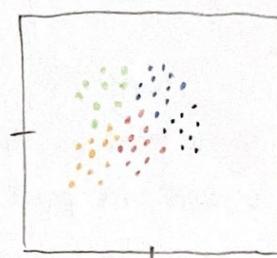
DIMENSIONALITY REDUCTION { * Principal Component Analysis (PCA) would be applied on all samples to reduce the dimensions of the feature space to a smaller significant set of features

* t-SNE (/t-dSNE/) : t-distributed stochastic neighbor embeddings

- non-linear dimensionality reduction
- data is separated in clusters in a low dimensional feature space, eg 2D.

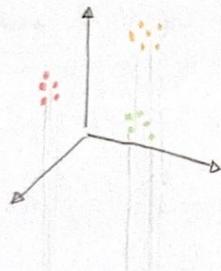
MNIST dataset
28x28
feature space

because t-distribution is used



2D feature space with 10 clusters (0-9 digits)

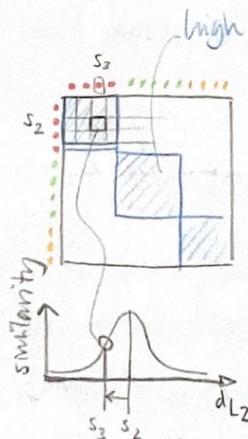
• Method intuition



- samples are projected onto a low dimensional space and similarity matrix is computed

- the matrix is expected to be sparse

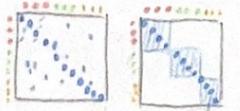
- the projected samples are moved step by step and similarity matrix is recomputed so that the final matrix resembles the original one from high dim space



high similarity values

- similarity matrix of all samples against all samples in high dimensional space
- similarity between samples s_i & s_j is the prob. of both belonging to same cluster, which is evaluated by the distance between the samples, eg L2 distance

$d_{L2} = \|s_i - s_j\|_2$: the scaled normal distribution of samples in cluster is computed



samples moved in low dim space

Occlusion, Saliency and Guided Backpropagation: What is the network looking at for the decision?

HEATMAPS: ROI relevance for class by experimenting with occlusion

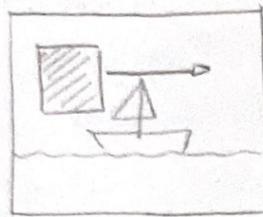


Image is occluded with a mask/box and its class computed for each box position (box is moved)



for each of the pixels in the box, the decrease of the class confidence when overlaying the box is accumulated

Relevant image areas that lead to the predicted class are computed and visualized encoded in color.

- very relevant region
- relevant region
- ...

SALIENCY MAPS

- Which pixels are most important for classifying an image?
- The gradient of the class score wrt pixel values is computed for each pixel
 - ↳ this creates a map similar to the heatmap
 - ↳ $\nabla C = \left[\frac{\partial C_1}{\partial p}, \frac{\partial C_2}{\partial p}, \dots \right]$
 - ↳ Image segmentation can be done this way: the saliency map is overlaid and thresholded to get a mask.

GUIDED BACKPROPAGATION

- Same idea as in the saliency map, but we look at the change in the output of a neuron/node, when a pixel value is changed Δp .
- Big neuron output changes indicate that the pixel has a big influence on the neuron
 - (- Basically the chain rule is applied as in backpropagation?)

Some Applications with Feature Maps

- Deep Dream: selected features are amplified in an image
- Style Transfer: content and style are captured with feature maps from 2 images and then they're combined

SKIN CANCER DETECTION (optional excusus) - Mini Project

Precision, Recall, Accuracy & CONFUSION MATRIX

		ACTUAL/TRUTH		PRECISION	
PREDICTED	+			$\frac{TP}{TP+FP}$	
	-				
		TP	FP type I error	$\frac{TP}{TP+FN}$	
		FN type II error	TN		
		ACTUAL +	ACTUAL -	ACCURACY = $\frac{TP+TN}{TP+FP+FN+TN} = \text{TOTAL}$	
				how well does our prediction perform?	

- FP: type I error: false alarm

"cancer prediction in cases where there's no cancer"

- FN: type II error: missed case

"sane prediction in cases where there's actually cancer"

- In medicine context, often appear:

• SENSITIVITY = RECALL \rightarrow "of all sick people, how many did we diagnose sick?"

• SPECIFICITY = True negative rate \rightarrow "of all healthy people, how many did we diagnose healthy?"

\rightarrow very dangerous!

$$\text{PRECISION} = \frac{TP}{TP+FP}$$

$$\text{RECALL} = \frac{TP}{TP+FN}$$

$$\text{ACCURACY} = \frac{TP+TN}{TP+FP+FN+TN}$$

"Among all patients PREDICTED to have cancer, how many actually have it?"

"Among all patients that ACTUALLY have cancer, how many were predicted to have it?"

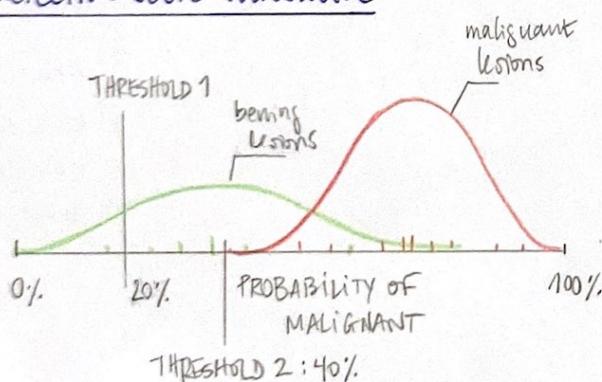
$$\text{True negative rate} = \frac{TN}{TN+FP}$$

SPECIFICITY

$$\text{False positive rate} = \frac{FP}{TN+FP}$$

actual negative

Detection Score Threshold



- If the distributions of the benign & malignant lesions for the variable PROBABILITY OF MALIGNANT are unknown, it's better to take a low threshold, because the cost of regret is much lower! Better to order new tests than missing a cancer diagnosis \rightarrow THRESHOLD 1

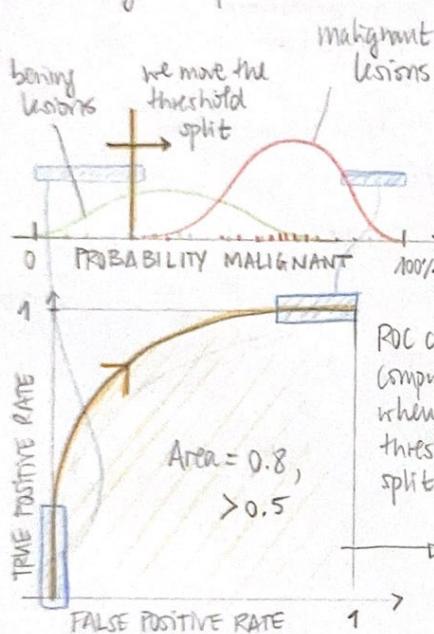
- If the distributions are known, we take THRESHOLD 2, which eliminates FNs or type II error!

ROC Curve = Receiver Operating Characteristic Curve

- Plot that illustrates the diagnostic ability of a binary classifier.
- There are several types of ROC plots; in many cases, they are created by sweeping the threshold along the classification probability and computing true/false positive/negative rates
- Analysis of 4 cases & 2 types of ROC curves

CASE 1: usual case

Distributions of malignant & benign overlap



CASE 2: perfectly disjoint - nice to have scenario

to have scenarios

100%
1
100%
1
100%
1
100%
1

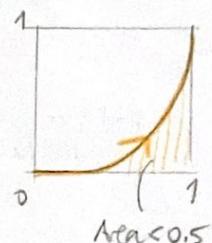
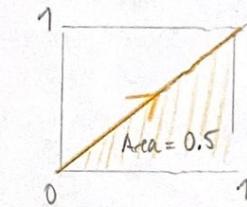
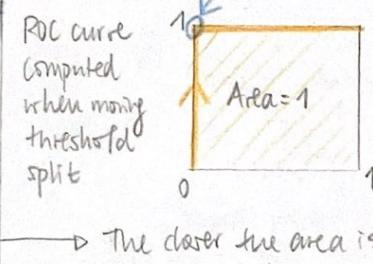
CASE 3: perfect overlap - classification doesn't work, random values (e.g., not trained)

values (e.g., not trained)

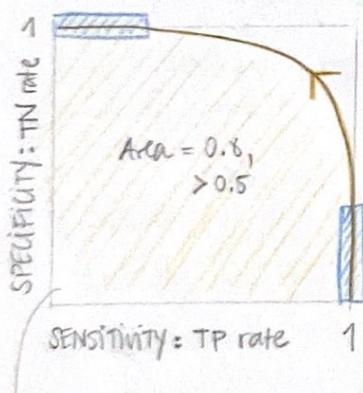
scenarios!

CASE 4: predicted probability is the inverse to the real scenarios!

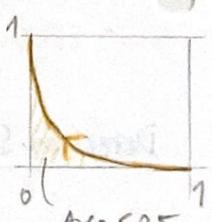
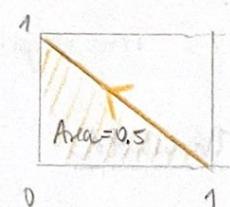
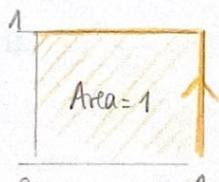
scenarios!



The closer the area is to 1, the better our model is and the better splitting threshold we can take



common ROC plot in medical context



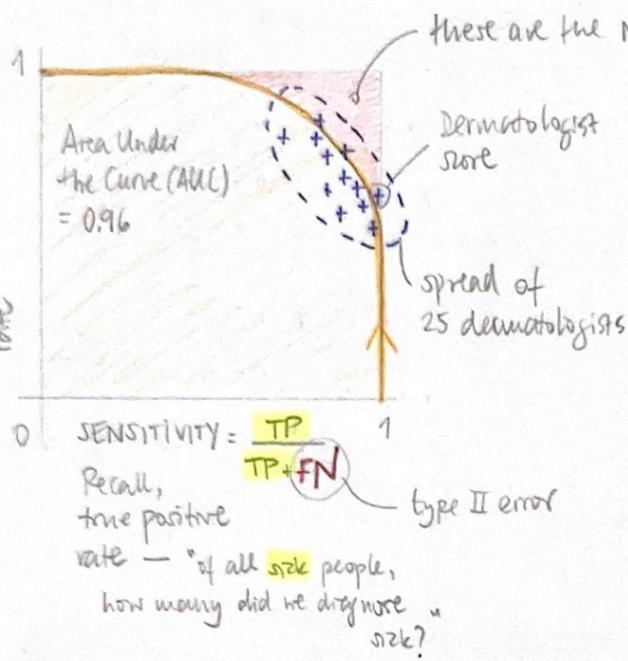
In this scenario, we would actively kill people...

Type I error

$\frac{TN}{TN+FP}$

SPECIFICITY : true negative rate

- Interpretation of the ROC curve (in the particular case of medical context)



these are the MISS DIAGNOSED cases

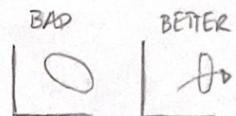
	SPECIFICITY	SENSITIVITY
HIGHER	More accurate ability to say who is healthy	More accurate ability to say who is sick
LOWER	We tend to diagnose healthy people as sick CONSERVATIVE, but more COSTLY	We miss to detect sick people BAD!!! (apparently cheaper)

"of all healthy people, how many did we diagnose healthy?"

Increasing Specificity makes the test more efficient

SENSITIVITY is critical and should be as high as possible

- Dermatologists are usually below the curve \Rightarrow less performance
- The spread of the dermatologists is quite large!
- The spread should be as right as possible (high SENSITIVITY)
 - ↳ some dermatologists are to the left: is that so because it's cheaper for the insurance in the short-term?



t-SNE Visualization (see p. 32, 33)

- Very intuitive projection of the feature space on a reduced space that allows for clustered visualization of the samples
- If visualized images of the clustered samples, similarities are clear to the human eye!!!

Visualization of what the network is looking at = SENSIBILITY ANALYSIS (see p. 34)

- What is the network looking at to take that classification decision?

- Sensibility analysis: we take samples, change them and observe how the classification output changes \rightarrow this gives us an idea of what's the net looking at

↓ METHODS: heatmaps, saliency maps, guided backpropagation \rightarrow TRY THEM & DRAW CONCLUSIONS! 37

Confusion Matrix - How good is our model?

		TRUTH	
		SRK	healthy
PREDICTION	SRK	TP	FP type I error
	healthy	FN type II error	TN

type I error: healthy patients diagnosed as SRK \rightarrow More tests, costly

type II error: sick patients diagnosed as healthy \rightarrow DANGER!

- Main diagonal should be 1 and rest 0 \Rightarrow ideally
- Everything off the diagonal is a confusion!
- Dermatologists had a more diffuse confusion matrix \Rightarrow more miss classifications

3 sicknesses \rightarrow TRUTH

		A	B	C	Σ
PREDICTION	A	0.8	0.08	0.3	$\Sigma \neq 1!$
	B	0.1	0.9	0.1	
	C	0.1	0.02	0.6	
Σ		1	1	1	

$P_{ij} = 0.1$ probability of patients having A but of being diagnosed with B

i: DIAGNOSIS
j: TRUTH

probability of getting misdiagnosed
if you suffer C: $0.4 = 0.1 + 0.3$

POPULAR NETWORKS (Personal Review)

AlexNet

Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton
U. Toronto

2012

CONVOLUTIONAL

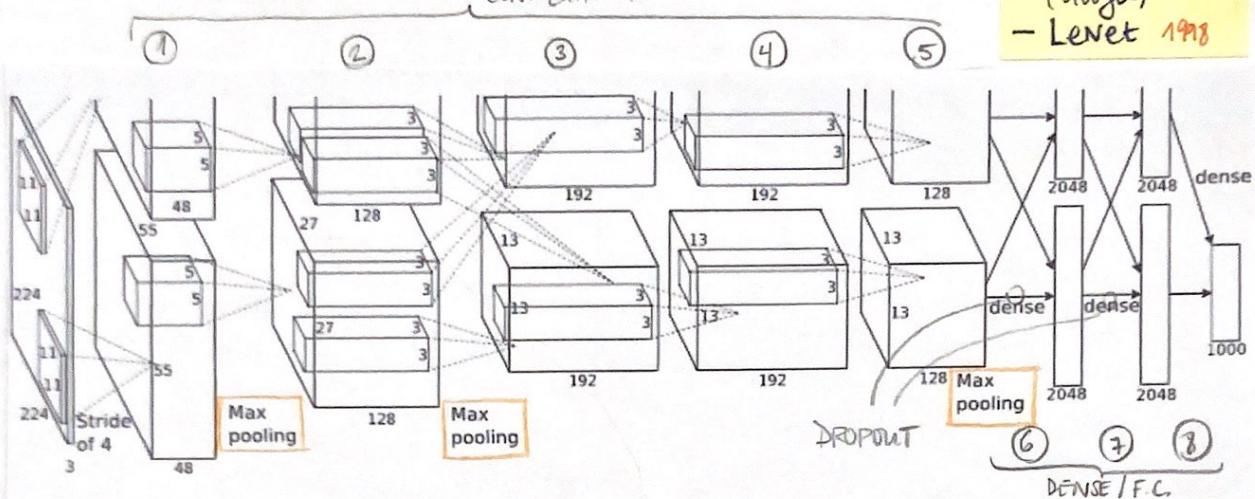


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

- Competed in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2012
 - ↳ top-5 error: 15.3% : 10.8 points better than best ever!
 - ↳ people started paying attention to AI/Deep Learning
- Properties
 - 2 GPUs used: image top & bottom → faster
 - ReLU: better than tanh(x) or sigmoid
 - 8 layers:
 - CONV → CONV → CONV → CONV → CONV → DENSE → DENSE → DENSE → 1000 classes
 - Max pooling ↑ Max pooling ↑ Max pooling ↑
 - 1 2 3 4 5 6 7 8
 - 60 M parameters, 650K neurons
 - Data augmentation used for training (to prevent overfitting)
 - Dropout with $p=0.5$ on the first 2 fully connected (dense) layers
 - Momentum 0.9, weight decay 0.0005

- AlexNet 2012
- VGG-16 2014
- ResNet50 2015 (Microsoft)
- Inception-v3 2015 (Google)
- Lenet 1998

LeNet

Lecun, Yann
1989 - 1998

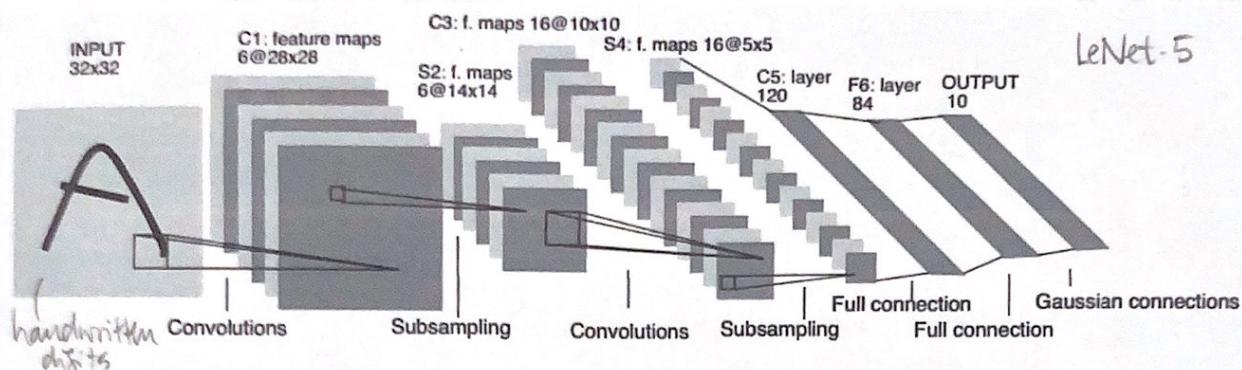
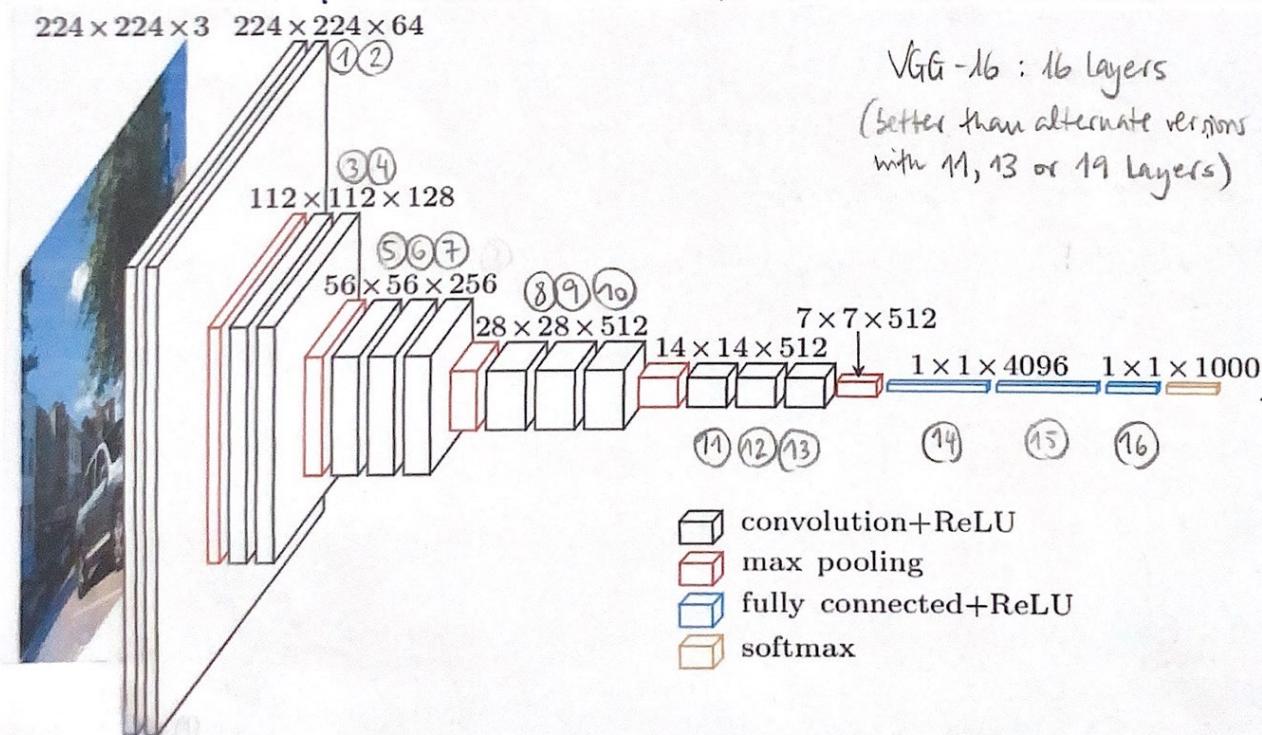


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical. (1998)

- Very old architecture, but the inspiration/base for later architectures like AlexNet
- LeNet was presented by Yann Lecun and the idea evolved between 1989-1998; the final implementation is LeNet-5 (the one that's usually meant)
- LeCun applied for the first time
 - Convolutional layers
 - Backpropagation

VGG-16

Very Deep Convolutional Networks for Large-Scale Image Recognition
Simonyan, Zisserman - U. Oxford, UK, 2015



VGG-16 : 16 layers
(better than alternate versions with 11, 13 or 19 layers)

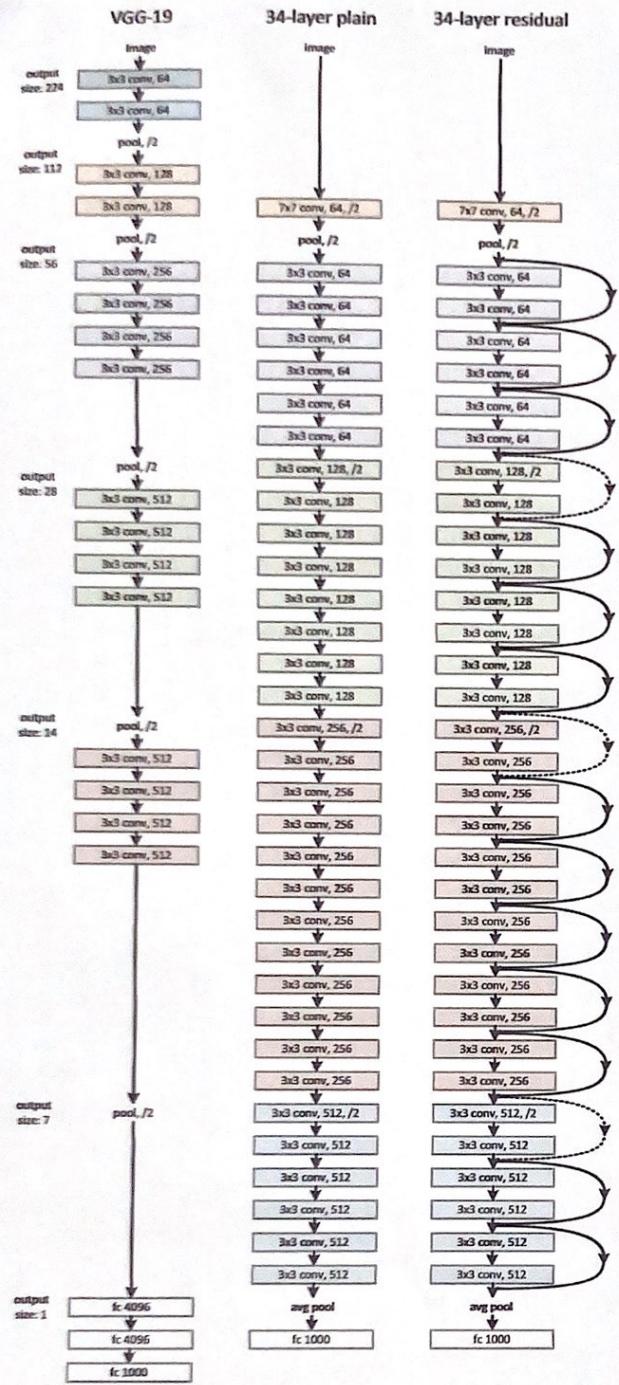
- Created by the Visual Geometry Group: VGG
- Properties
 - They discovered smaller filters (3×3) are enough \Rightarrow less params \Rightarrow faster & less images needed!
 - They tried increasing and local normalization response (LNR) and less overfitting!
 - proposed by AlexNet
 - neuron output inhibition strategy inspired in human biology; ReLU is not bounded and normalizing groups of outputs was supposed to yield better results
 - BUT: nowadays not used anymore, because not efficient; much better: Dropout
 - RESULT (BEST): VGG Batch Normalization with 16 layers (better than 19 Layers!)
 - \hookrightarrow ablation study: remove layers and see what happens
 - Input images were 224×224 , but multi-scale training was performed to generalize to different object scales.
 - Take bigger picture sizes and crop to 224×224
 - This improves accuracy
 - They did something with the final dense layers during testing (but not training) — they were replaced by 1×1 convolutions, ie spatial averages ???

ResNet

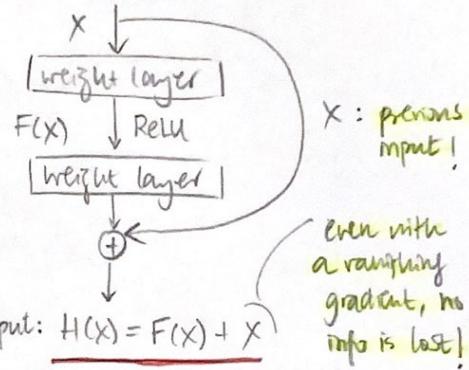
He, Zhang, Ren, Sun - Microsoft

2015 (from Medium article by Sik-Ho Tsang)

- ResNet can have very deep networks (up to 152 layers)
 - \hookrightarrow this can be achieved by learning the residual representation instead of the actual signal / feature maps of the images
- They introduce skip or shortcut connections: input from previous layer is taken without any modification
 - \hookrightarrow They call plain networks all previous that don't do that (AlexNet, VGG, etc) and can't achieve so many layers because of the vanishing / exploding gradients problem
- VANISHING/EXPLODING GRADIENT Problem of plain networks
 - It's not uncommon to have small or large derivatives in the network
 - $\partial E / \partial w$ is a multiplication of derivatives (s.p. 19), due to the chain rule
 - If the network is deep, this product is VERY SMALL/LARGE! (backpropagation)



— The skip / shortcut connection solves the problem of the vanishing gradient:



BUT the weights are optimized during learning with the **RESIDUAL MAPPING**:

$$F(X) = H(X) - X$$

- They additionally use a **BOTTLENECK design** with **1x1 convolutions**
- Bottleneck designs reduce complexity
 - Adding 1x1 convolutions we go from 34 layers to 50 layers: ResNet 50

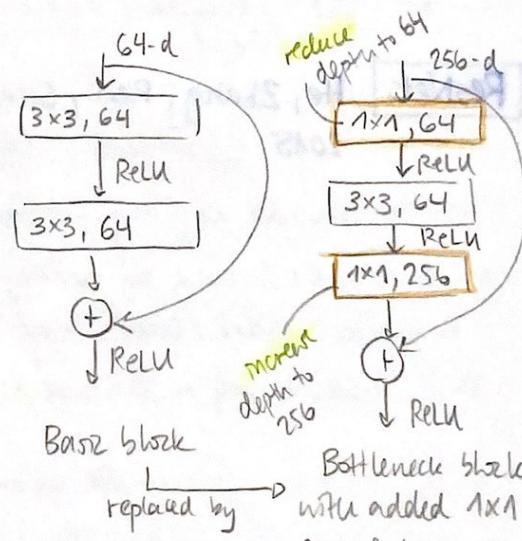
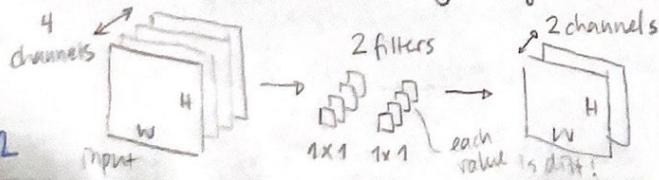


Figure 3. Example network architectures for ImageNet. Left: the VGG-19 model [41] (19.6 billion FLOPs) as a reference. Middle: a plain network with 34 parameter layers (3.6 billion FLOPs). Right: a residual network with 34 parameter layers (3.6 billion FLOPs). The dotted shortcuts increase dimensions. Table 1 shows more details and other variants.

- A 1x1 convolution layer reduces the number of connections
(= parameters) without degrading network



- a 1x1 convolution can be used to modify the depth of the signal:
reduce / increase

- Result / Conclusion

- Networks can be deeper while
 - having less parameters \Rightarrow faster to train and use]
 - Increasing accuracy] \rightarrow in other architectures, increasing layers up from a certain point (eg VGG 16), decreased accuracy due to vanishing gradient issues
- As we increase layers { accuracy \uparrow } \downarrow speed \downarrow \Rightarrow ResNet50 is often a good alternative nowadays (2020)

other shallower architectures have much more parameters!

Inception v3

Szegedy, Vanhoucke, Ioffe, Shlens, Wojna - Google

2015 (from Medium article by Sik-to Tsang)

- There are several versions

- The first version (2015), was called GoogLeNet
- Later versions introduced BATCH NORMALIZATION and extended FACTORIZATION

(see p. 31)

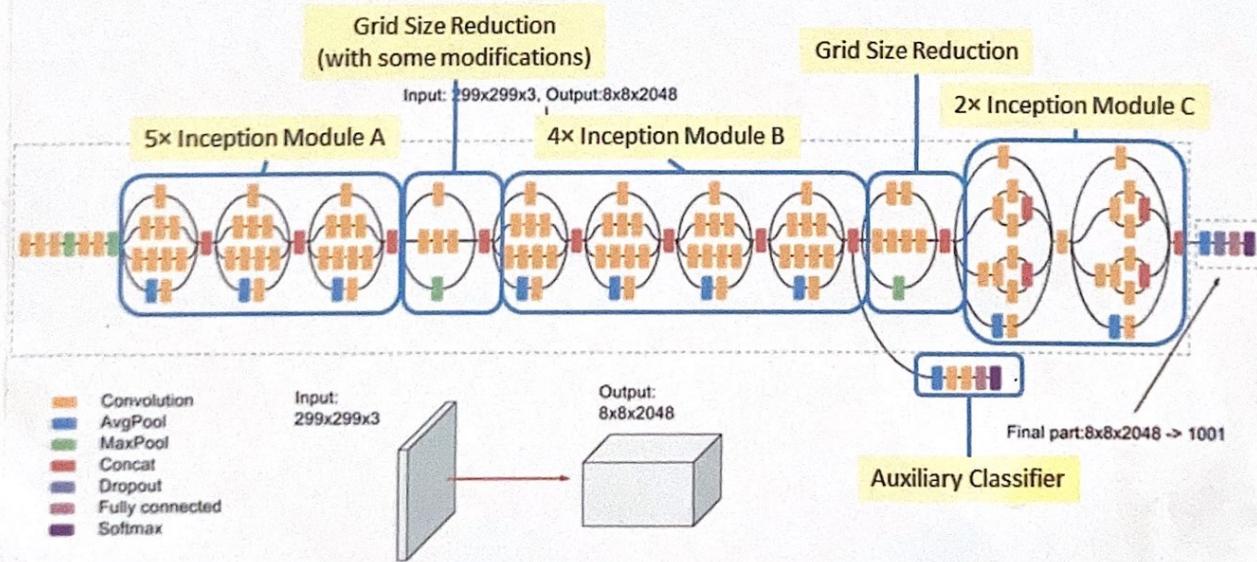
output of each batch normalized (with mean + std) to avoid shifting of weights

larger filters (eg 5x5) are replaced by smaller ones (eg 3x3) that work in parallel and then result is concatenated \rightarrow this reduced number of params without degrading network efficiency

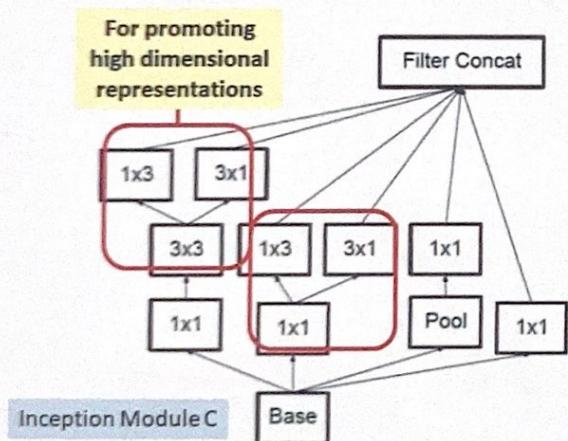
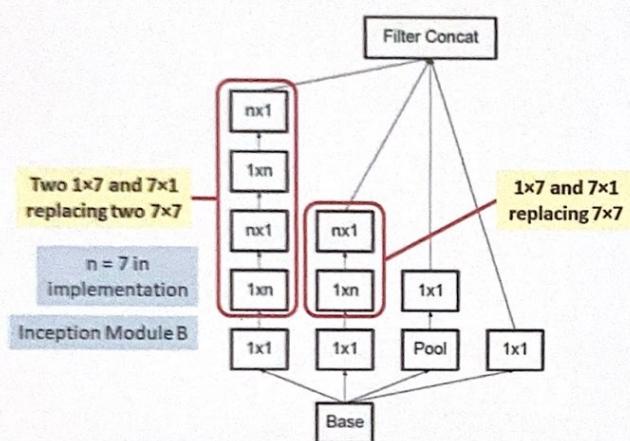
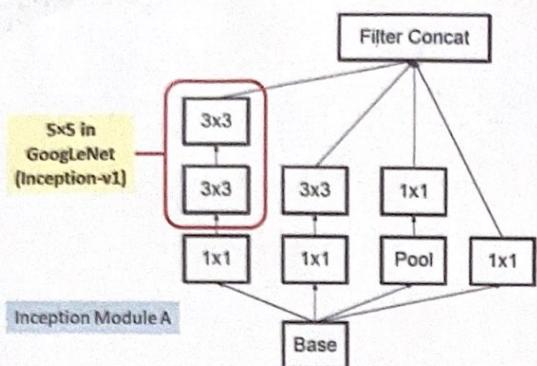
Factorization example

• 5x5 filter: $5 \times 5 = 25$ params

• 2 3x3 filters: $2 \times 3 \times 3 = 18$ params



- They define 3 types of factorization modules, here informally named A, B, C (i.e. msg)



- They have an extra branch from which an auxiliary classifier emerges
→ it's used as a regularizer (How?)

- They also use a grid size reduction module

- AlexNet & VGGNet use feature map downsampling with max pooling
- Maxpooling after a convolution can result either coarse (pooling) or expensive (convolution operation alone anyway)
- Their proposed grid reduction module processes feature maps in 2 parallel branches using convolutions and pooling and it leads to a **LESS EXPENSIVE & STILL EFFICIENT** network.

- Conclusion (my observations)

- What they aim is to reduce the cost of training / using the network while still having a complex & deep architecture; FACTORIZATION & GRID SIZE REDUCTION are used for that.

DenseNet

Huang, Liu, van der Maaten, Weinberger, 2018/7
Best Paper Award CVPR 2017

Summarized from
the notes by
Mukul Khanna
at medium

- Issue for deep networks: VANISHING GRADIENT: with very deep networks, the gradient becomes small when back-propagating it (recall chain rule) and first layers don't learn efficiently
 - Several architectures have been introduced to tackle this, e.g. ResNets, Inception.
 - All of them create some kind of paths between layers

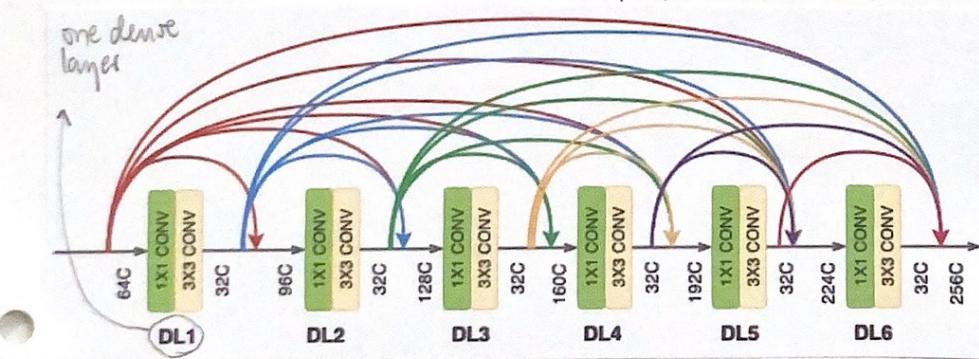
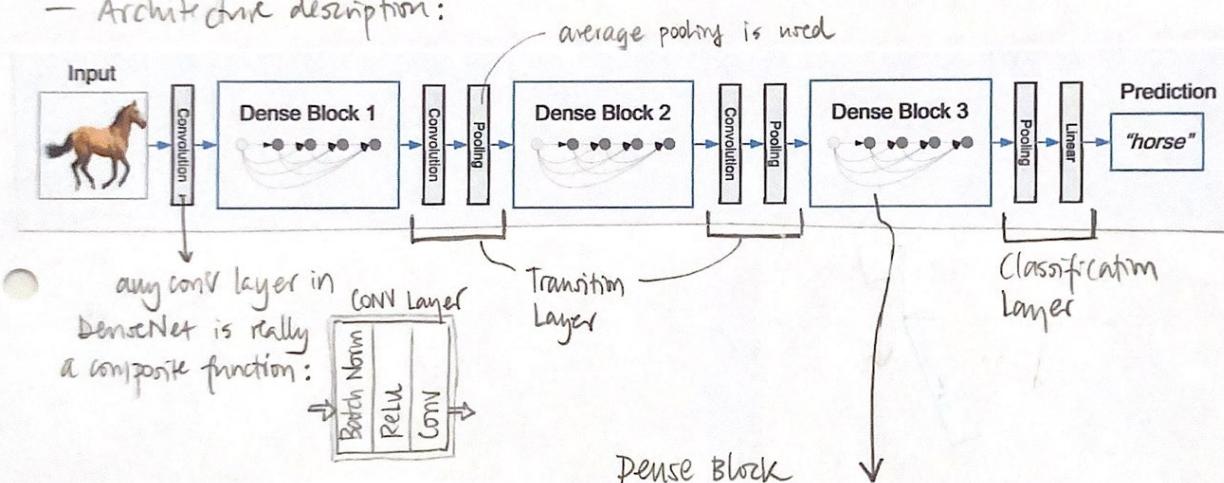
Main idea:

- The architecture is composed of dense blocks with layers
- Each layer from a dense block receives feature maps from all preceding layers and these are fused through CONCATENATION, not summation (like in ResNets)

Consequences / contributions

- Vanishing gradient is alleviated
- Number of parameters is considerably reduced!
- Features are learnt collectively: there are reused and propagated

Architecture description:

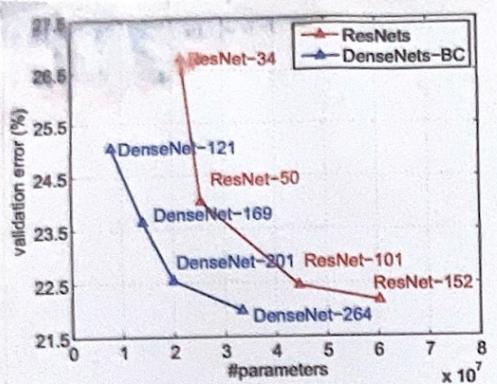


- There are several DenseNets, all have 4 dense blocks, 3 transition layers, 1 classification layer

- The difference is the # of layers within a block

- Example: DenseNet-121:

Block	1	2	3	4
Dense Layers	6	12	24	16



Results: Comparison to ResNet

DenseNets require much less parameters to achieve the same accuracy !!!

recall this architecture also targets the VANISHING GRADIENT problem!

Udacity CV | 1 | INTRO CV 1

- Image representation
 - RGB channels
 - Pixels & image access & histograms
 - Color spaces
 - Basic classification
- Fourier
- Convolutional filters
 - low-pass: Gauss
 - high-pass: Sobel
 - general
- Canny edges
- Hough transform: lines & circles
- Haar cascade: classifiers trained for edge detection

Udacity CV | 1 | INTRO CV 2

- Harris corner detection
- Morphological operations
- Contours: contour features (shape info)
- Image moments
- K-means clustering
- Feature vectors / descriptors
- ORB
- HOG
- Deep Learning intro

UDACITY CV | 1 | Deep Learning

- Perceptron
- MLP
- Feed-forward
- Error / loss
- Gradient descent
- Backpropagation
- Training
 - random start
 - other activation functions
 - SGD
 - learning rate
 - momentum
- Pytorch
 - split: train / test
 - overfitting
 - early stopping
 - regularization
 - dropout