

What is learning?

Image Source: Internet

Can You Recognize these Pictures ?



- If Yes, How do you Recognize it?

Image Source: Internet

Descriptors / Feature Vectors





08 02 22 97 38 15 00 40 00 75 04 05 07 78 52 12 50 77 85 98
49 49 99 40 17 81 18 57 60 87 17 40 98 43 69 44 05 56 62 00
81 49 31 73 55 79 14 29 93 71 40 67 54 03 30 03 49 13 36 65
52 70 95 23 04 60 11 42 65 01 69 56 01 32 56 71 37 02 36 91
22 31 16 71 51 67 03 89 41 92 36 54 22 40 40 28 66 33 13 80
24 47 33 00 99 03 45 02 44 75 33 53 78 36 84 20 35 17 12 50
02 98 81 28 64 23 67 10 26 38 40 67 59 54 70 66 18 38 64 70
67 26 20 68 02 62 12 20 95 63 94 39 63 08 40 91 66 49 94 21
24 55 58 05 66 73 99 26 97 17 78 78 94 03 14 88 34 89 63 72
21 36 23 09 75 00 76 44 20 45 35 14 00 61 33 97 34 31 33 95
78 17 53 28 22 75 31 67 15 94 03 80 04 62 16 14 09 53 56 92
16 39 05 42 96 35 31 47 55 58 88 24 00 17 54 24 36 29 85 57
86 56 00 48 35 71 89 07 05 44 44 37 44 60 21 58 51 54 17 58
19 80 81 68 05 94 47 69 28 73 92 13 86 52 17 77 04 89 55 40
04 52 08 83 97 35 99 16 07 97 57 32 16 26 26 79 33 27 98 66
03 94 48 87 57 62 20 72 03 46 33 67 46 55 12 32 63 93 53 69
04 42 16 73 35 44 19 11 24 94 72 18 08 46 29 32 40 62 76 36
20 69 36 41 72 30 23 88 34 04 03 69 82 67 59 85 74 04 36 16
20 73 35 29 78 31 90 01 74 31 49 71 46 01 16 23 57 05 54
01 70 54 71 83 51 54 69 16 92 33 48 61 43 52 01 89 23 67 48

What the computer sees

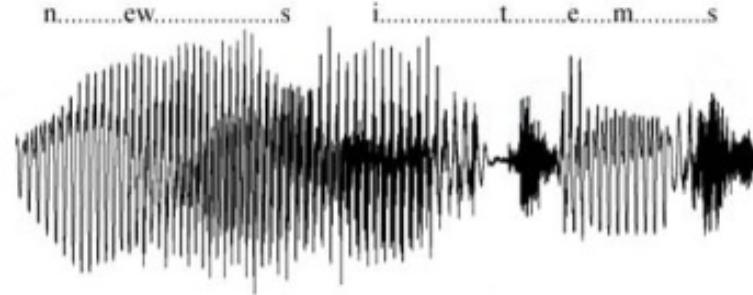


Image Source: Internet

Descriptors / Feature Vectors



Image Source: Internet

Descriptors / Feature Vectors

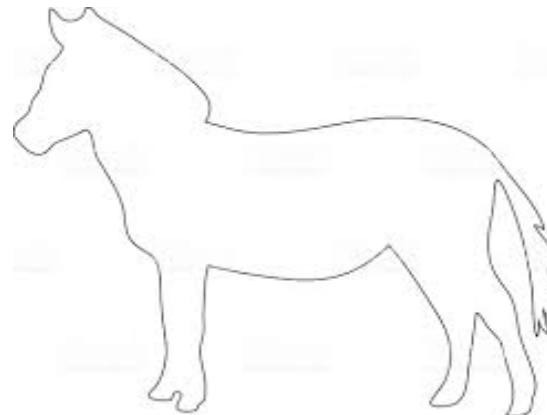
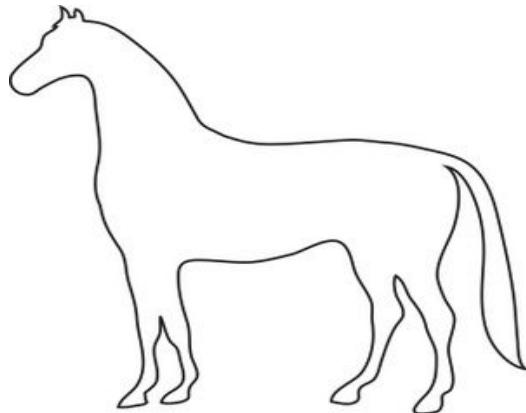


Image Source: Internet

Descriptors / Feature Vectors



Image Source: Internet

Challenges

Image Source: Internet

Viewing Angle



Image Source: Internet

Pose



Image Source: Internet

Illumination



Image Source: Internet

Intraclass Variation



Image Source: Internet

Distortion and Occlusion



Image Source: Internet

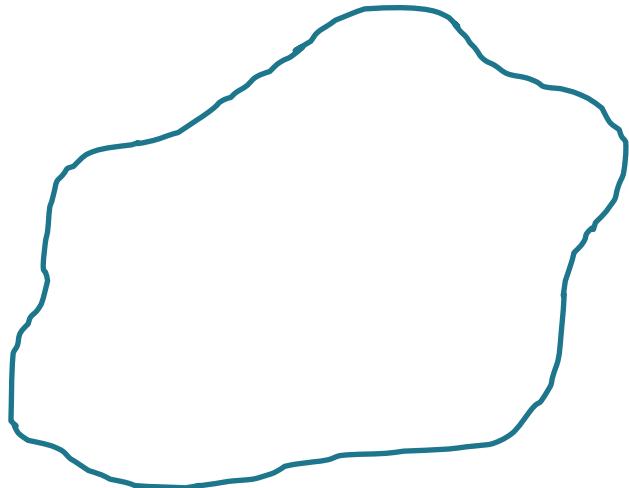
Power of Deep Learning



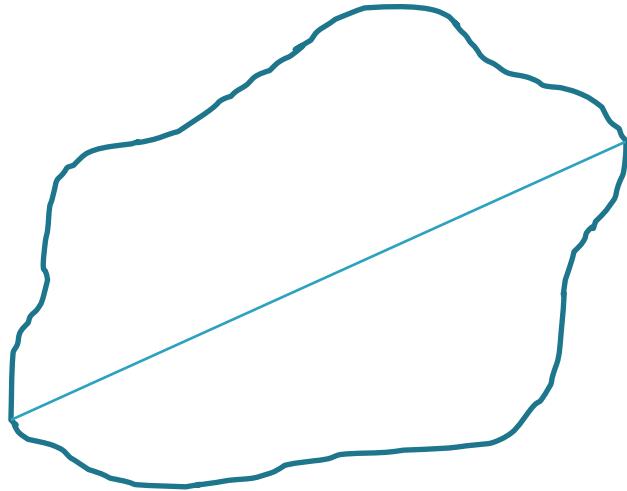
Boundary Descriptors



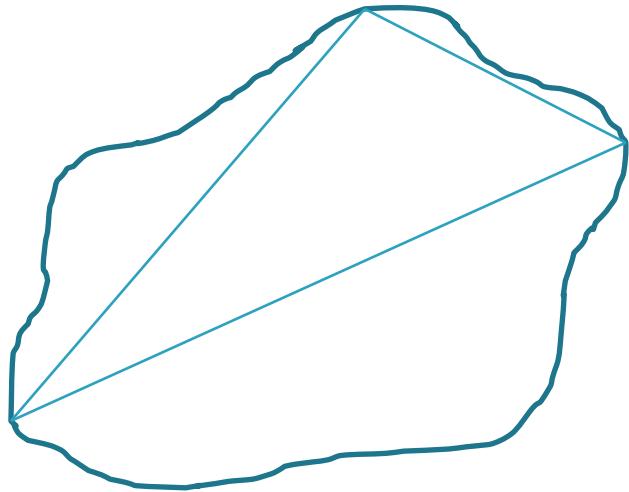
Boundary Descriptors



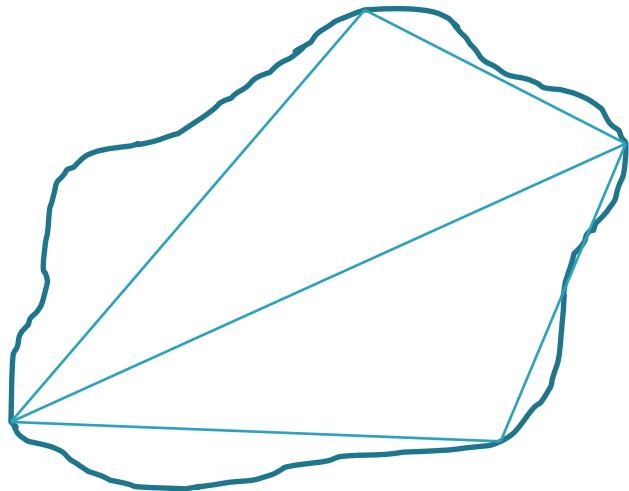
Boundary Descriptors



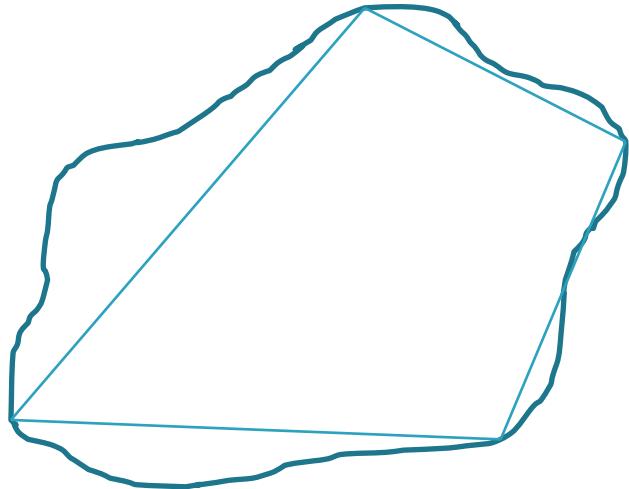
Boundary Descriptors

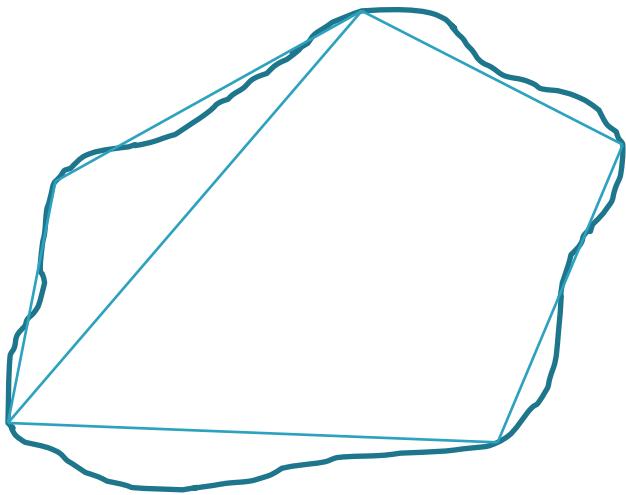


Boundary Descriptors

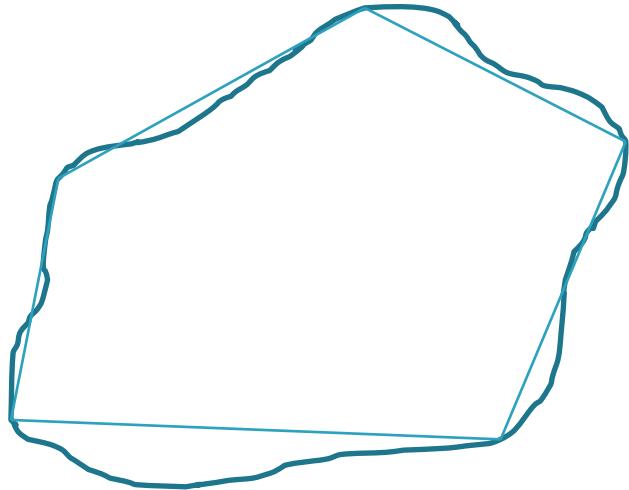


Boundary Descriptors

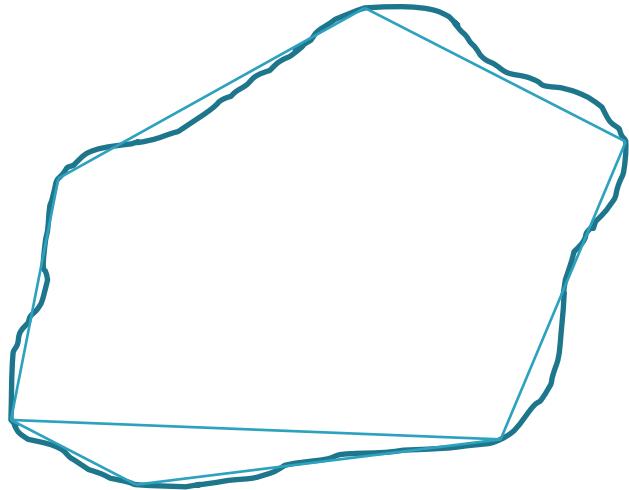




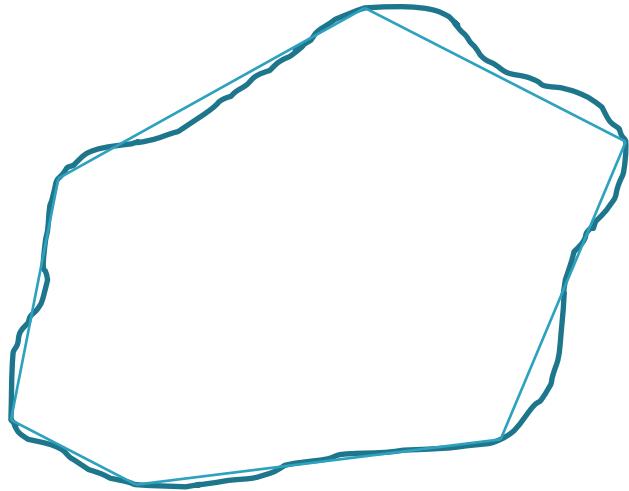
Boundary Descriptors



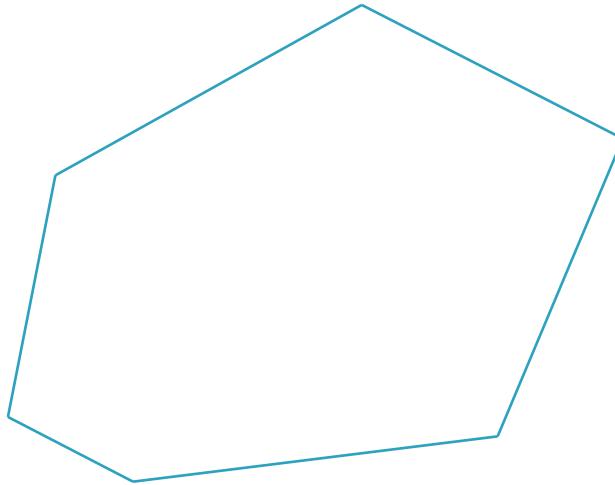
Boundary Descriptors



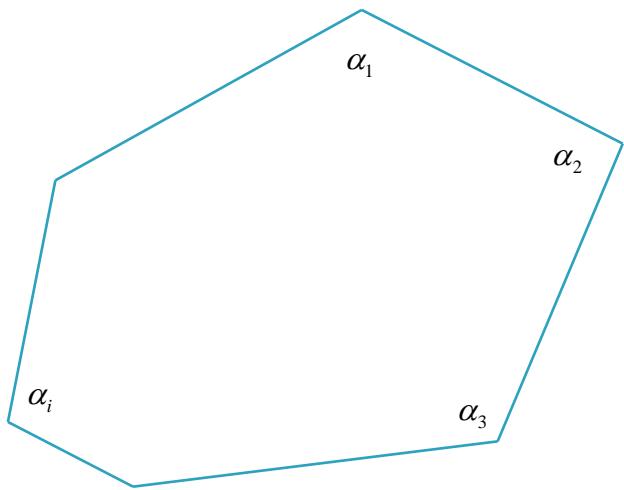
Boundary Descriptors



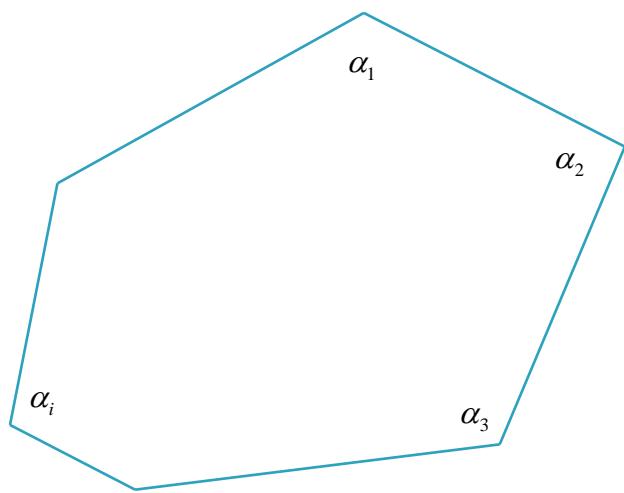
Boundary Descriptors



Boundary Descriptors



Boundary Descriptors


$$\alpha_1 \alpha_2 \alpha_3 \dots \alpha_{i-k} \dots \alpha_i \alpha_{i+1} \dots$$

$$\alpha_i = \sum_{k=1}^p \beta_k \alpha_{i-k}$$

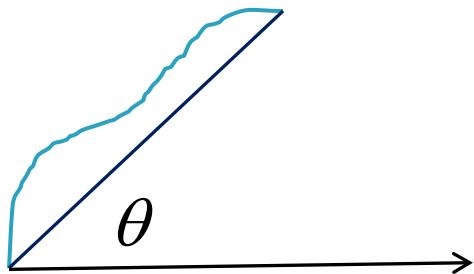
$$\beta_1 \beta_2 \beta_3 \dots \beta_p$$

Descriptors (features) of polygonal shape

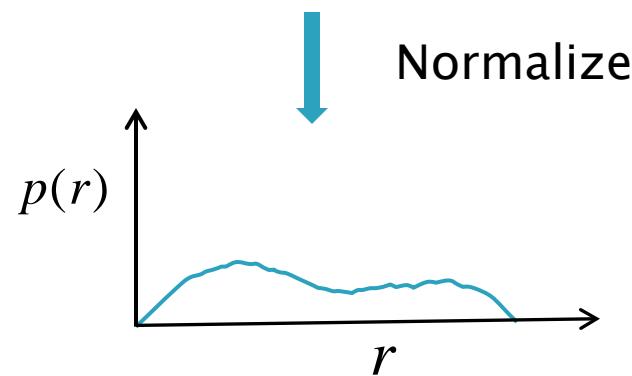
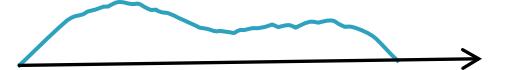
Boundary Descriptors



Boundary Descriptors



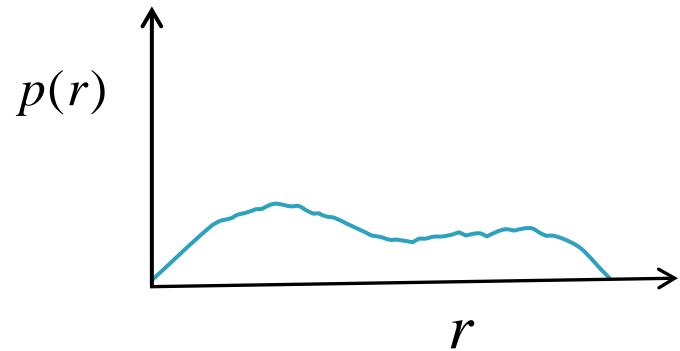
Rotate
→



Normalize
↓

Boundary Descriptors

Statistical Moments



$$\sum_{i=1}^k (r_i - \hat{r})^2 p(r_i) \quad \text{Spread}$$

$$\sum_{i=1}^k (r_i - \hat{r})^3 p(r_i) \quad \text{Skewness}$$

$$\sum_{i=1}^k (r_i - \hat{r})^4 p(r_i) \quad \text{Kurtosis}$$

Region Descriptors

Intensity Descriptor

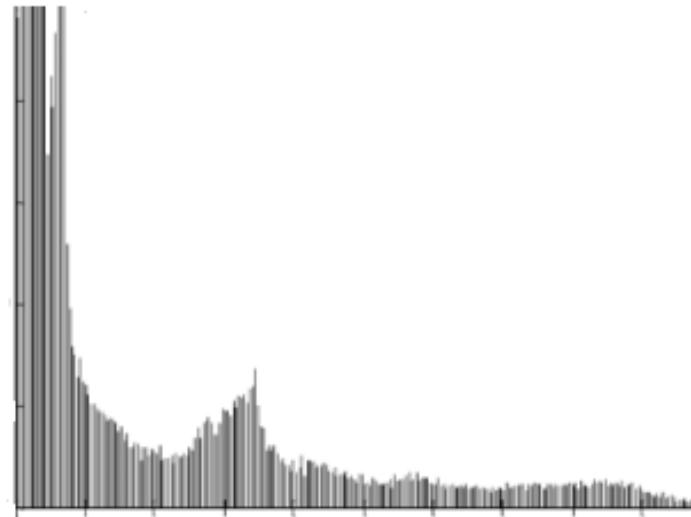


Image Source: Internet

Intensity Descriptor

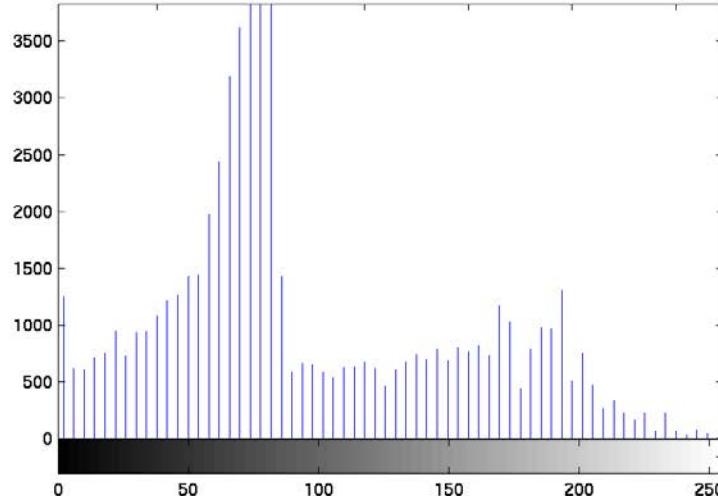


Image Source: Internet

Colour Feature

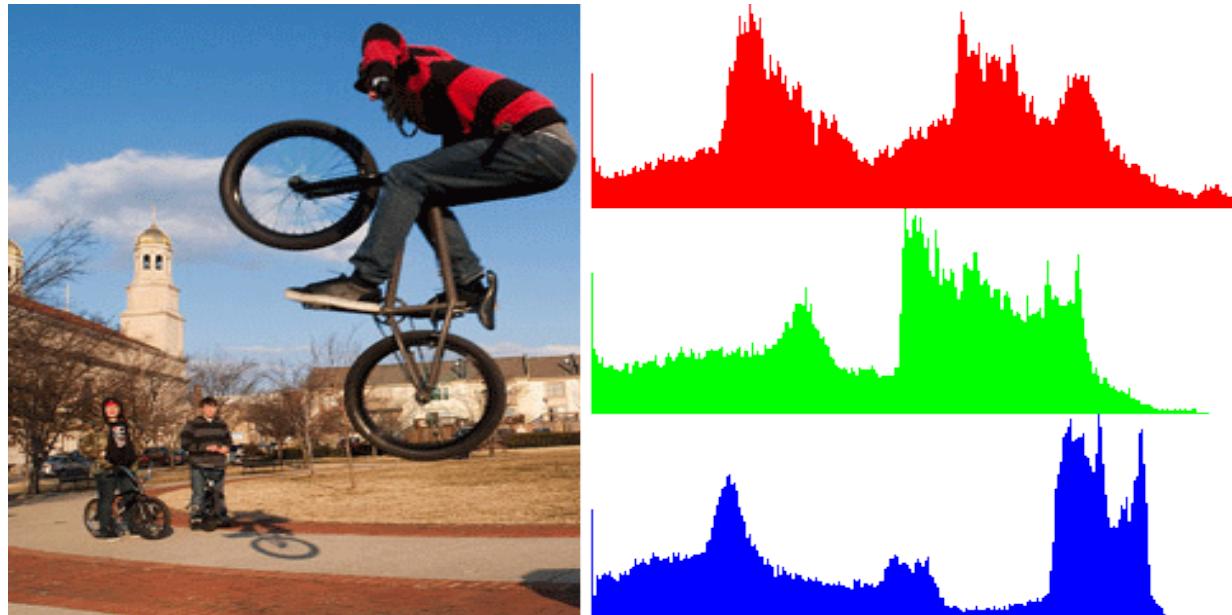
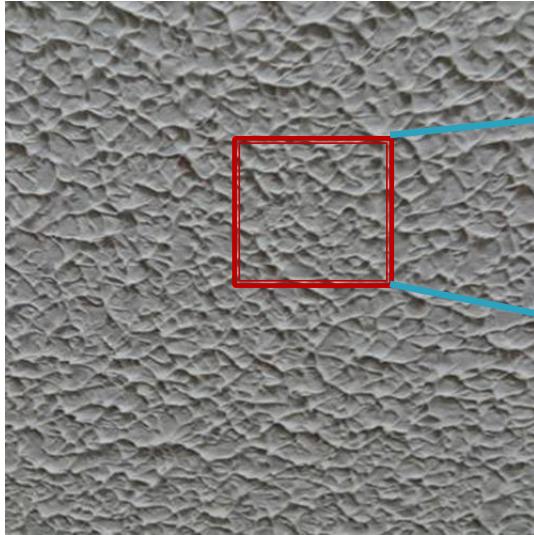


Image Source: https://billmill.org/the_histogram.html

Texture Descriptors

Pixel Domain/ Co-occurrence Matrix



150	100	115	109	112	100	145	140
110	112	120	135	125	120	132	133
152	99	129	130	122	135	98	100
147	138	142	95	108	136	110	125
99	127	149	138	138	129	108	129
128	125	139	115	120	145	137	131
146	159	150	130	147	139	143	127
140	120	128	98	100	106	115	119

Pixel Domain/ Co-occurrence Matrix

10	9	7	9	5	8	11	9
6	5	15	12	4	6	3	2
9	3	2	10	6	8	4	5
8	2	4	3	7	5	6	1
2	0	11	8	10	9	8	2
8	4	7	1	6	0	7	6
2	3	8	9	11	6	3	9
7	2	8	8	6	12	6	7

Pixel Domain/ Co-occurrence Matrix

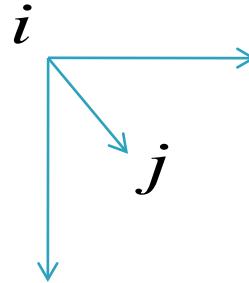
7	6	7	5	5	6	4	6
6	5	1	4	4	2	3	2
7	3	2	5	6	7	4	5
5	2	4	3	7	5	6	1
2	0	4	7	3	5	7	2
5	4	7	1	6	0	7	6
2	3	4	5	6	6	3	5
7	2	3	5	6	3	6	7

$$A_{l,\theta}(i, j)$$

Pixel Domain/ Co-occurrence Matrix

7	6	7	5	5	6	4	6
6	5	1	4	4	2	3	2
7	3	2	5	6	7	4	5
5	2	4	3	7	5	6	1
2	0	4	7	3	5	7	2
5	4	7	1	6	0	7	6
2	3	4	5	6	6	3	5
7	2	3	5	6	3	6	7

$$A_{1,45^0}(i, j)$$



0	0	0	1	0	0	0	1
0	?	?	?	?	1	?	?
?	?	?	1	1	?	?	?
?	?	?	?	?	?	?	?
?	?	?	?	?	?	0	?
?	?	?	?	?	?	?	?
?	?	?	2	?	?	?	?
?	?	?	?	?	3	?	?

Co-occurrence matrix based descriptors

Maximum Probability

$$\max_{i,j}(c_{ij})$$

Element Difference Moment

$$\sum_i \sum_j (i - j)^k C_{i,j}$$

Inverse Element Difference Moment

$$\sum_i \sum_j C_{i,j} / (i - j)^k \quad i \neq j$$

Uniformity

$$\sum_i \sum_i C_{ij}^2$$

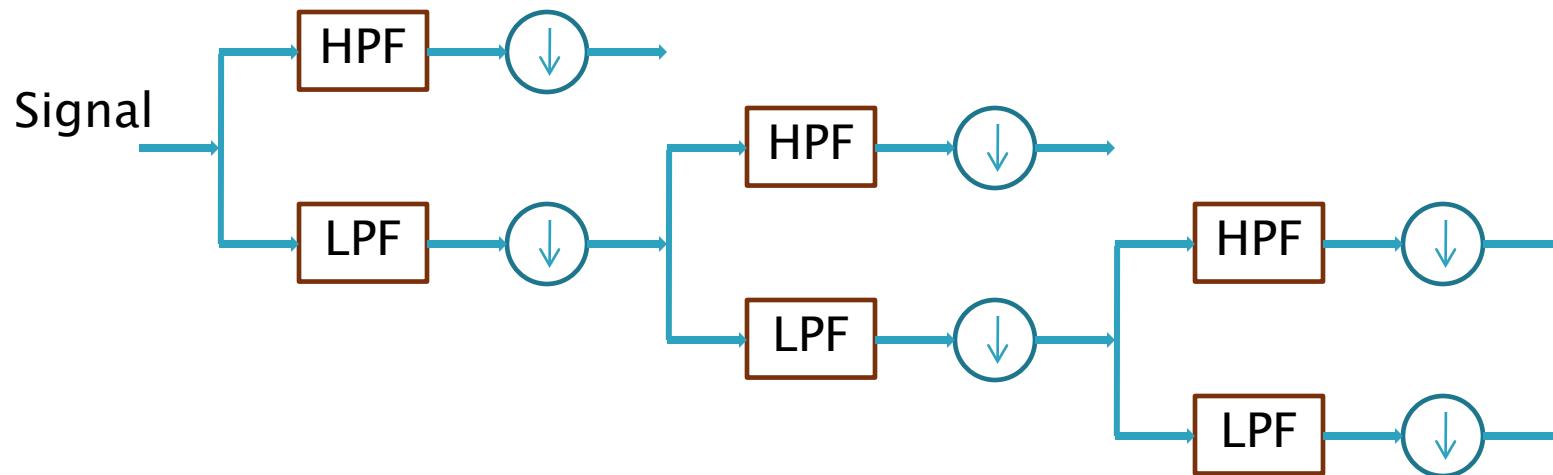
Entropy

$$-\sum_i \sum_j c_{ij} \log_2 C_{ij}$$

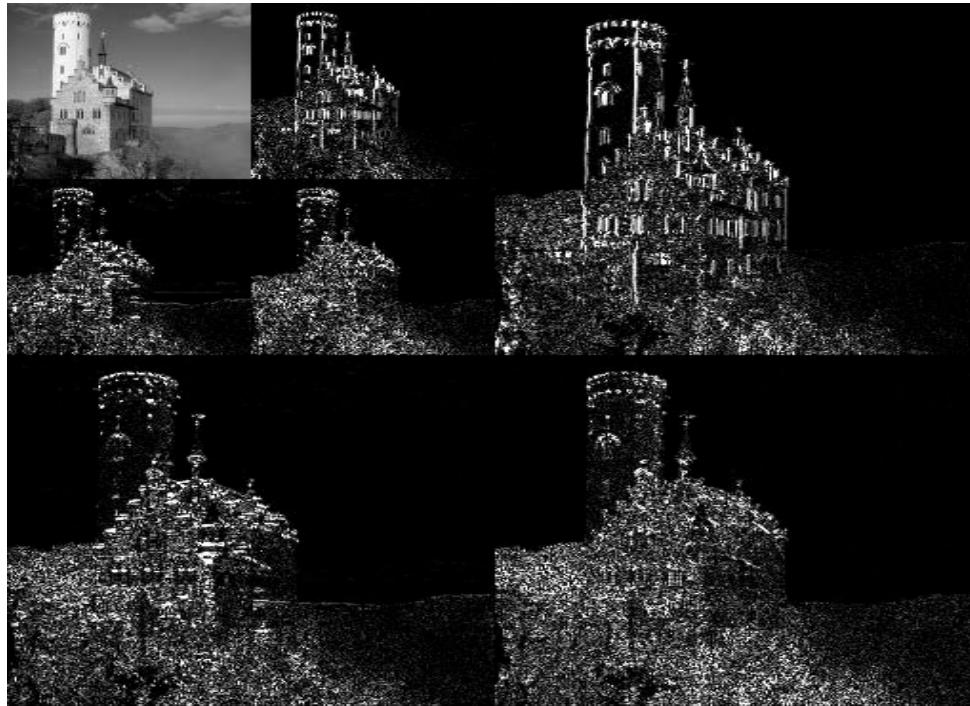
Transform Domain Features



Wavelet Transform



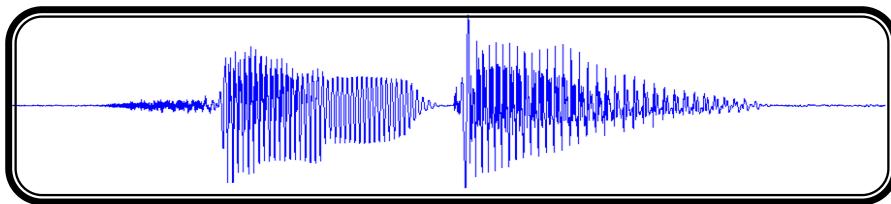
Wavelet Transform



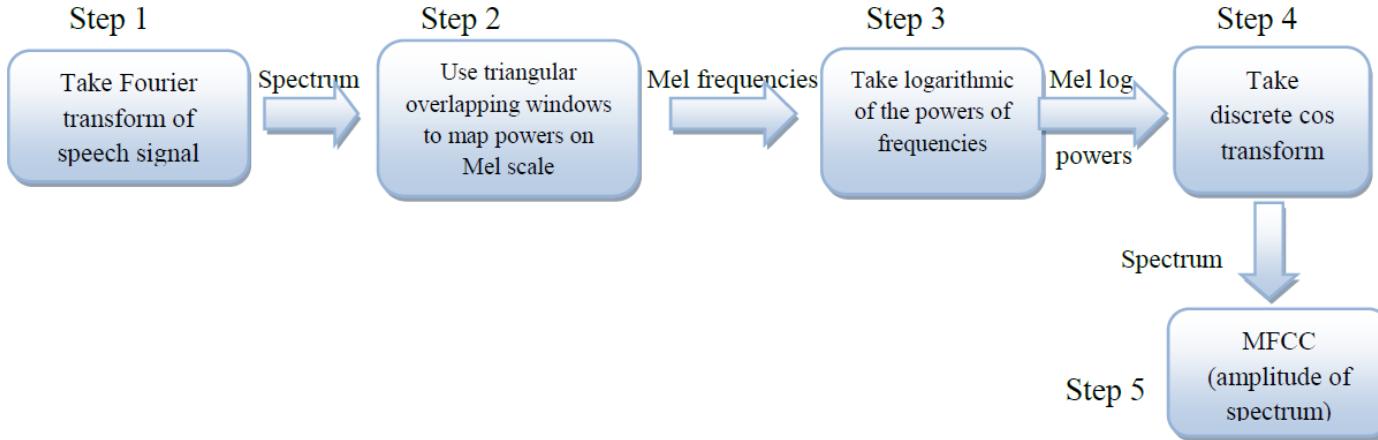
Source: Wikipedia

Audio

Time Domain Feature – LPC



Spectral Domain- MFCC

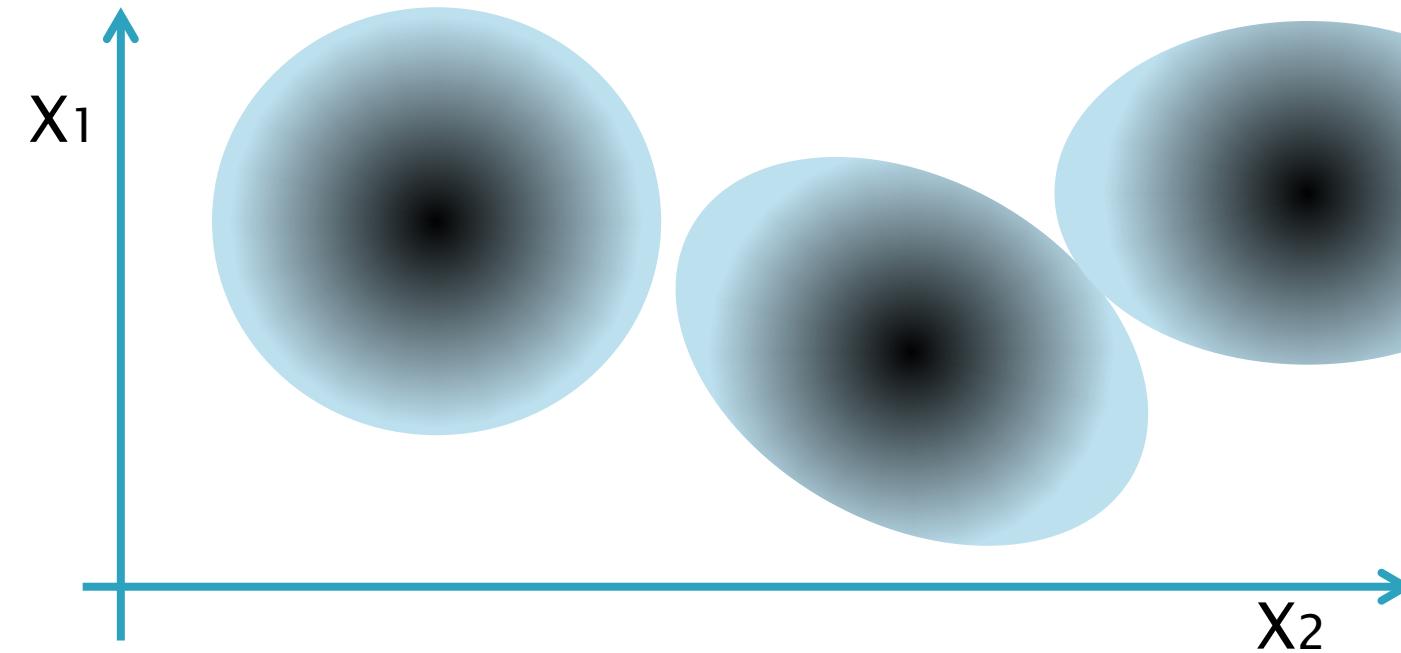


Feature Space Representation



Image Source: Internet

Feature Space Representation



Feature Space Representation

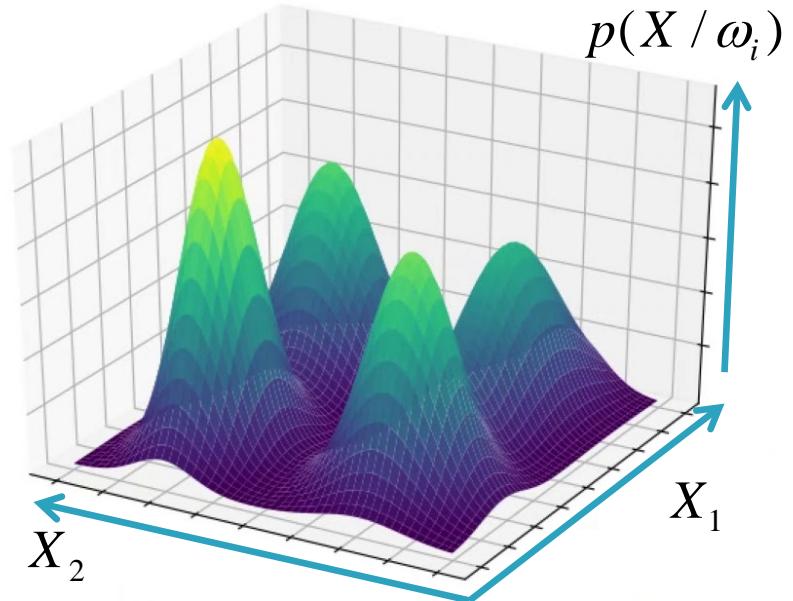


Image Source: Internet

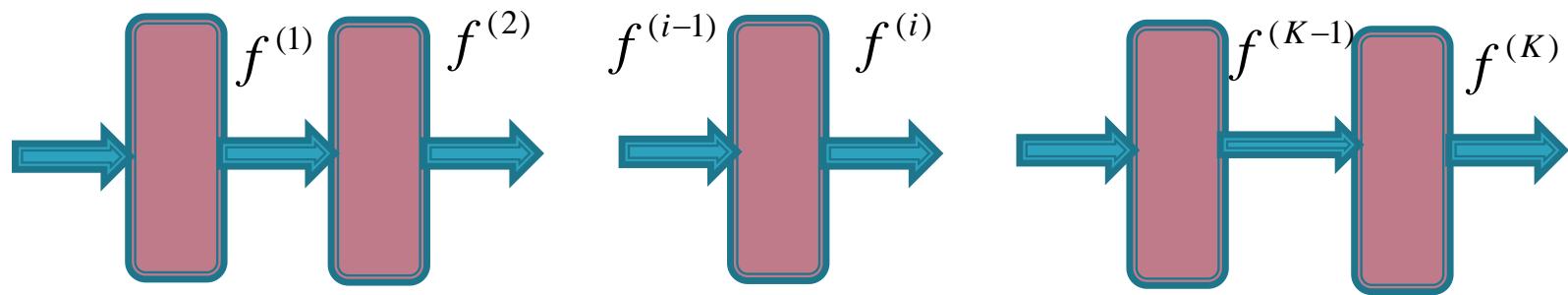
Bayesian Learning

Class conditional density $\Rightarrow p(X / \omega_i)$

For classification $\Rightarrow p(\omega_i / X)$

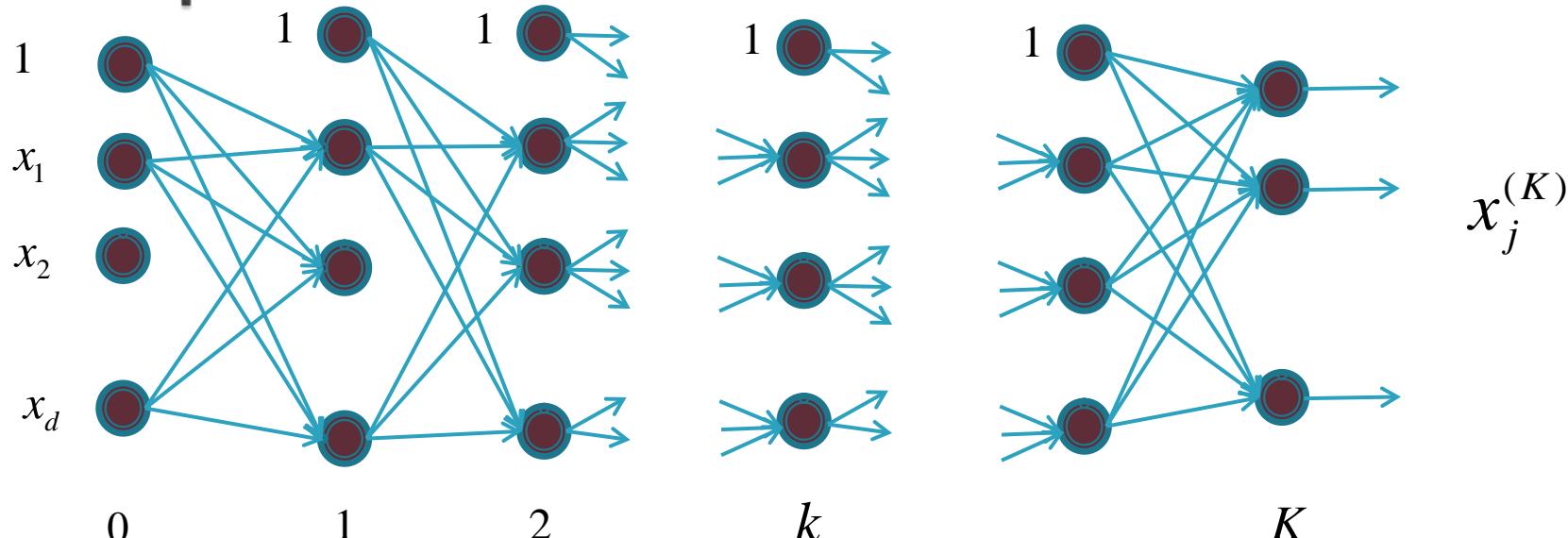
$$p(\omega_i / X) = \frac{p(X / \omega_i)P(\omega_i)}{p(X)}$$

Neural Network Function



$$f^{(K)}(f^{(K-1)} \dots (f^{(i)} \dots (f^{(2)}(f^{(1)}(X)))))$$

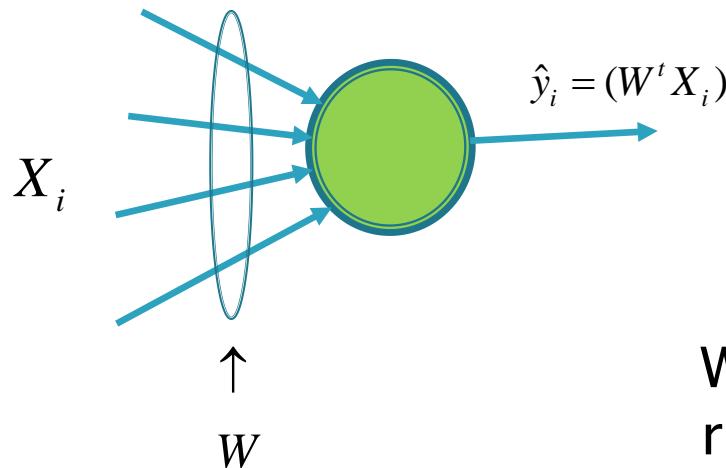
Multilayer Perceptron



$M_k \rightarrow$ No. of nodes in k^{th} layer

Back Propagation Learning

Single Layer Network- Single Output without nonlinearity



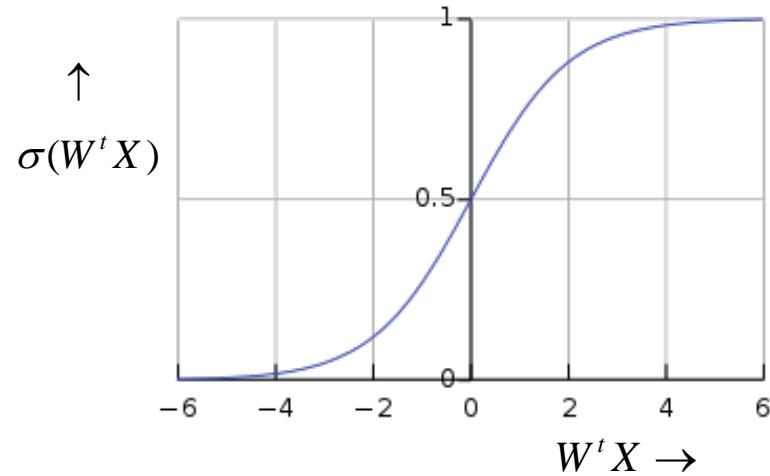
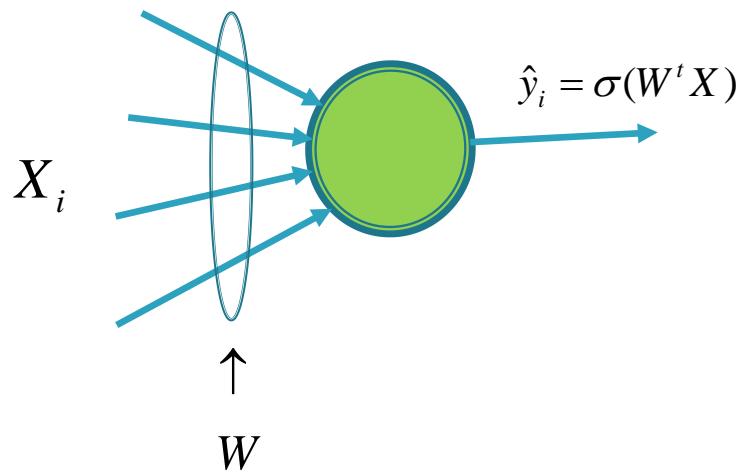
$$E = \frac{1}{2} \sum_{i=1}^N (W^t X_i - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$\nabla_W E = \sum_{i=1}^N (\hat{y}_i - y_i) X_i$$

Weight updation rule

$$W \leftarrow W - \eta \sum_{i=1}^N (\hat{y}_i - y_i) X_i$$

Single Layer Network– Single Output with nonlinearity



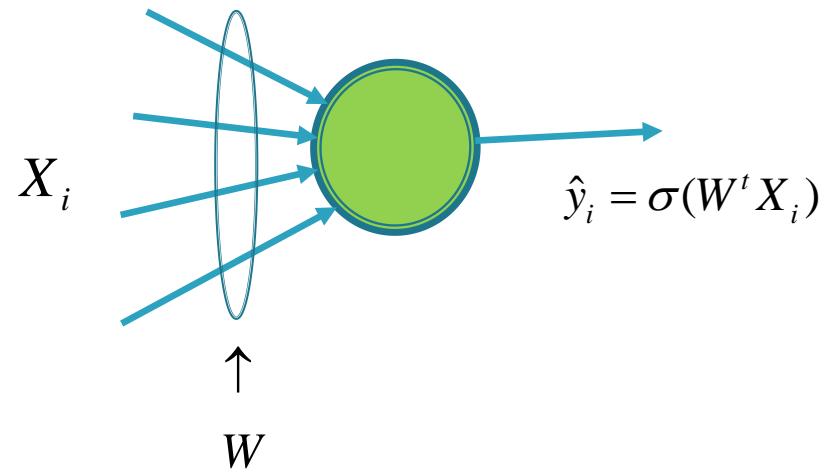
Single Layer Network– Single Output with nonlinearity

$$E = \frac{1}{2}(\hat{y}_i - y_i)^2 = \frac{1}{2}(\sigma(W^t X_i) - y_i)^2$$

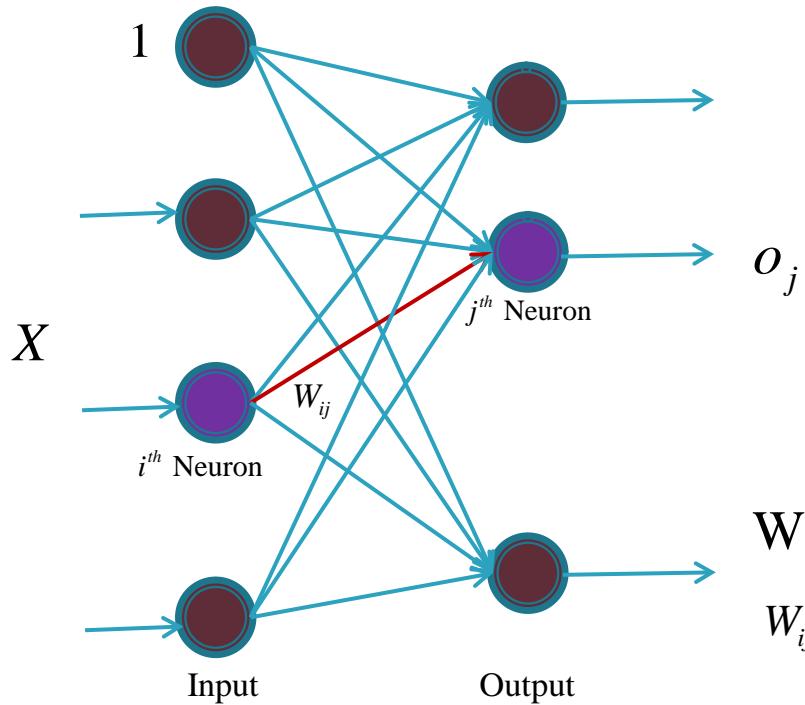
$$\frac{\partial E}{\partial W} = \hat{y}_i(1 - \hat{y}_i)(\hat{y}_i - y_i)X_i$$

Weight updation rule \Rightarrow

$$W \leftarrow W - \eta \hat{y}_i(1 - \hat{y}_i)(\hat{y}_i - y_i)X_i$$



Back Propagation Learning:- Single Layer Multiple Output



$$o_j = \frac{1}{1 + e^{-\theta_j}} \quad \theta_j = \sum_{i=1}^D W_{ij} x_i$$

$$E = \frac{1}{2} \sum_{j=1}^M (o_j - t_j)^2$$

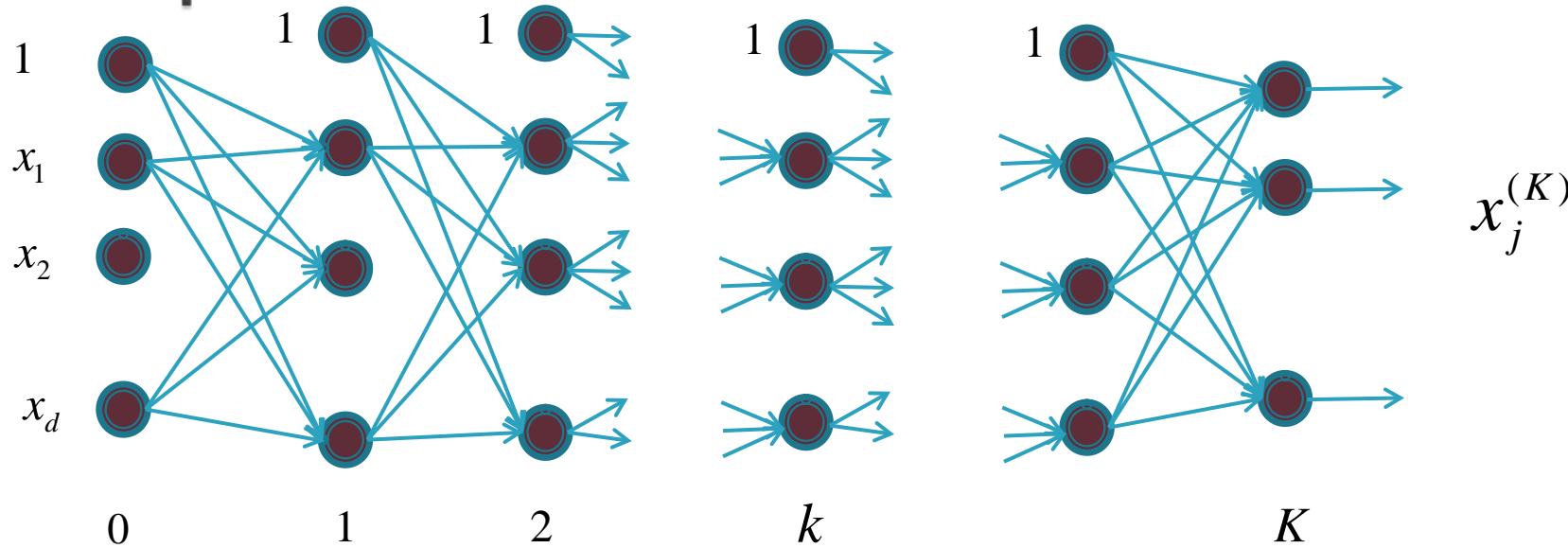
$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial \theta_j} \cdot \frac{\partial \theta_j}{\partial W_{ij}}$$

$$= (o_j - t_j) o_j (1 - o_j) x_i$$

Weight updation rule \Rightarrow

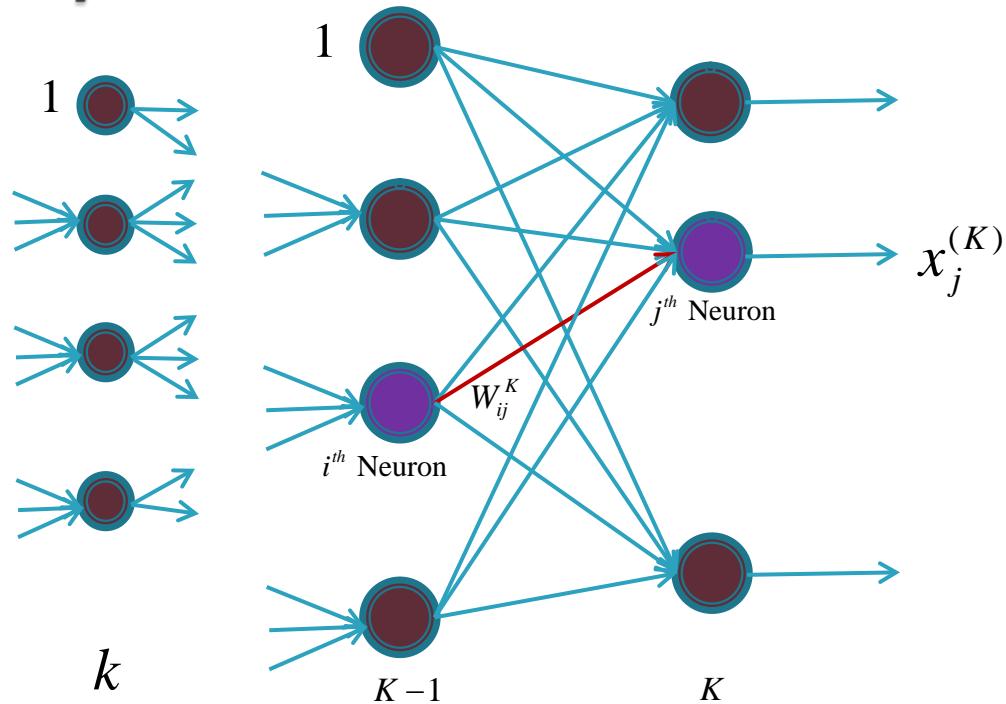
$$W_{ij} \leftarrow W_{ij} - \eta (o_j - t_j) o_j (1 - o_j) x_i$$

Multilayer Perceptron



$M_k \rightarrow$ No. of nodes in k^{th} layer

Back Propagation Learning:- Output Layer



$$x_j^K = \frac{1}{1 + e^{-\theta_j^K}} \quad \theta_j^K = \sum_{i=1}^{M_{K-1}} W_{ij}^K x_i^{K-1}$$

$$E = \frac{1}{2} \sum_{j=1}^{M_K} (x_j^K - t_j)^2$$

Back Propagation Learning:- Output Layer

Find W_{ij}^K that minimizes $E = \frac{1}{2} \sum_{j=1}^{M_K} (x_j^K - t_j)^2$

Gradient Descent $\nabla_{W_{ij}^K} E$

Back Propagation Learning:- Output Layer

$$\frac{\partial E}{\partial W_{ij}^K} = \frac{\partial E}{\partial x_j^K} \cdot \frac{\partial x_j^K}{\partial \theta_j^K} \cdot \frac{\partial \theta_j^K}{\partial W_{ij}^K}$$

$$= (x_j^K - t_j) x_j^K (1 - x_j^K) x_i^{K-1}$$

Let $\delta_j^K = x_j^K (1 - x_j^K) (x_j^K - t_j)$

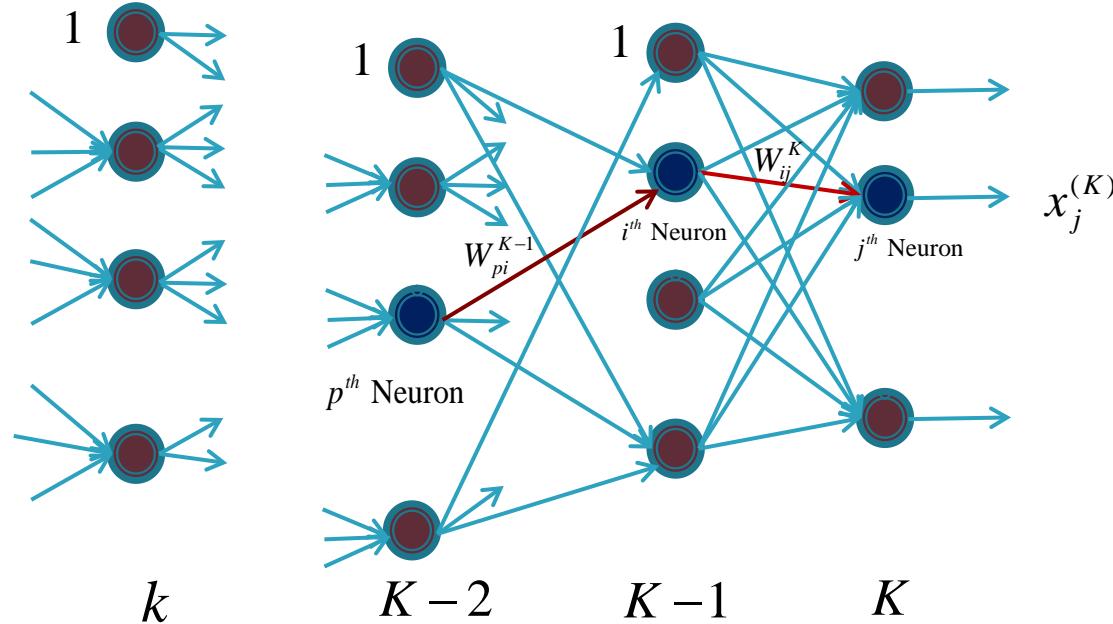
$$\Rightarrow \frac{\partial E}{\partial W_{ij}^K} = \delta_j^K x_i^{K-1}$$

$$x_j^K = \frac{1}{1 + e^{-\theta_j^K}} \quad \theta_j^K = \sum_{i=1}^{M_{K-1}} W_{ij}^K x_i^{K-1}$$

Weight updation rule
Output Layer

$$W_{ij}^K \leftarrow W_{ij}^K - \eta \delta_j^K x_i^{K-1}$$

Back Propagation Learning:- Hidden Layer



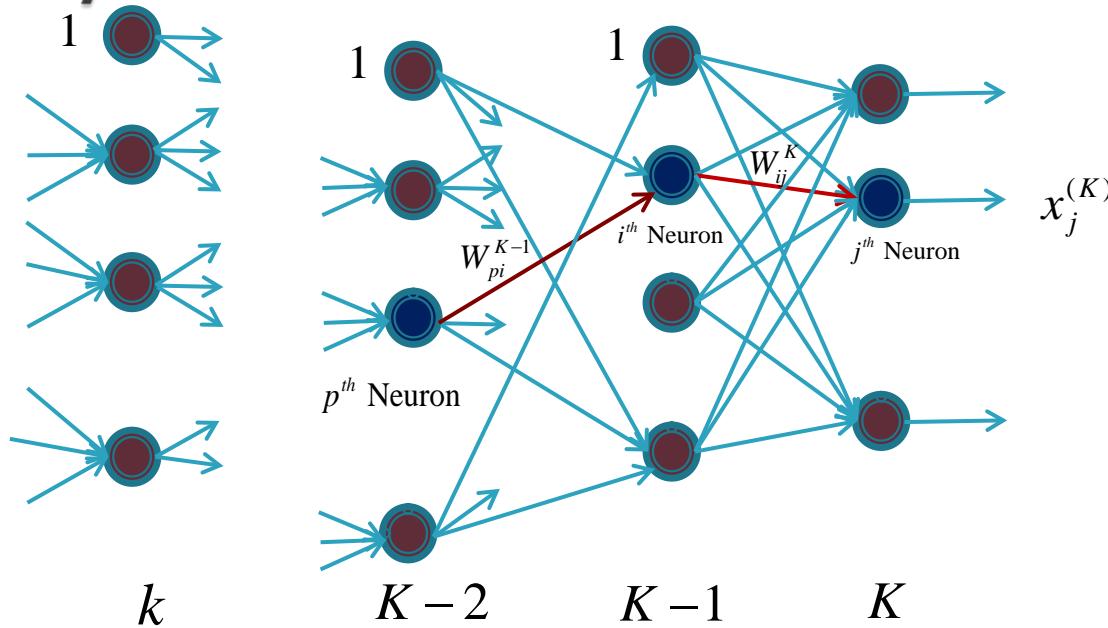
$$E = \frac{1}{2} \sum_{j=1}^{M_K} (x_j^K - t_j)^2$$

Back Propagation Learning:- Hidden Layer

Find W_{pi}^{K-1} that minimizes $E = \frac{1}{2} \sum_{j=1}^{M_K} (x_j^K - t_j)^2$

Gradient Descent $\Rightarrow \nabla_{W_{pi}^{K-1}} E$

Back Propagation Learning:- Hidden Layer



$$x_i^{K-1} = \frac{1}{1 + e^{-\theta_i^{K-1}}}$$
$$\theta_i^{K-1} = \sum_{p=1}^{M_{K-2}} W_{pi}^{K-1} x_p^{K-2}$$

Back Propagation Learning:- Hidden Layer

$$\frac{\partial E}{\partial W_{pi}^{K-1}} = \frac{\partial E}{\partial x_i^{K-1}} \cdot \frac{\partial x_i^{K-1}}{\partial W_{pi}^{K-1}}$$

$$= \frac{\partial E}{\partial x_i^{K-1}} \cdot \frac{\partial x_i^{K-1}}{\partial \theta_i^{K-1}} \cdot \frac{\partial \theta_i^{K-1}}{\partial W_{pi}^{K-1}}$$

$$= \frac{\partial E}{\partial x_i^{K-1}} \cdot x_i^{K-1} (1 - x_i^{K-1}) \cdot x_p^{K-2}$$

$$x_i^{K-1} = \frac{1}{1 + e^{-\theta_i^{K-1}}}$$

$$\theta_i^{K-1} = \sum_{p=1}^{M_{K-2}} W_{pi}^{K-1} x_p^{K-2}$$

Back Propagation Learning:- Hidden Layer

$$\frac{\partial E}{\partial x_i^{K-1}} = \frac{\partial E}{\partial x_j^K} \cdot \frac{\partial x_j^K}{\partial \theta_j^K} \cdot \frac{\partial \theta_j^K}{\partial x_i^{K-1}}$$

$$= \sum_{j=1}^{M_K} (x_j^k - t_j) x_j^K (1 - x_j^K) W_{ij}^K$$

$$= \sum_{j=1}^{M_K} \partial_j^K W_{ij}^K$$

$$E = \frac{1}{2} \sum_{j=1}^{M_K} (x_j^K - t_j)^2$$

$$x_j^K = \frac{1}{1 + e^{-\theta_j^K}} \quad \theta_j^K = \sum_{i=1}^{M_{K-1}} W_{ij}^K x_i^{K-1}$$

$$\delta_j^K = x_j^K (1 - x_j^K) (x_j^K - t_j)$$

Back Propagation Learning:- Hidden Layer

$$\frac{\partial E}{\partial W_{pi}^{K-1}} = \frac{\partial E}{\partial x_i^{K-1}} \cdot x_i^{K-1} (1 - x_i^{K-1}) \cdot x_p^{K-2} = x_i^{K-1} (1 - x_i^{K-1}) \cdot x_p^{K-2} \sum_{j=1}^{M_K} \partial_j^K W_{ij}^K$$

Putting $\delta_i^{K-1} = x_i^{K-1} (1 - x_i^{K-1}) \sum_{j=1}^{M_K} \partial_j^K W_{ij}^K$

Weight updation rule

Last but Output Layer

$$W_{pi}^{K-1} \leftarrow W_{pi}^{K-1} - \eta \delta_i^{K-1} x_p^{K-2}$$

Back Propagation Learning:- Hidden Layer

For any hidden layer weight W_{ij}^k

Putting $\delta_i^k = x_i^k (1 - x_i^k) \sum_{j=1}^{M_{k+1}} \partial_j^{k+1} W_{ij}^{k+1}$

Weight updation rule

$$W_{ij}^k \leftarrow W_{ij}^k - \eta \delta_j^k x_i^{k-1}$$

Machine Learning

vs

Deep Learning

Convolution

Kernel

1 D → 2A+1

$$y(n) = \sum_{p=-A}^A w(p)x(n-p)$$

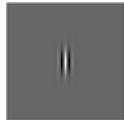
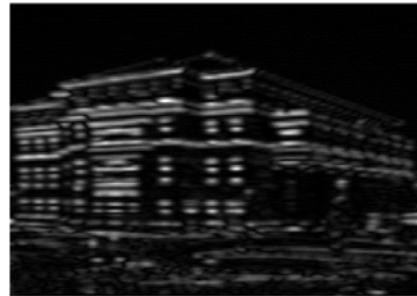
2 D →

(2A+1) × (2A+1)

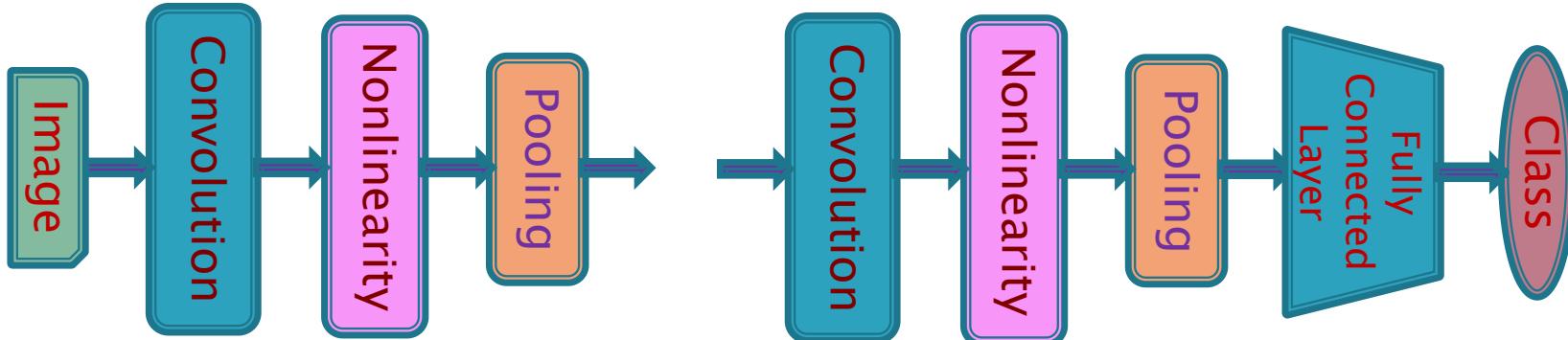
$$y(m, n) = \sum_{p=-A}^A \sum_{q=-A}^A w(p, q)x(m-p, n-q)$$

Finite Convolution Kernel

Feature at a point is local in nature

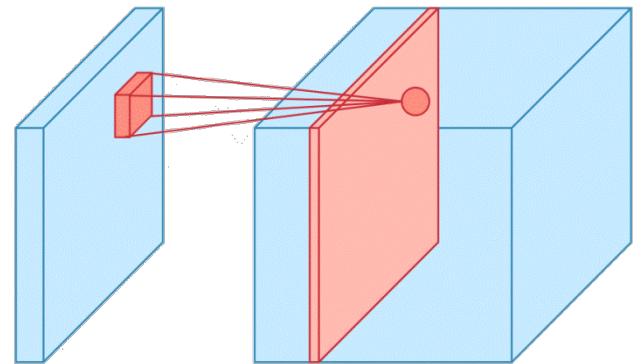


CNN Architecture



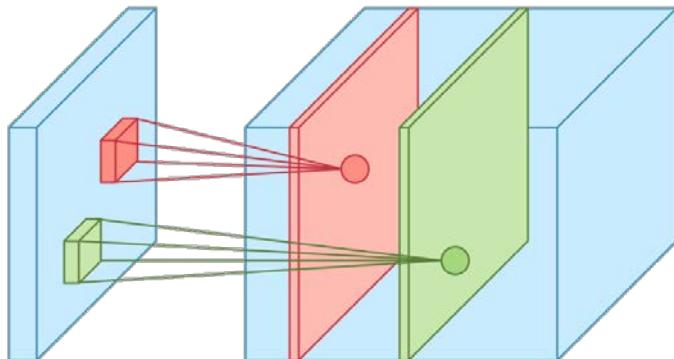
3 D Convolution– Visualization

- The kernel strides over the input Image.
- At each location $I(m, n)$ compute $f(m, n) = \sum \sum w(p, q)I(p - m, q - n)$ collect them in the feature map.
- The animation shows the sliding operation at 4 locations, but in reality it is performed over the entire input.



Animation:- Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

3 D Convolution- Visualization



- ▶ Red and green boxes are two different featured maps obtained by convolving the same input with two different kernels.
- ▶ The feature maps are stacked along the depth dimension as shown.

Figure: Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

3 D Convolution- Visualization

- ▶ An RGB Image of size $32 \times 32 \times 3$
- ▶ 10 Kernels of size $5 \times 5 \times 3$
- ▶ Output featuremap of size $32 \times 32 \times 10$

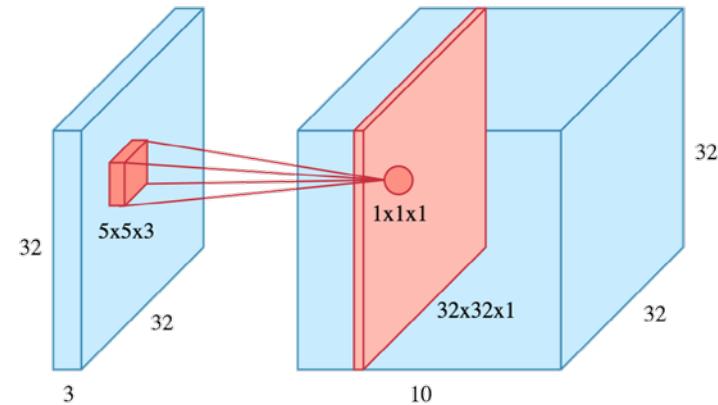


Figure: Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

Nonlinearity

- ReLU is an element wise operation (applied per pixel) and replaces all negative pixel values in the feature map by zero

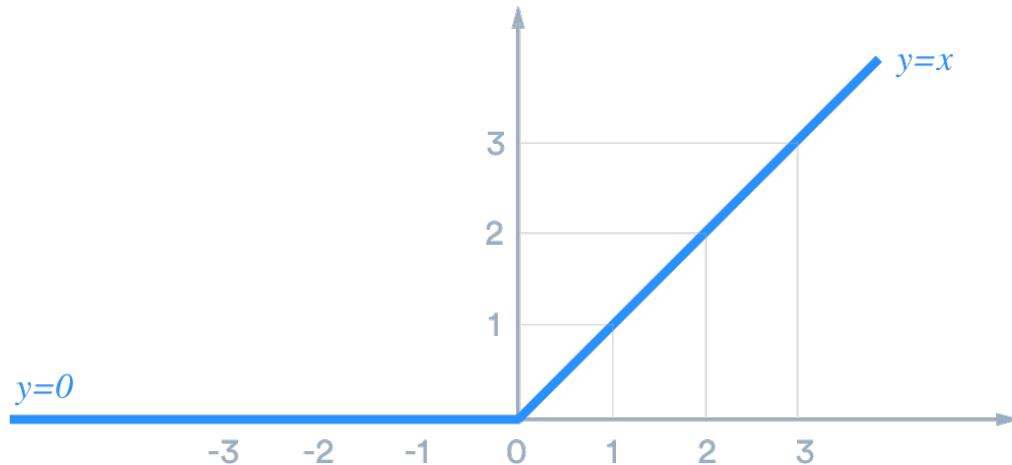


Figure: Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

Pooling

- ▶ Replaces the output of a node at certain locations with a summary statistic of nearby locations.
- ▶ Spatial Pooling can be of different types: Max, Average, Sum etc.
- ▶ Max Pooling report the maximum output within a rectangular neighborhood.
- ▶ Pooling helps to make the output approximately invariant to small translation.
- ▶ Pooling layers down sample each feature map independently, reducing the height and width, keeping the depth intact.
- ▶ In pooling layer stride and window size needs to be specified

Pooling

- Figure below is the result of max pooling using a 2×2 window and stride 2. Each color denotes a different window. Since both the window size and stride are 2, the windows are not overlapping

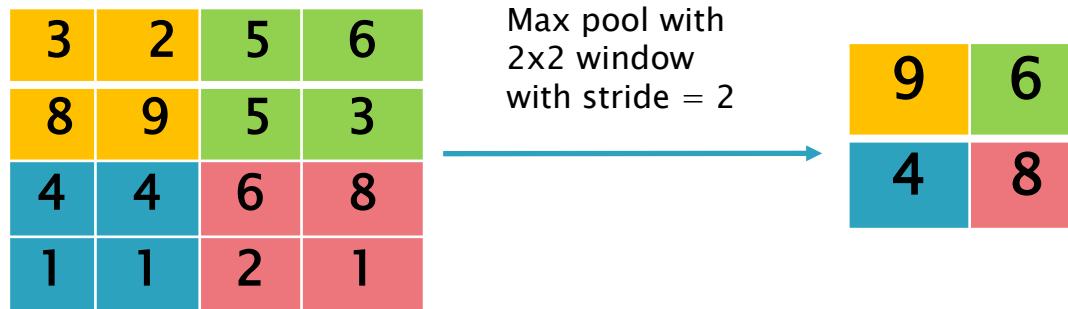


Figure: Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

Pooling

- ▶ Pooling reduces the height and the width of the feature map, but the depth remains unchanged as shown in figure
- ▶ Pooling operation is independently carried out across each depth

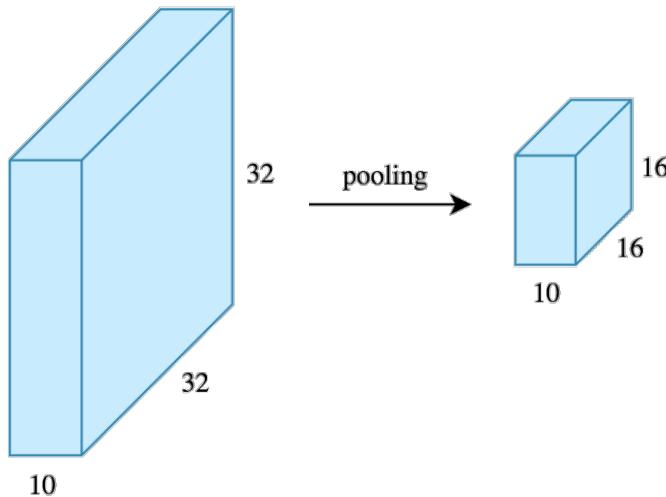
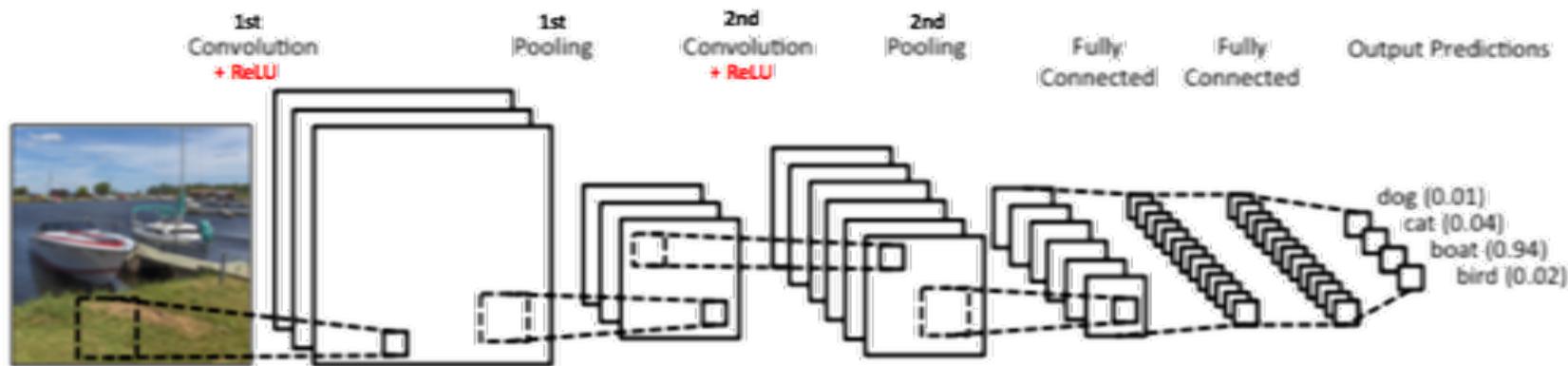
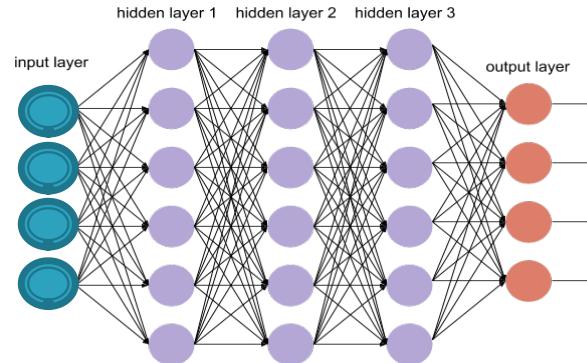
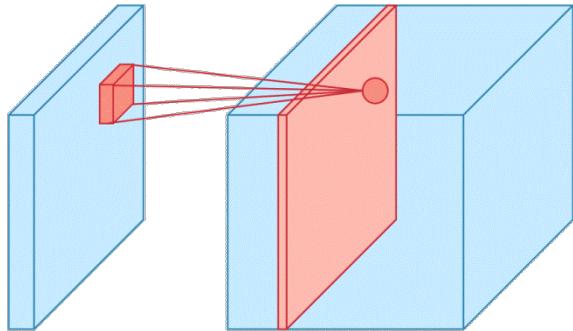


Figure: Arden Dertat <https://towardsdatascience.com/applied-deep-learning-part-4-convolutional-neural-networks-584bc134c1e2>

CNN Architecture



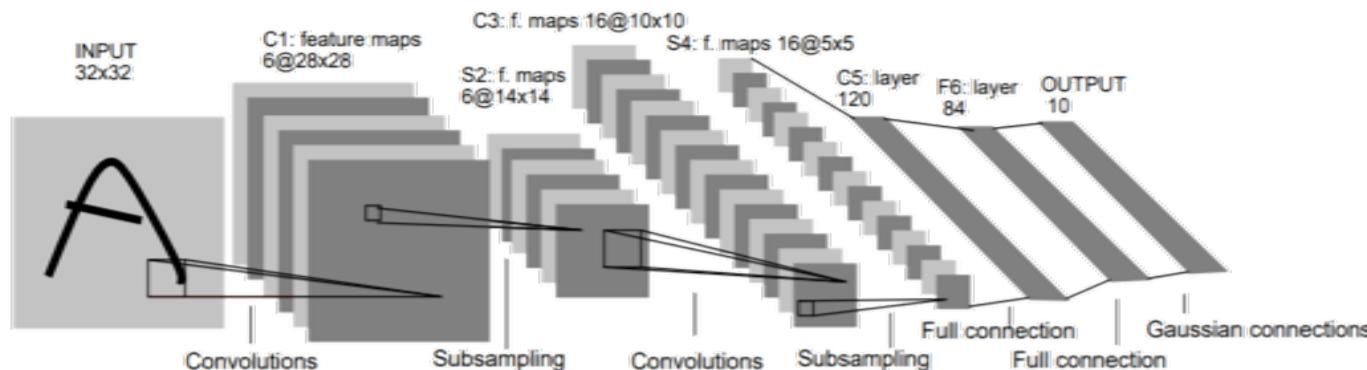
MLP vs CNN



- ❑ Sparse Connectivity: Every node in the Convolution Layer receives input from a small number of nodes in the previous layer (Receptive Field), needing smaller number of parameters.
- ❑ Parameter Sharing: Each member of the Convolution Kernel is used at every position of the input, dramatically reducing the number of parameters.
- ❑ This makes CNN much more efficient than MLP.

LeNet 5

- ▶ Proposed by Yann LeCun, Leon Bottou, Yoshua Bengio and Patrick Haffner for handwritten and machine-printed character recognition.
- ▶ Used by many Banks for recognition of hand written numbers on cheques.
- ▶ This architecture achieves an error rate as low as 0.95% on test data



Yann LeCun, Leon Bottou, Yoshua Bengio and Patrick Haffner, "Gradient -Based Learning Applied to Document Recognition", Proc. IEEE, Nov. 1998

IMAGENET Large Scale Visual Recognition Challenge (ILSVRC)



ILSVRC

- ▶ IMAGENET Large Scale Visual Recognition Challenge.
- ▶ Evaluates algorithms for Object Detection and Image Classification on large image database.
- ▶ Helps researchers to review state of the art Machine Learning techniques for object detection across a wider variety of objects.
- ▶ Monitor the progress of computer vision for large scale image indexing for retrieval and annotation.
- ▶ Database contains large number of Images from 1000 categories.
- ▶ More than 1000 images in every category.

ILSVRC

- ▶ Every year of the challenge the forum also organizes a workshop at one of the premier computer vision conferences.
- ▶ The purpose of the workshop is to disseminate the new findings of the challenge.
- ▶ Contestants with the most successful and innovative techniques are invited to present their work.

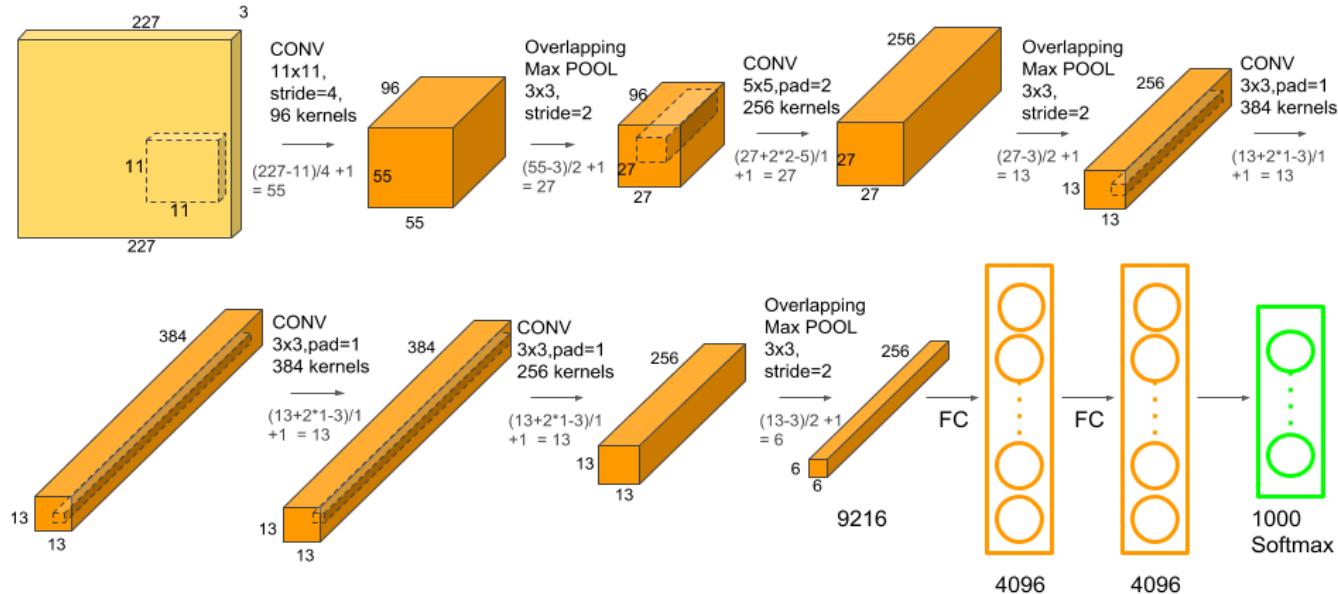
AlexNet

ILSVRC 2012 Winner

Krizhevsky Alex, Ilya Sutskever and Geoffrey E. Hilton, “Imagenet Classification with deep convolutional neural networks”, Advances in Neural Information Processing Systems, 2012

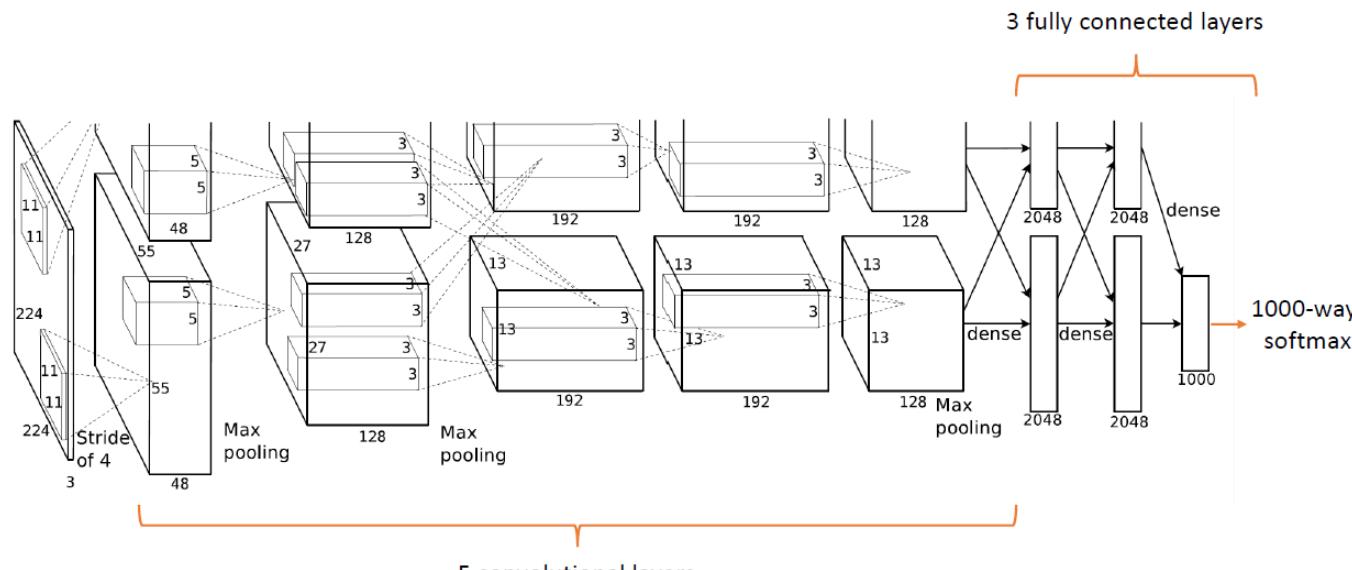
AlexNet

ILSVRC 2012
Winner



AlexNe

†



6

Krizhevsky Alex, Ilya Sutskever and Geoffrey E. Hilton, "Imagenet Classification with deep convolutional neural networks", Advances in Neural Information Processing Systems, 2012

AlexNet

- 60 Million parameters and 650000 neurons.
- The network is split into two pipelines and was trained on two GPU.
- Input Image size 256 x 256 RGB.
- Grey scale images to be replicated to obtain 3-Channel RGB
- Random crops of size 227 x 227 are fed to the input layer of AlexNet.
- Stochastic Gradient Descent with Momentum Optimizer.
- Top-5 error rate 15.3%.

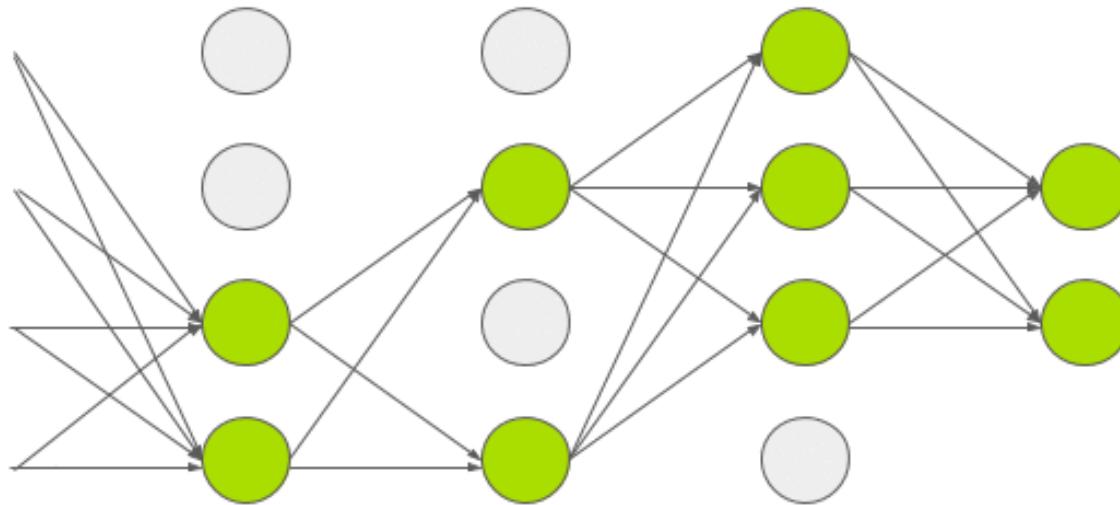
Krizhevsky Alex, Ilya Sutskever and Geoffrey E. Hilton, “Imagenet Classification with deep convolutional neural networks”, Advances in Neural Information Processing Systems, 2012

Reducing Overfitting

- Train the network with different variants of the same image helps avoiding overfitting
 - ❖ Generate additional data from existing data (Augmentation)
 - ❖ Data augmentation by mirroring
 - ❖ Data Augmentation by random crops
- Dropout Regularization

Krizhevsky Alex, Ilya Sutskever and Geoffrey E. Hilton, “Imagenet Classification with deep convolutional neural networks”, Advances in Neural Information Processing Systems, 2012

AlexNet

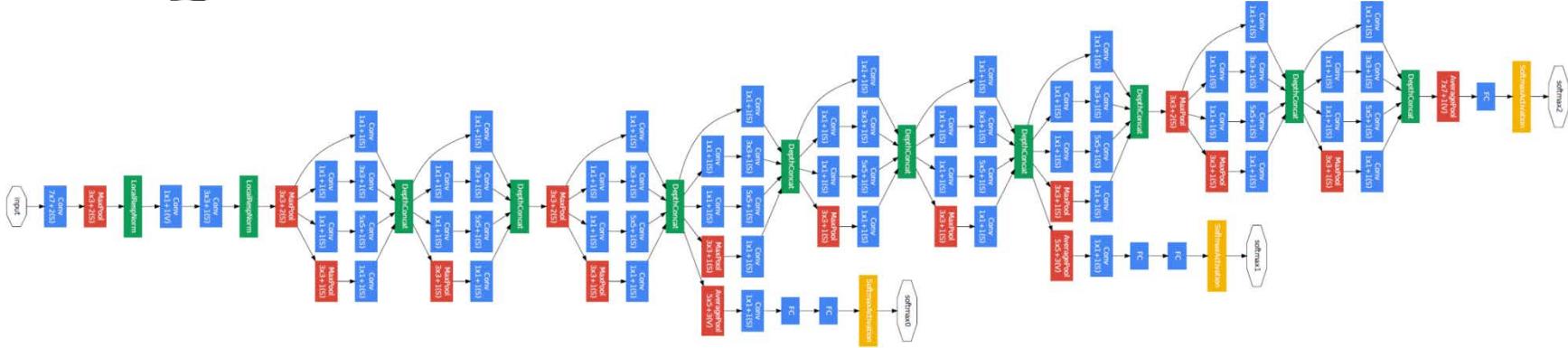


<https://www.learnopencv.com/understanding-alexnet/>

GoogLeNet



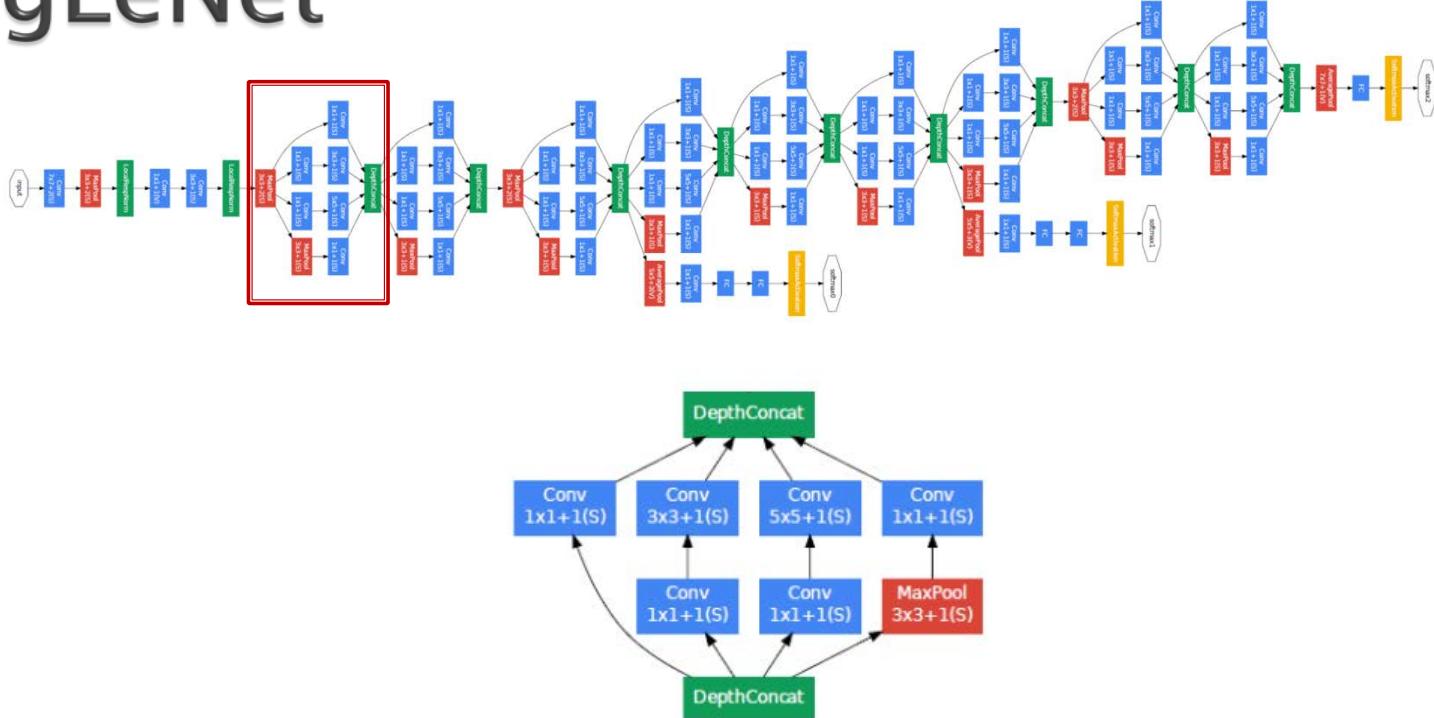
GoogLeNet



- ❖ 22 Layers with parameters
- ❖ 27 Layer including Maxpool layers

Convolution Layer
 Maxpool Layer
 Feature Concatenation
 Softmax Layer

GoogLeNet



Krizhevsky Alex, Ilya Sutskever and Geoffrey E. Hilton, “Imagenet Classification with deep convolutional neural networks”, Advances in Neural Information Processing Systems, 2012

Discriminative vs. Generative Model

Image Source: Internet

Discriminative Model

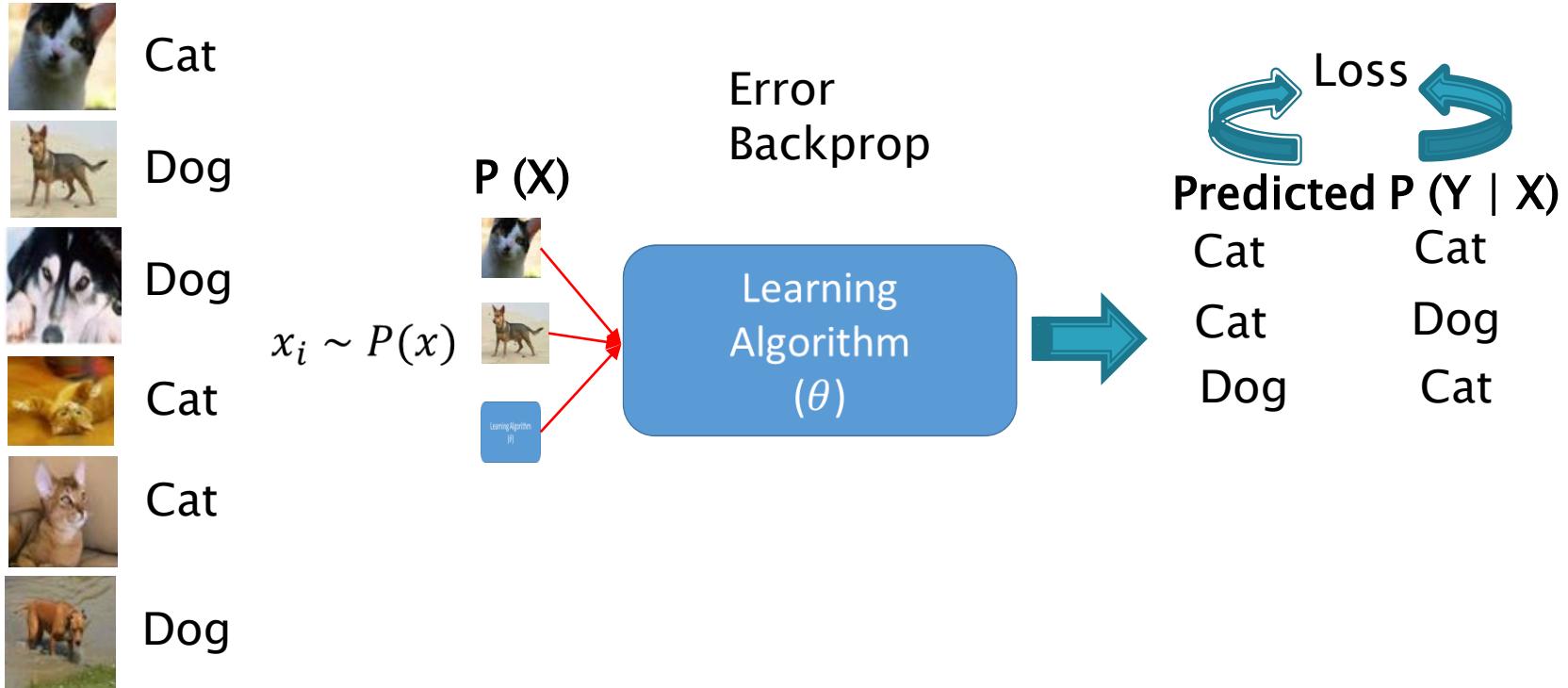


Image Source: Internet

Generative Model

“What I can not create, I do not understand.”

– Richard Feynman

- Collect a large amount of data in some domain
- Train a model to generate data like it.

Autoencoder

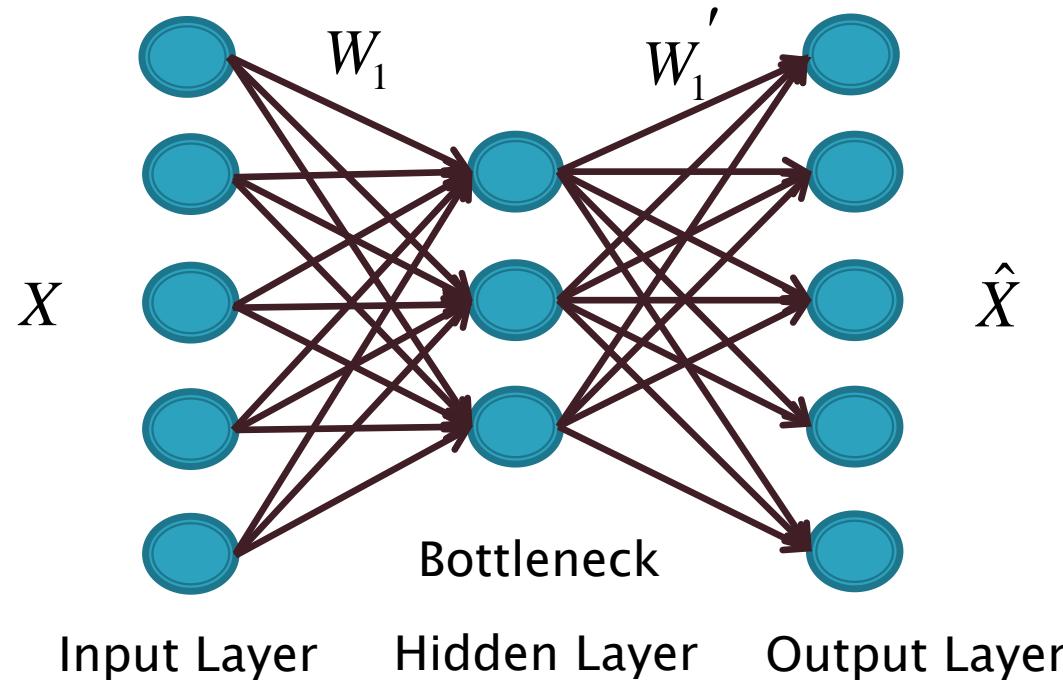
- ❖ Unsupervised Learning where Neural Networks are subject to the task of representation learning.
- ❖ Impose a bottleneck in the network
- ❖ The bottleneck forces a compressed knowledge representation of the input.

Autoencoder

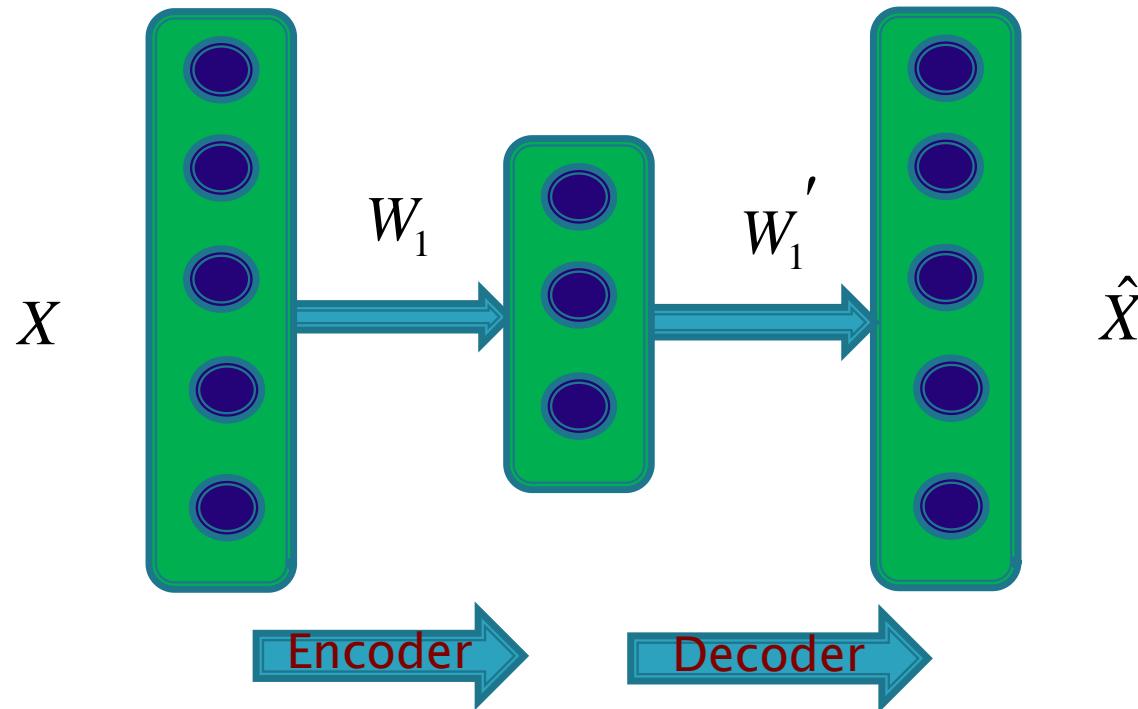
Assumption:

- High degree of correlation/structure exists in the data.
- For uncorrelated data (input features are independent), then compression and subsequent reconstruction would be difficult.

Autoencoder



Autoencoder



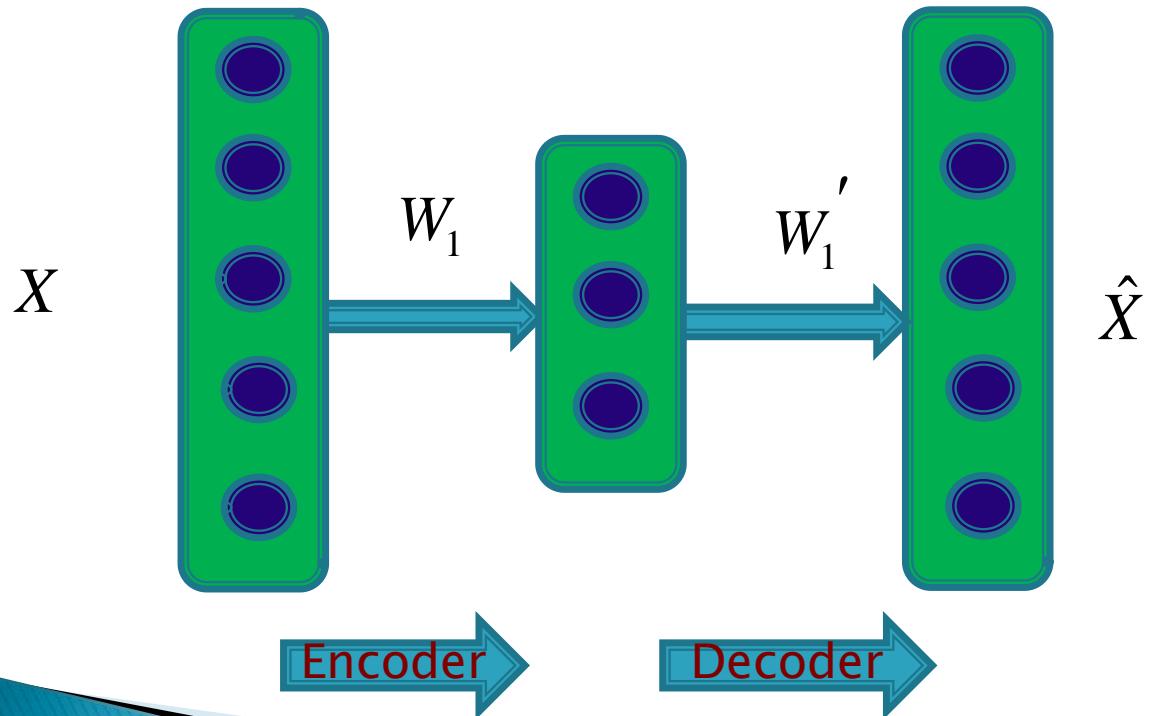
Expectation

- ❑ Sensitive enough to input for accurate reconstruction
- ❑ Insensitive enough that it does not memorize or overfit the training data



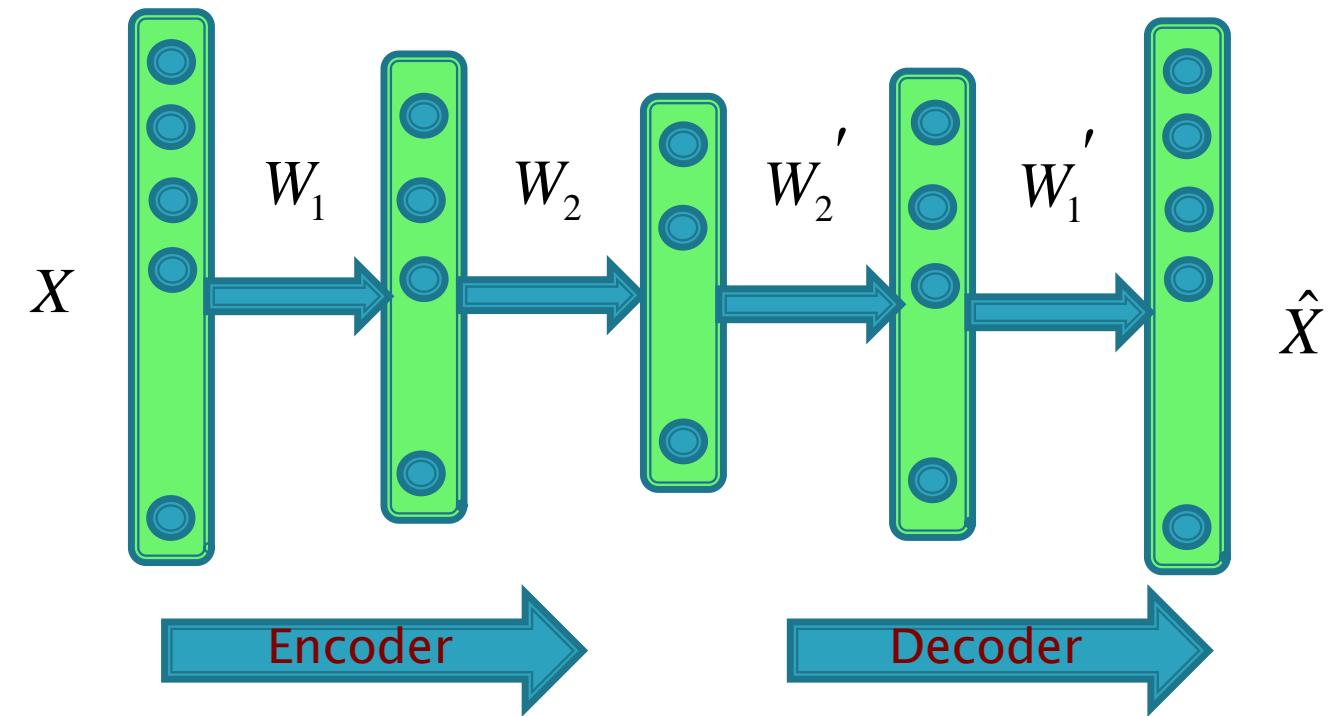
Loss Function $\Rightarrow L(X, \hat{X}) + \text{Regularizer}$

Undercomplete Autoencoder



$$L(X, \hat{X}) = \frac{1}{2} \sum_N \|X - \hat{X}\|^2$$

Stacked Autoencoder



$$L(X, \hat{X}) = \frac{1}{2} \sum_N \|X - \hat{X}\|^2$$