

Date:

Convolution Neural Network

Sun Mon Tue Wed Thu Fri Sat

AIMS AND OBJECTIVES:

CNNs:

- Sparse connectivity
- Shared weights → reduced parameters.
- Captures spatial hierarchies
- Automatically learn hierarchical features from images
- Efficiently process larger images
- Superior performance on image related tasks

Consist of: Interconnected neurons, convolutional layers, pooling layers

edges → combination of → objects.
edges / shape

Building blocks:

- Convolutional Layers: Apply filters to extract features from input data
- Padding: Adjust spatial dimensions of the output volume to preserve resolution.
- Batch normalization: Normalize inputs of each layer to improve training speed & stability

Date: _____

Sun Mon Tue Wed Thu Fri Sat

- Stride: Determine the step size at which filters convolve.
- Pooling layer: Reduce spatial dimensions of the inputs while retaining information.
- Dropout layer: Randomly deactivates neurons during training to prevent overfitting.
- Fully connected layers: Connect every neuron in one layer to every neuron in next layer. for classification/regression
- Activation function: Introduce non-linearity to enable the network to learn complex features

Convolution:

- feature extraction
 - Sliding filter \rightarrow kernel.
 - ~~filter requires values~~ 1
 - Summation (filter value \times pixel value) \rightarrow one value of feature map.
-
- Different filter size = different levels of abstractions.
-
- filters start with random values,
 - Adjusted using backpropagation to reach optimal filters

for classification using deep Network:

- feature maps is flattened & passed to NN.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Spatial Exploitation based CNNs:

- Small filter \rightarrow fine grained
- Large sized filter \rightarrow coarse grained.

LeNet: 1998:

- first CNN $\xrightarrow{\text{traditional NN}}$
 - Considered each pixel as an input which increase the computation
 - Neighboring filter pixels are correlated.
 - LeNet used two to convolution filters
 - filters shared the learning across the image
 - Didn't reduce the computations but learned features from raw images
 - 5 layers 60K to 70K parameters 5×5 filters
- Conv - Pool - Conv - Pool - flatten - Fc - Fc

Date: _____

Sun Mon Tue Wed Thu Fri Sat

AlexNet (2012)

- first deep CNN
- was overfitting cuz of more layers so:
 - 1- ReLU was used to improve convergence by removing vanishing gradient.
 - 2- Some transformational units were skipped Randomly
 - 3- overlapping subsampling by local response normalization

5x5 filters 8 layers 60-70 million parameter.

Convolution - max pool - Norm - Convolution - max pool - Norm - convolution - convolution - max pool - fc - fc - fc

Input = $227 \times 227 \times 3$

$1 + \frac{\text{input size} - \text{kern size}}{\text{Stride}} + 2 \text{ (padding)}$

Conv L1: received 96 units 11x11 filter

Stride = 4

$$\frac{227 - 11}{4} + 1 = 55$$

$$55 \times 55 \times 96$$

$$\text{Parameters: } (11 \times 11 \times 3) \times 96 = 35K$$

Pooling layer preserves feature map depth which means depth grows at every layer.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Zifret: 2013:

- understand δ represent the performance of CNN
- we start with output of a convol layer instead of an input image.
- ~~we apply~~ ~~filtering~~
- The reverse order of convol pooling layer helps to convert an output to an image. This helps us understand which parts of input image are responsible for neuron activation.
- Experiments on AlexNet showed neurons in 1st & 2nd layers were ~~dead~~ dead
- They found distortions δ errors in 2nd layer output.
- They minimized the learning δ .

6.7-70 million parameters

Reduces number channels in bottleneck layer

Date:

Sun Mon Tue Wed Thu Fri Sat

VGG: 2014:

- Showed the relation of depth with representation capacity of NN
- Used small filter size
- Small filter mean low computational complexity
- Regulates complexity by placing 1×1 conv filter between layers
- Max pooling by padding

~ 138 million parameters - ~ 16 layers - 3×3

143 million parameters - 19 layers - 3×3

GoogleNet 2014:

- AKA Inception V1
- Reduce computation by improve accuracy.
- Introduced 'inception block'. Split, transform & merge.
- Different size filters. | Parallel convolution paths.
- Learns diverse types of variations present in same category of images.
- Regulates using 1×1 conv layer.
- Sparse connections.
- Omit feature maps that were irrelevant

Date:

Sun Mon Tue Wed Thu Fri Sat

- Global avg pooling in last layer
- Batch Normalization, RMSProp optimizer.
- 8m No fc layers
- Introduced auxiliary layers.
- Needs to be customized from module to module.
- Due to bottleneck useful information is also lost
-
- 5 million parameters.

Summary:

Parameter	Layers	Filter	Stride
→ Lenet: 60K-70K p	5(2c,3fc)	5x5	1
→ AlexNet: 60M-70M	8(5c,3fc)	1- 11x11 2- 5x5 346- 3x3	4 c1 1 c2,3,4,5
→ Zfnet: 60M-70M	8(5c,3fc)	7x7 3x3	2 c1 1 c2,3,4,5
→ VGG 16	138M	16 (13c,3fc)	3x3
→ VGG 19	145M	19 (16c,3fc)	3x3
→ GoogleNet	5M	22(20,2fc)	1x1 3x3 5x5

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Depth Based CNNs:

- Emphasize that increasing the depth of network will capture hierarchical representation of input data
- Deeper network can learn more complex features
- > Universal approximation theorem states a single layer with infinite neurons can approximate any function.
- Deeper network can maintain expressive power at reduced cost
- Problems:
 - Slow training
 - convergence speed
- = *

Highway Networks: 2015

- gating mechanisms to control ~~pass~~ information flow.
- Multipath CNN, cross layer connectivity
- 50 layers
- Normal CNN performance decrease even with 10 hidden layer
- Highway Network performed significantly better even with 900 layers.
- Gate control flow of information, when to pass info by when to carry.

Date:

Sun Mon Tue Wed Thu Fri Sat

forward propagation: gates determine how much information to flow through. gates are implemented as sigmoid functions. (output only 1)

Backpropagation: gradients of loss of all parameters are calculated including gate weights.

- big Network learns to adjust gate weights such that it effectively controls the flow of information.

→ Transform Gate: Decides how much of the input should be transformed. Allows a larger portion of the input to pass when input data contain complex features

→ Carry Gate: Determines how of input should be carried over to next layer without any transformation. when input features are already very informative or data contain important features that should be preserved

- Both gates are learned during training

Date: _____

Sun Mon Tue Wed Thu Fri Sat

ResNet: 2015

- Residual Learning
- 152 layers and still less computational complexity than any previous Network.
- 50/101/152 layer has less error than 34 layer plain Net.

Inception V3, V4 by Inception ResNet

- Improved versions of Inception V1, V2.
- V3 → Reduce computational cost
 - used 1×7 and 1×5 filter by 1×1 convolution for bottleneck.
- Combined Residual learning by inception blocks filter concatenation replaced with residual connection.
- Inception Resnet has more width by depth than normal Inception V4.
- Inception ResNet converges faster
- Residual connection accelerates training.

Auxiliary classifier: Added classifier in the middle of network to improve performance

Date: _____

Stem: initial operations

Sun Mon Tue Wed Thu Fri Sat

→ GoogleNet: 2014 22 layers 5 million parameters high comp

Inception v1: - Modules with multiple filter sizes.

→ Inception V2: 2015 42 layers 11 million parameters Mod-high comp

- Batch Normalization

- factorized 7×7 filters.

→ Inception V3: 2015 48 layers 22 million parameters Mod-high comp

- Auxiliary classifiers

- Improved inception modules

- factorized 7×7 filters.

→ Inception V4: 2016 76 layers 42 million parameters high comp

- Stem with 3×3 convolutions.

- improved architecture

→ Inception ResNet V2: 2016 164 layers 55 million parameters high comp

- Utilizes residual connection

- improved architecture

* Residual blocks decrease spatial dimension of feature maps

- This helps to extract higher level abstract features

- condense feature maps

- Down sample to manage computation complexity.

- Pooling, 1x1 conv, big stride

Date:

Sun Mon Tue Wed Thu Fri Sat

Multipath Based CNN:

- Multipath or shortcut connections connects one layer to next another by skipping some intermediate layers to allow the specialized flow of information across layers
- Solves vanishing gradient problems
- Zero padded, projection based, dropout, skip connections 1x1 connections
- Multiple path with different receptive fields

Highway networks:

- gating mechanism
- allows training many layers (even 900 layer)

when gate is open input is passed to next layers

when gate is closed then input is assigned one next layer (skipped layer)

ResNet:

- Residual connection:
- Adds input to output of layer

Depth problems: Time intensive, feature reuse,
vanishing/exploding gradient

fraction of improvement require many layers

Date:

Sun Mon Tue Wed Thu Fri Sat

DenseNet 2017:

- DenseNet connectivity between layers
- Each layer connected to every other layer in feed-forward
- Concatenates feature maps across

Width based Multicongnition CNN:

- width is also essential parameter
- width defines the learning power of NN.
- Deep Networks some units weren't learning so they shifted focus to wide networks
- Increase number of channels of extra layer

Wide ResNet 2016:

- Some blocks in deep Networks were not contributing enough
- introduced additional factor k to increase width
- Twice the number of parameters than ResNet
But can be better than deeper networks
- dropout between layers rather than in Residual blocks
- 50 layer

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Pyramidal Net 2017:

- In deep networks the depth of feature maps increase while the spatial dimensions decrease cuz of pooling layers. (down sampling)
- loss in spatial dimensions limits the learning ability
-
- In Pyramidal net width is increased gradually ^{per} ~~across~~ associated residual unit
- enables to cover all possible locations.
- uses 2 types of widening:
 - 1- additive: increases linearly
 - 2- multiplicative: increases geometrically.
- Pyramidal pooling layers: feature maps pooled at multiple scales to capture global by local information
- Multiscale feature fusion: fuse features from different scales using skip connections.
- Hierarchical architectures: progressively learns features at different levels of abstractions
low level \rightarrow more abstract
- Shared parameters by parallel processing so more efficient

Date:

Sun Mon Tue Wed Thu Fri Sat

XceptionNet 2017:

- extreme inception architecture
- Modified original Inception block by making it wider
- replacing different scales by single 3×3 scale by 1×1 convolutions
- Replaces standard convolution with depth wise convolution
 - depthwise convolution: operate on each input channel separately applying single feature filter per channel.
 - point-wise: combines the outputs of depth wise allowing for cross channel interaction.
- More computationally efficient
- No intermediate pooling layers, directly connects input to each convolution operation

ResNeXt: 2017: Aggregated Residual Transformation Network

→ Number of paths.

- cardinality: additional dimension, size of the set transformations
- used: deep homogeneous topology of VGG
 - Simplified GoogleNet architecture (split, transform, merge)
 - Residual learning

feature map wise statistic / fm motifs / fm descriptor / channel wise descriptor
all come (output of Squeeze operation) Date:

Sun Mon Tue Wed Thu Fri Sat

Feature Map Exploitation Based CNN:

- Some feature maps import little-to-no role in object discrimination.
- Some feature sets may create noise which leads to overfitting
- Selection of feature maps can help in improving generalization
- Compress or enhance feature maps.

Squeeze by excitation: 2018:

- Squeeze operation: reduce the spatial dimension of the input feature map to a single spatial point (global avg pool)
Spatial information \rightarrow channel wise descriptor.
- Excite operation: learns a set of weights that capture importance of each feature map. channel.
 - learn by using a small fully connected layer
 - The weights are then applied to each feature map channel to modulate its importance

feature Re-calibrate: By performing squeeze by excite ops SENet recalibrate feature maps in a channel wise manner.

- This allows it to focus on informative features while suppressing less relevant ones

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Integration with CNN: can be added as a lightweight module after convolayers

- Act as attention mechanisms

- Suppress less important features by giving higher weightage to class specifying feature maps
- Can be added to any CNN before convolutional layer.
- gets global view of feature maps

Competitive Squeeze by Excitation 2018:

- In ResNet it only considers the residual information for determining weights of each feature map
- Addressed this problem by generating feature maps-wise motifs (statistics) from both residual and identity mapping based feature maps
- global representation of fm generated using global average pooling

Inner Imaging:

- relevance between feature maps estimated by establishing a competition between feature descriptors of residual by identity mapping.
- Not only models the relation between residual feature maps but also maps their relation with the identity feature maps

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Channel exploitation based CNN:

distinguishing features

- The lack of diversity of class discernable information in the input affect CNN's performance
- Introduced channel boosting using auxiliary learners.
- exploits channel wise relationships in feature maps
- attention mechanisms to adaptively recalibrate feature maps and boost informative content.

Channel Boosted CNN 2018:

- proposed idea of boosting number of input channels for improving the ~~capacity~~ representational capacity of the network
- artificially create extra channels AKA auxiliary channels

$$I_B = g_k(I_c[A_1, \dots, A_m])$$

$$f_i^k = g_c(I_B, k_i)$$

I_c → original inputs

A_m → artificial inputs

g_k → combiner function

I_B → Boosted input

F_i^k → feature map generated by convolving I_B with filter k

- for improving representation of data power of TL by deep generative teacher learners were exploited.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

- Generative learners learns patterns of input data to create similar data
- Autoencoders are used as generative learners
- Transfer learners used to create encoded input representation
 - knowledge gained from one task is used to improve another task
- CB-CNN encodes the channel boosting phase into generic block
 - inserted at start of deep Network
- Suggested that auxiliary channels can be inserted at any layer in deep architecture.
boosting block

* Attention Based CNNs:

- In addition to learning multiple hierarchies of abstractions focusing on features relevant to the context also plays significant role in image localization.
- view scenes in succession of partial glimpse and pay attention to context relevant parts
- RNN by LSTMs exploit attention modules for generation of sequential data

- Idea of attention also help CNN to recognize objects even from cluttered ^{background} ~~Received~~ by complex scenes
- let CNN incorporate attention mechanisms to capture long range dependencies by focus on relevant regions in image +
- They also improve network ability to attend to informative parts of the input

Residual Attention Neural Network: 2017:

- proposed to improve feature representation of the network
- ~~Make networks object-aware~~
- Make network capable of learning object-aware features
- Built by stacking residual blocks with attention module.
- Branched into trunk by mask branches.
adopts Topdown by bottom up learning

Bottom up: feed forward structure produces low resolution feature maps with strong semantic information

Topdown: ~~parallel~~ globally optimizes the network in such a way that it gradually outputs the maps to input & produces dense feature to make inference of each pixel.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

$$g_{am}(F_L^k) = g_{cm}(f_L^k) \cdot g_{tm}(f_L^k)$$

$g_{cm} \rightarrow$ object aware softmax for feature map at layer L

- assigns attention towards objects

$\overrightarrow{g_{tm}}$

- RAN can be made efficient by for cluttered, complex & noisy images by adding attention modules.

- Adaptively assign weight to each feature map based on their relevance in the layer.
- Learning of deep hierarchical structure was supported by residual units
- Mixed, channel, by spatial attentions were incorporated

Convolutional Block Attention 2018:

- SENet only considers the contribution of feature maps in image classification, but ignores the spatial locality of the object in images
- CBAM - infers attention maps sequentially by first applying feature map(channels) attention and then spatial attention to find the refined feature maps

Date: _____

Sun Mon Tue Wed Thu Fri Sat

- generally 1×1 convolution by pooling layer are used for spatial attention.
- CBAM concatenates avg pooling operation with max pooling which generate strong spatial attention map.
- feature map statistics were modelled using combination of max pooling and global max average pooling
- Max pooling could provide clues about distinctive object features
- Global average pooling returns suboptimal inference of feature map attention.
- The exploitation of both pooling improves representational power of network. They not only focus on the important part but also increase the representational power of the selected feature maps
- formulation of 3D attention map via serial learning reduces parameters and computation cost
- CBAM can be inserted in CNN.

Date: _____

Sun Mon Tue Wed Thu Fri Sat

Concurrent Squeeze by Excitation 2018:

- Extended SENet by incorporating the effect of spatial information in the combination with feature maps (channels) information
- They introduced 3 modules:
 - CSE • Squeezing partially by existing feature maps information
 - SSE • Squeezing feature map by exciting spatial information
 - ScSSE • concurrent squeeze by excitation of spatial by fm info
- AE based convolutional NN was used for segmentation
- modules were inserted after the encoders and decoders
-

CSE: • similar concept as SE Block
• Scaling factor is derived based on combination of fm

SSE: • spatial locality is given more importance
• different combination of feature maps are selected by exploited Spatially to use for segmentation.