

# Convolutional Neural Networks

①

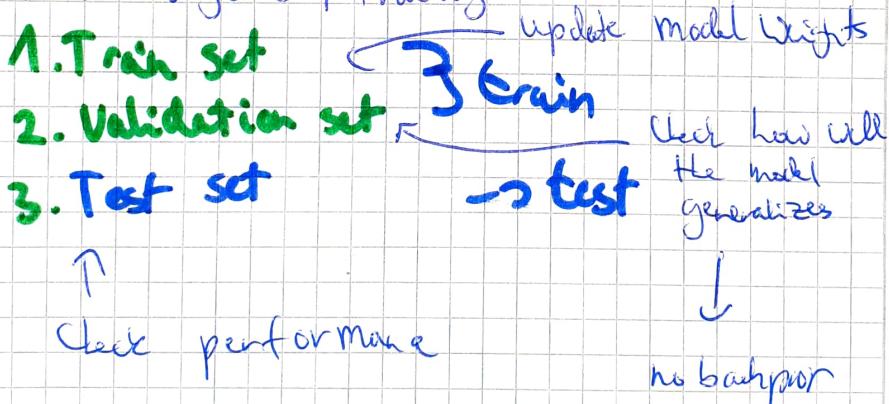
- Normalize Inputs so Images are in some Range (e.g. 8bit)

↳ mapping of colors, (camera settings, ...)

## Train batch based:

1. clean gradients of optimized variable
2. forward pass, compute predicted outputs by passing inputs to model
3. calculate loss
4. Backward pass: compute gradients of the loss w/ respect to our params
5. perform optimization step.
6. update average training loss

When do you stop training?



## MLP vs CNN

- MLP: only FCL
- MLP: 2D info is lost
- too many parameters

## CNN

- sparsely connected
- multi input
- > MLP in the end

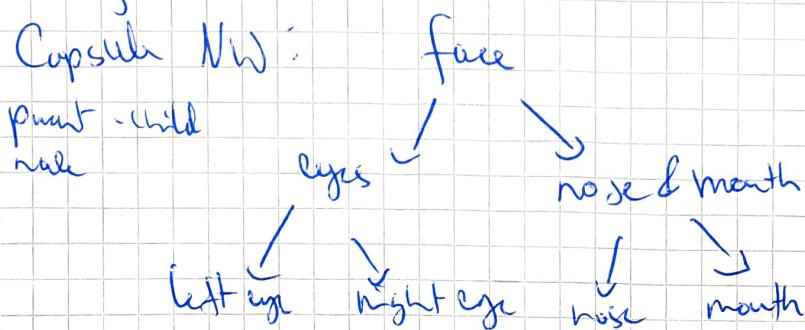
## Filters

②

- Learned by the Network to detect edges
- Detect faces more complex & Stages, the deeper the network
- One filter Dimension depends on image depth
- move filter above over image
- Multiple filters are used, stack them like images (dataspaces)
- Locality Connected, parameter sharing
  - Stride  $\Rightarrow$  how do we slide image?
  - Padding  $\rightarrow$  for edges to not lose information
- Pooling: Reduces Dimensionality
  - $\hookrightarrow$  MaxPool (takes max. value of window)
  - $\hookrightarrow$  also have stride
  - $\hookrightarrow$  no parameters!  $\rightarrow$  Dropout doesn't make sense
  - $\hookrightarrow$  avg-pool, max-pool, ...
  - $\hookrightarrow$  Information is lost (Spatial)

## Alternative to pooling: Capsule Networks:

$\rightarrow$  extract spatial information between parts/objects

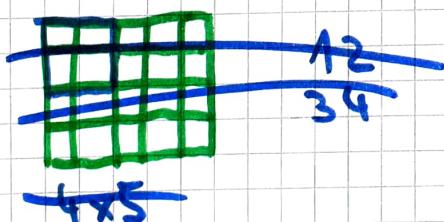


# Capsules

- Collection of nodes, contain info about a specific part & it's properties
- Each capsule outputs a  $\vec{v}$  vector w/ some magnitude & orientation
  - Vector allows to build up a parse tree to recognize whole objects consisting of small parts
    - ↳ Magnitude should be high, no matter the orientation
  - Dynamic Routing is used to find best connections from one layer of capsules to the next → even good when objects overlap

## Pre-processing

- Images need to be resized



130 x 130 input       $\rightarrow$       Weights per Filter:  $F^2 \cdot D_{in} (\text{e.g. } K \text{ filters})$

Output Size:

$$(W - F + 2P) / S + 1$$

$\uparrow$        $\uparrow$        $\uparrow$   
Width      Filter Size      Paddles

Width Filter Size

④

## Augmentation

- Scale invariant
- Rotation invariant
- Big Pooling → position not as relevant anymore

⇒ Add augmented images

↳ Translation

↳ Rotation

↳ fight Overfitting

→ torchvision transforms library

→ AlexNet → Dropout & ReLU

→ VGG, simple, elegant  $3 \times 3$  Conv,  $2 \times 2$  Pool

→ ResNet, Shortcut Connections

↳ vanishing Gradient Problem

→ Visualizing Feature Maps

↳ Edges

↳ stripes

↳

Objects

Deep Dream

# Transfer Learning

- pre trained Networks
- transfer it to other task
- save computation
- last filters are most specific
  - ↳ Only those need to be trained

Depends on:

- size of data set
- similarity of data

In cases:

Large data, different data

Small data, different data

Small data, similar data

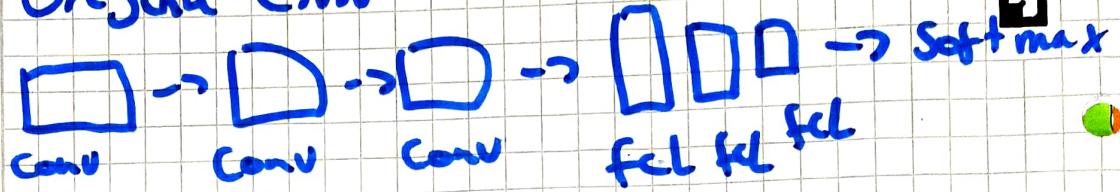
Large data, similar data

Guide for TL

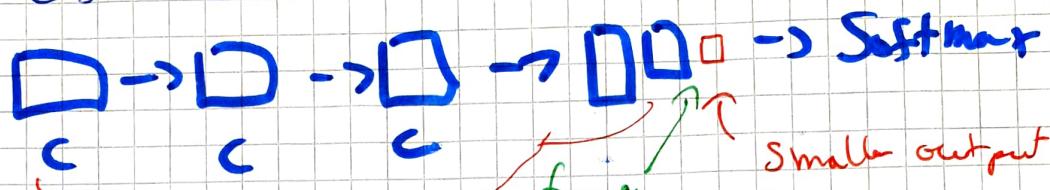
To	From	fine-tune	fine-tune or retrain
Large Data	Small Data	End of ConvNet	Start of ConvNet
Similar	Similar	different	
Similar to Training Data			

→ Network might need to be tweaked  
(Output layers)

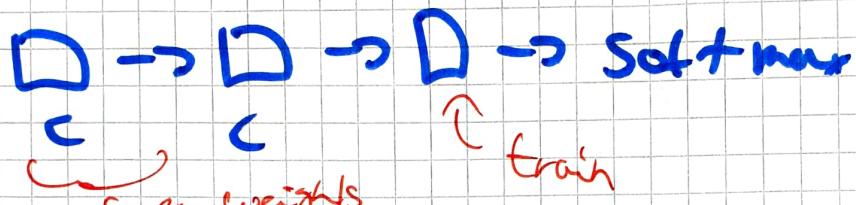
## Original CNN



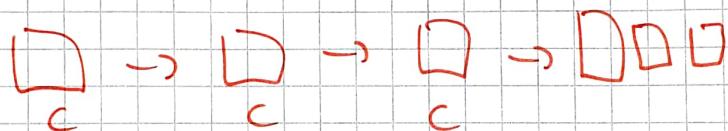
Case Small Data set, Similar Data



(use small data set, different data)

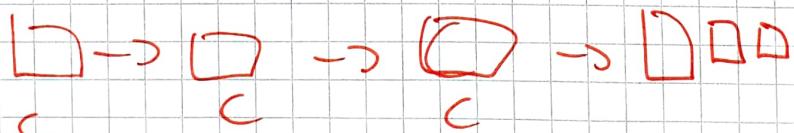


(ex Large Data, Similar Data)



retrain entire Network  $\rightarrow$  pre-trained weights as a start  
in conv layers

Case Large Data, different Data



fine-tune on retrain entire Network

$\rightarrow$  initialize <sup>with</sup> pre-trained weights might be good

## Freeze weights:

for param in vgg16.features.parameters():  
param.requires\_grad = False

7

## Training loop w/ Validation

Valid-loss-min = np.inf

for epoch in range(1, n\_epochs):

train-loss = 0.0

Valid-loss = 0.0

for batch\_i, (data, target) in enumerate(train\_loader):

if train-on-gpu:

data, target = data.cuda(), target.cuda()

optimizer.zero\_grad()

output = vgg16(data)

loss = criterion(output, target)

loss.backward()

optimizer.step()

train-loss += loss.item() \* data.size(0)

# Validation

for data, target in valid\_loader:

output = model(data)

loss = criterion(output, target)

valid-loss += loss.item() \* data.size(0)

train-loss = train-loss / len(train\_loader.dataset)

valid-loss = valid-loss / len(valid\_loader.dataset)

If valid-loss <= valid-loss-min:

valid-loss-min = valid-loss

model =

torch.save(model.state\_dict(), "model.pt")

Save best  
model



8

# Weight initialization

## Constant weight

↳ close to zero (nn.init)

→ all zero weights:

→ bad accurate scores

→ CEL, backpropagation is hard to do,  
sources difficult to tell (less than 10%)

## Instead: random uniform distribution

→ Initialize each weight uniform

→ create function to access weights & change them

Good range:  $\rightarrow [-y, y]$  where

$$y = 1/\sqrt{h} \quad h = \text{number of classes}$$

## Normal distribution

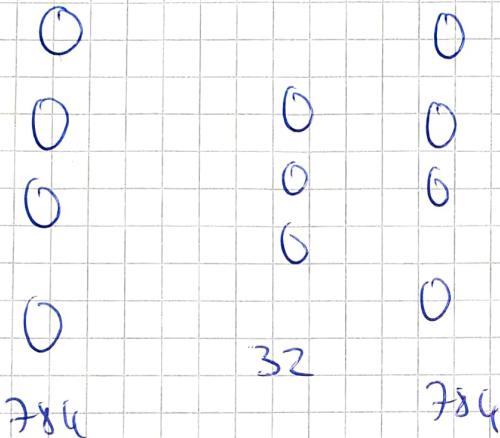
→ Equally distribution, Bell curve

→ more values closer to zero

# Autoencoder

## Encoder $\leftrightarrow$ Decoder

- learns efficient compression / decompression
- good for denoising
- not concerned with labels
- output dimension = input dimensions



MSE loss is better!

Adam optimizer is better!

for data in train\_loader:

→ labels don't matter!

→ Compare images to reconstructions

↳ sometimes blurry → Convolutional Autoencoder

# Conv Autoencoder

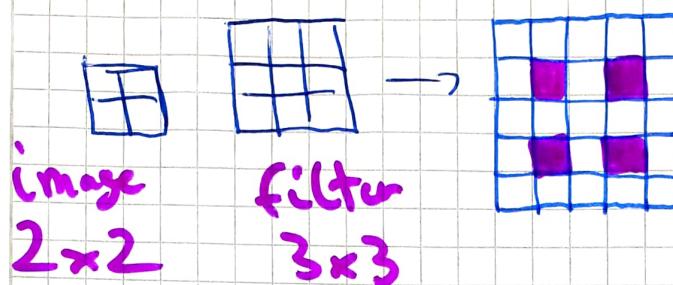
- encoder like normal NN
- decoder: Reverse Down-sampling → Upsampling
- Interpolation like Nearest Neighbor → Only copying → learn Up-sampling → Transpose Convolutional Layers

# Transpose Conv Layer

⑩

→ Stride = downsample factor

→ We have large filter than image



Stride = 2

⇒ upsample  $2 \times 1$

→ Can throw artifacts, upsampling with Nearest Neighbor or Bilinear Interpolation

(or: Kernel size =

Overlap)

→ upsample with replacement of upsample 2x2 & normal conv instead of transpose conv layer

→ Denoising

Denoising

→ Deeper Networks are required

→ increased depth by 1 conv layer  
in both Encoder & Decoder  
parts

# Style Transfer

- Merge Image with Style of another one
- CNN extracts features, later layers:  
objects
- Style: feature space for color & shapes
- See which features (color, shapes) are similar  
across filters
- Style & content image

VGG19 is used

- Conv Stacks (4 convs together)

(content)  $\approx$  (content)  
Representation Representation

from content image

from target image



from conv 4-2

Content Loss

$$L_{content} = \frac{1}{2} \sum (T_C - C_C)^2 \rightarrow \text{Minimize}$$

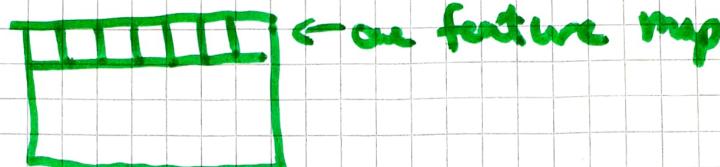
Should be close in representation even with  
different style!

Weights are not updated, only target image is  
updated (using backprop to minimize  $L_{content}$ )

12

Same for target & style image $\rightarrow$  Multi scale image in all 5 conv. layers

Gram Matrix

 $\rightarrow$  flatten feature maps to go from 3D to 2D $\rightarrow$  non localized information

$$\begin{matrix} 16 \\ 8 \end{matrix} \times \begin{matrix} 16 \\ 8 \end{matrix} = \begin{matrix} 8 \\ 8 \end{matrix}$$

Transpose                                          Gram Matrix

 $\Rightarrow$  Multiply features in each map $\rightarrow$  rates style of target image

so

## Style Loss

Mean Squared Error between style &amp; target image

Gram Matrices

$$L_{\text{style}} = \alpha \sum_i w_i (T_{s,i} - S_{s,i})^2$$

↑                      ↑  
Target Image      Style Image

Total loss  $\alpha L_{\text{style}} + L_{\text{content}}$   
 $\alpha$  content weight  
 $\beta$  style weight

- requires grad - (False)
- freeze weights

13

## Style Content Representations:

conv 1-1, conv2-1, conv3-1, conv4-1, conv4-2,  
conv5-1

### Style Loss:

- iterate through style layers
- same as content loss, but multiplied by weights

### Content Loss:

=  $\text{Each term} \left( \text{target\_features}['\text{conv4\_2}'] - \text{content\_features}['\text{conv4\_2}']^3 \right)^{\star\star 2}$

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

↑

Sensitivity

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

↑

Specificity but with TP at bottom

Input is  $224 \times 224$ :

~~222 → 220~~

~~220 → 216~~