

Cheat Sheet – Convolutional Neural Network

Convolutional Neural Network:

The data gets into the CNN through the input layer and passes through various hidden layers before getting to the output layer. The output of the network is compared to the actual labels in terms of loss or error. The partial derivatives of this loss w.r.t the trainable weights are calculated, and the weights are updated through one of the various methods using backpropagation.

CNN Template:

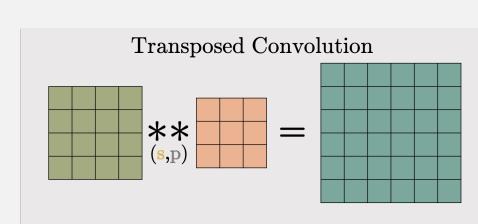
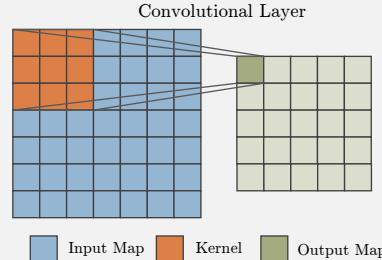
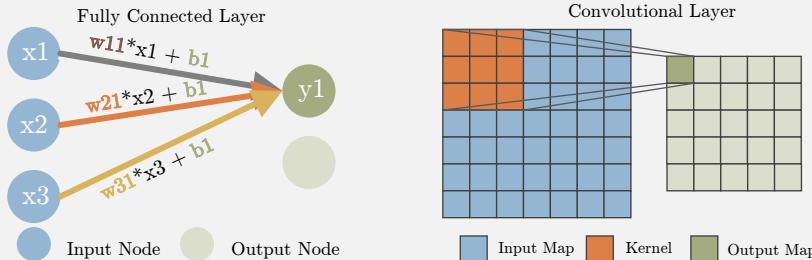
Most of the commonly used hidden layers (not all) follow a pattern

1. Layer function: Basic transforming function such as convolutional or fully connected layer.

a. Fully Connected: Linear functions between the input and the

a. Convolutional Layers: These layers are applied to 2D (3D) input feature maps. The trainable weights are a 2D (3D) kernel/filter that moves across the input feature map, generating dot products with the overlapping region of the input feature map.

b. Transposed Convolutional (DeConvolutional) Layer: Usually used to increase the size of the output feature map (Upsampling) The idea behind the transposed convolutional layer is to undo (not exactly) the convolutional layer



2. Pooling: Non-trainable layer to change the size of the feature map

a. Max/Average Pooling: Decrease the spatial size of the input layer based on selecting the maximum/average value in receptive field defined by the kernel

b. UnPooling: A non-trainable layer used to increase the spatial size of the input layer based on placing the input pixel at a certain index in the receptive field of the output defined by the kernel.

3. Normalization: Usually used just before the activation functions to limit the unbounded activation from increasing the output layer values too high

a. Local Response Normalization LRN: A **non-trainable layer** that square-normalizes the pixel values in a feature map within a local neighborhood.

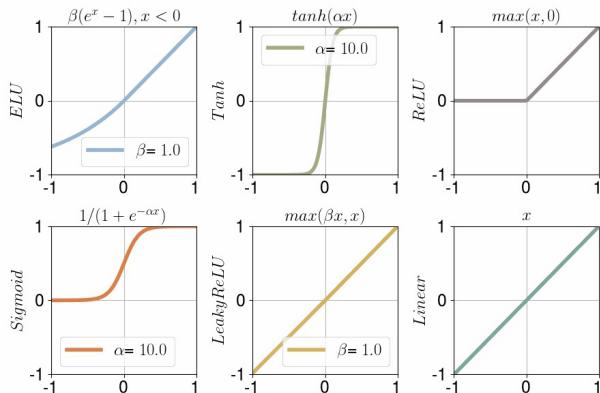
b. Batch Normalization: A trainable approach to normalizing the data by learning scale and shift variable during training.

3. Activation: Introduce non-linearity so CNN can efficiently map non-linear complex mapping.

a. Non-parametric/Static functions: Linear, ReLU

b. Parametric functions: ELU, tanh, sigmoid, Leaky ReLU

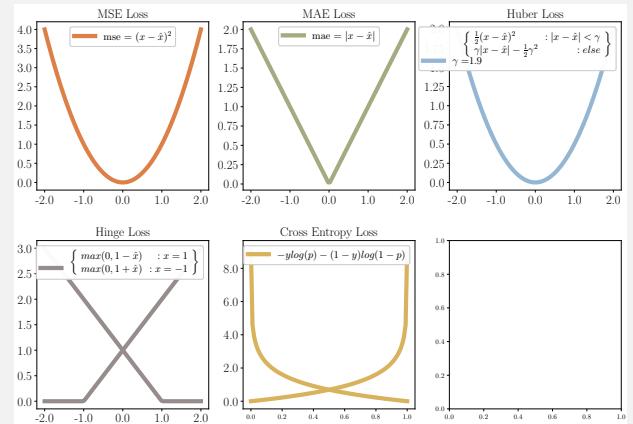
c. Bounded functions: tanh, sigmoid



5. Loss function: Quantifies how far off the CNN prediction is from the actual labels.

a. Regression Loss Functions: MAE, MSE, Huber loss

b. Classification Loss Functions: Cross entropy, Hinge loss



Cheat Sheet – Famous CNNs

AlexNet – 2012

Why: AlexNet was born out of the need to improve the results of the ImageNet challenge.

What: The network consists of 5 Convolutional (CONV) layers and 3 Fully Connected (FC) layers. The activation used is the Rectified Linear Unit (ReLU).

How: Data augmentation is carried out to reduce over-fitting, Uses Local response normalization.

VGGNet – 2014

Why: VGGNet was born out of the need to reduce the # of parameters in the CONV layers and improve on training time

What: There are multiple variants of VGGNet (VGG16, VGG19, etc.)

How: The important point to note here is that all the conv kernels are of size 3x3 and maxpool kernels are of size 2x2 with a stride of two.

ResNet – 2015

Why: Neural Networks are notorious for not being able to find a simpler mapping when it exists. ResNet solves that.

What: There are multiple versions of ResNetXX architectures where 'XX' denotes the number of layers. The most used ones are ResNet50 and ResNet101. Since the vanishing gradient problem was taken care of (more about it in the How part), CNN started to get deeper and deeper

How: ResNet architecture makes use of shortcut connections to solve the vanishing gradient problem. The basic building block of ResNet is a Residual block that is repeated throughout the network.

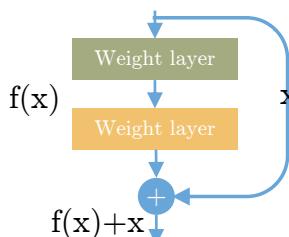


Figure 1 ResNet Block

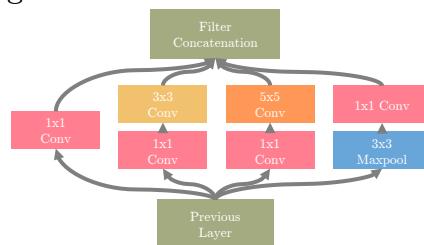


Figure 2 Inception Block

Inception – 2014

Why: Larger kernels are preferred for more global features, on the other hand, smaller kernels provide good results in detecting area-specific features. For effective recognition of such a variable-sized feature, we need kernels of different sizes. That is what Inception does.

What: The Inception network architecture consists of several inception modules of the following structure. Each inception module consists of four operations in parallel, 1x1 conv layer, 3x3 conv layer, 5x5 conv layer, max pooling

How: Inception increases the network space from which the best network is to be chosen via training. Each inception module can capture salient features at different levels.

AlexNet Network - Structural Details										
Input	Output	Layer	Stride	Pad	Kernel size	in	out	# of Param		
227x227x3	55x55x96	conv1	4	0	11	11	3	96	34944	
55x55x96	27x27x96	maxpool1	2	0	3	96	2	96	614656	
27x27x96	27x27x256	conv2	1	2	5	96	256	128	614656	
27x27x256	13x13x256	maxpool2	2	0	3	256	256	0		
13x13x256	13x13x384	conv3	1	3	3	256	384	885120		
13x13x384	13x13x384	conv4	1	1	3	384	384	1327488		
13x13x384	13x13x256	conv5	1	1	3	384	256	884992		
13x13x256	6x6x256	maxpool5	2	0	3	256	256	0		
		f6				1	1	9216	4096	37752832
		f7				1	1	4096	4096	16781312
		f8				1	1	4096	1000	4097000
		Total								62378.344

VGG16 - Structural Details												
#	Input	Image	output	Layer	Stride	Kernel	in	out	Param			
1	224	224	3	224x224	64	conv3-64	1	3	3	64	1792	
2	224	224	64	224x224	64	conv3064	1	3	3	64	36928	
3	112	112	64	112x112	128	maxpool	2	2	2	64	64	
4	112	112	128	112x112	128	conv3-128	1	3	3	64	128	
5	112	112	64	112x112	128	maxpool	2	2	2	128	128	
6	56	56	128	56x56	256	conv3-256	1	3	3	128	256	
7	56	56	256	56x56	256	conv3-256	1	3	3	256	509080	
8	28	28	256	28x28	512	maxpool	2	2	2	256	256	
9	28	28	512	28x28	512	conv3-512	1	3	3	512	512	
10	28	28	512	28x28	512	conv3-512	1	3	3	512	2359808	
11	14	14	512	14x14	512	maxpool	2	2	2	512	512	
12	14	14	512	14x14	512	conv3-512	1	3	3	512	2359808	
13	14	14	512	14x14	512	conv3-512	1	3	3	512	512	
14	14	512	7	512	512	maxpool	2	2	2	512	0	
15	1	1	25088	1	1	4096	fc	1	1	25088	4096	102764544
16	1	1	4096	1	1	4096	fc	1	1	4096	4096	16781312
		Total									138,423,208	

ResNet18 - Structural Details										Param		
#	Input	Image	output	Layer	Stride	Pad	Kernel	in	out	Param		
1	227x227	3	112x112	64	conv1	2	1	7	7	64	9472	
2	112x112	64	56	56	maxpool	2	0	5	3	64	0	
3	56	56	64	56	64	conv2-1	1	1	3	64	64	
4	56	56	64	56	64	conv2-2	1	1	3	64	64	
5	56	56	64	56	64	conv2-4	1	1	3	64	64	
6	56	56	64	56	64	conv2-8	1	1	3	128	128	
7	28	28	128	28x28	128	conv3-2	1	1	3	128	147584	
8	28	28	128	28x28	128	conv3-4	1	1	3	128	147584	
9	28	28	128	28x28	128	conv3-8	1	1	3	128	147584	
10	28	28	128	28x28	128	conv3-16	1	1	3	256	256	
11	14	14	256	14x14	256	conv3-32	1	1	3	256	256	
12	14	14	512	14x14	512	conv3-512	1	1	3	512	512	
13	14	14	512	14x14	512	conv3-512	1	1	3	512	512	
14	14	512	7	512	512	maxpool	2	2	2	512	0	
15	1	1	25088	1	1	4096	fc	1	1	25088	4096	102764544
16	1	1	4096	1	1	4096	fc	1	1	4096	4096	16781312
		Total									11,511,784	

GoogLeNet - Structural Details										Param	
#	Input	Image	output	Layer	Input Layer	Stride	Pad	Kernel	in	out	Param
1	227x227x3	3	110x110x64	conv1	conv1	2	0	3	64	64	0
2	110x110x64	64	56x56x64	maxpool	conv1	2	0	3	64	64	0
3	56x56x64	64	56x56x192	conv1x1	maxpool	1	0	1	1	64	1480
4	56x56x192	192	28x28x192	conv2-1	conv1x1	2	0	3	192	192	118,924
5	28x28x192	192	28x28x384	conv2-2	conv2-1	2	0	3	192	192	118,924
6	28x28x384	384	28x28x384	conv2-3	conv2-2	2	0	3	192	192	118,924
7	28x28x384	384	28x28x384	conv2-4	conv2-3	2	0	3	192	192	118,924
8	28x28x384	384	28x28x384	conv2-5	conv2-4	2	0	3	192	192	118,924
9	28x28x384	384	28x28x384	conv2-6	conv2-5	2	0	3	192	192	118,924
10	28x28x384	384	28x28x384	conv2-7	conv2-6	2	0	3	192	192	118,924
11	28x28x384	384	28x28x384	conv2-8	conv2-7	2	0	3	192	192	118,924
12	28x28x384	384	28x28x384	conv2-9	conv2-8	2	0	3	192	192	118,924
13	28x28x384	384	28x28x384	conv2-10	conv2-9	2	0	3	192	192	118,924
14	28x28x384	384	28x28x384	conv2-11	conv2-10	2	0	3	192	192	118,924
15	28x28x384	384	28x28x384	conv2-12	conv2-11	2	0	3	192	192	118,924
16	28x28x384	384	28x28x384	conv2-13	conv2-12	2	0	3	192	192	118,924
17	28x28x384	384	28x28x384	conv2-14	conv2-13	2	0	3	192	192	118,924
18	28x28x384	384	28x28x384	conv2-15	conv2-14	2	0	3	192	192	118,924
19	28x28x384	384	28x28x384	conv2-16	conv2-15	2	0	3	192	192	118,924
20	28x28x384	384	28x28x384	conv2-17	conv2-16	2	0	3	192	192	118,924
21	28x28x384	384	28x28x384	conv2-18	conv2-17	2	0	3	192	192	118,924
22	28x28x384	384	28x28x384	conv2-19	conv2-18	2	0	3	192	192	118,924
23	28x28x384	384	28x28x384	conv2-20	conv2-19	2	0	3	192	192	118,924
24	28x28x384	384	28x28x384	conv2-21	conv2-20	2	0	3	192	192	118,924
25	28x28x384	384	28x28x384	conv2-22	conv2-21	2	0	3	192	192	118,924
26	28x28x384	384	28x28x384	conv2-23	conv2-22	2	0	3	192	192	118,924
27	28x28x384	384	28x28x384	conv2-24	conv2-23	2	0	3	192	192	118,924
28	28x28x384	384	28x28x384	conv2-25	conv2-24	2	0	3	192	192	118,924
29	28x28x384	384	28x28x384	conv2-26	conv2-25	2	0	3	192	192	118,924
30	28x28x384	384	28x28x384	conv2-27	conv2-26	2	0	3	192	192	118,924
31	28x28x384	384	28x28x384	conv2-28	conv2-27	2	0	3	192	192	118,924
32	28x28x384	384	28x28x384	conv2-29	conv2-28	2	0	3	192	192	118,924
33	28x28x384	384	28x28x384	conv2-30	conv2-29	2	0	3	192	192	118,924
34	28x28x384	384	28x28x384	conv2-31	conv2-30	2	0	3	192	192	118,924
35	28x28x384	384	28x28x384	conv2-32	conv2-31	2	0	3	192	192	118,924
36	28x28x384	384	28x28x384	conv2-33	conv2-32	2	0	3	192</td		