

# Machine Learning

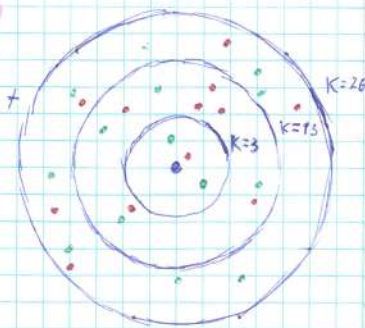
## Types of ML Problems:

- **Unsupervised Learning** - the machine did not use marked down data. Could find Hidden patterns.
- **Supervised Learning** - the machine use marked down data as training one. Could predict output due to previous training data experience.
- **Reinforcement Learning** - the machine use real-time simulation (system) where it have "rewards" and "punishments". Could simulate complicated behaviour and learn by own mistakes.



## K.N.N. ("K" - Nearest Neighbors)

**Idea:** Consider input data as point  $\langle \text{vector} \rangle$  in  $n$ -dimensional space. Then use "K" as Hyperparameter to find K-Nearest Neighbours. Simply count amount of points of each class among neighbours, so the answer is the greater amount class (for classification). Use mean ( $\bar{X}$ ) of all neighbours (for Regression).



## Problems:

- **Curse of Dimensionality** - The more dimensionality is, the greater point Dispersion is  $\rightarrow$  less points would be in "K" area.
- **Feature Scale** - Features have different Metrics, so one feature might be much more important than the other:  
 $[0 \dots 2000](m), [0 \dots 2](Km) \rightarrow 2m \sim 2Km$  **Wrong!**

## Distance Metrics:

**Euclidean:**  $D(\bar{a}, \bar{b}) = \sqrt{\sum_{i=1}^n (\bar{a}_i - \bar{b}_i)^2}$

**Manhattan:**  $D(\bar{a}, \bar{b}) = \sum_{i=1}^n |\bar{a}_i - \bar{b}_i|$

**Minkowski:**  $D(\bar{a}, \bar{b}) = \left( \sum_{i=1}^n |\bar{a}_i - \bar{b}_i|^p \right)^{\frac{1}{p}}$



# Learning

## Accuracy Metrics:

It is a performance indicators (functions) that helps to evaluate the correctness of the model:

### Classifier:

$$\text{Precision} = \frac{\text{Correct Pred. (class, "n")}}{\text{total (class, "n")}}$$

$$\text{Recall} = \frac{\text{correct Pred. (class, "n")}}{\text{total Pred. (class, "n")}}$$

$$\text{Accuracy} = \frac{\text{Pred. Correct}}{\text{total pred.}}$$

$$F1 = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

### Regressor:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \text{mean}(y))^2}$$

## Loss $L(y, \hat{y})$ :

It is a function that shows how did prediction " $\hat{y}$ " differs from correct answer " $y$ ". The lessier it is, the better predicted result you have.

## Over fitting:

Situation when model loss on training data  $\rightarrow 0$  and loss on real data  $\rightarrow 100$ .



## Split Data:

all available data

Training	Validation	Testing
----------	------------	---------

all available data.

Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Testing
--------	--------	--------	--------	--------	---------

split 1	1	2	3	4	5	Testing
split 2	1	2	3	4	5	
split 3	1	2	3	4	5	
split 4	1	2	3	4	5	
split 5	1	2	3	4	5	

## Data Normalization:

ALWAYS Normalize The Features!

$$x_{\text{new}} = \frac{x - \mu}{\sigma}$$

$\mu$  - mean value of feature.

$\sigma$  - Standard deviation of feature.

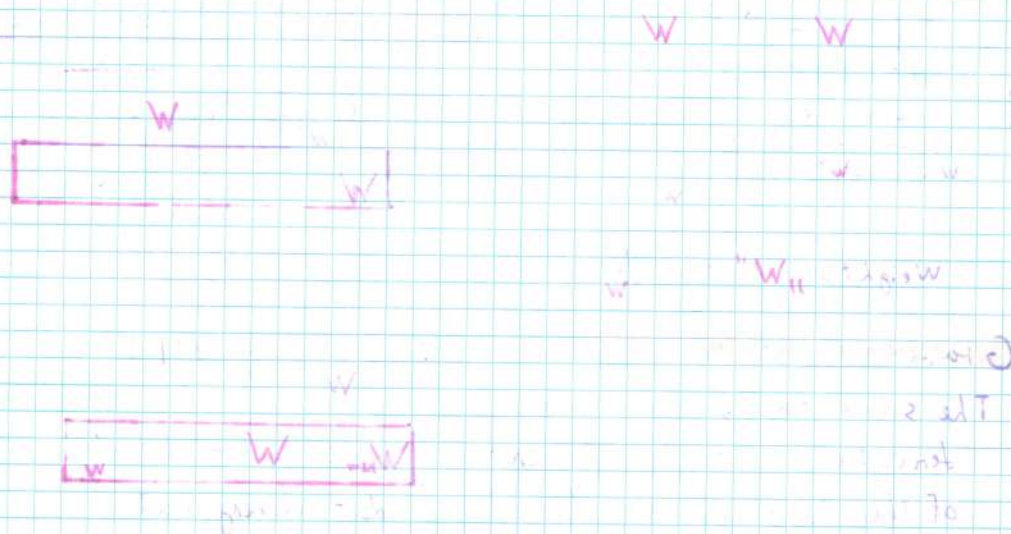
## Regularization:

A method of adding special constraints to the loss function to decrease model complexity and prevent overfitting.

$$L_1 \text{ (Lasso)} = \dots + \beta \sum_{i=1}^m |w_i|$$

$$L_2 \text{ (Ridge)} = \dots + \beta \sum_{i=1}^m w_i^2$$

$$\text{Elastic Net} = \dots + \beta_1 \sum_{i=1}^m |w_i| + \beta_2 \sum_{i=1}^m w_i^2$$





## INPUT:

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & x_{21} & \dots & x_{2n} \\ 1 & x_{31} & \dots & x_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{mn} \end{bmatrix}$$

features

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix}$$

correct answers

## Model:

$$W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix}$$

weights

## Output:

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_m \end{bmatrix}$$

prediction result

## Train:

$$X \cdot W = \hat{Y}$$

$$L(W) = \frac{1}{m} (XW - Y)^T (XW - Y) = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad \# \text{Mean Square Error}$$

## Analytical Method (not used):

The main idea is to minimize loss function "L" by finding such weights "W" to  $\frac{\partial L}{\partial W} = 0$

$$\frac{\partial L}{\partial W} = \frac{2}{m} X^T (XW - Y) = 0$$

$$W = (X^T X)^{-1} X^T Y \sim X^+ Y$$

## Gradient descent:

The same thing but consider derivative vector as gradient of the Loss function. (Repeat iteratively till the acceptable result)

$$\frac{\partial L}{\partial W} \sim \text{grad}(L)$$

$$W_{\text{new}} = W - \alpha \cdot \frac{\partial L}{\partial W}$$

$\alpha$  - Learning Rate

## If dependencies are Non-Linear:

Simply add more Polynomial features and normalize them.

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$X' = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$



## The Idea:

quite the same as in Linear Regression, But this one is used in Classification and have belonging class probability as an output. Also, for each class there is own Weights Vector.

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & x_{21} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{mn} & \dots & x_{mn} \end{bmatrix} \quad \begin{matrix} \text{Individuals} \\ \text{features} \end{matrix}$$

$$W = \begin{bmatrix} w_{00} & w_{01} & \dots & w_{0n} \\ w_{10} & w_{11} & \dots & w_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k0} & w_{k1} & \dots & w_{kn} \end{bmatrix} \quad \begin{matrix} \text{for each feature} \\ \text{for each class} \end{matrix}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \quad \begin{matrix} \text{Probab. of belong} \\ \text{to each class} \end{matrix}$$

## Train:

$$\text{logits}(X_{\text{individual}}) = X_{\text{individual}} \cdot W \sim [1 \ x_1 \ x_2 \dots x_n] \cdot \begin{bmatrix} w_{00} & \dots & w_{0n} \\ \vdots & \ddots & \vdots \\ w_{k0} & \dots & w_{kn} \end{bmatrix}$$

$$\text{Softmax}(\text{logits}(X_{\text{individual}})) = \hat{Y}$$

# The most possible class type is one with the highest probability.

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} \rightarrow \text{SoftMax} \left( \frac{e^z}{\sum e^z} \right) \rightarrow \begin{bmatrix} 0.02 \\ 0.9 \\ 0.09 \\ 0.01 \\ 0.02 \end{bmatrix}$$

## Loss:

$$L(W) = -\frac{1}{m} \sum_{i=1}^m \begin{cases} \log(1 - \hat{y}_i) & \text{if } y_i = 0 \\ \log(\hat{y}_i) & \text{if } y_i = 1 \end{cases}$$

# Log Loss - return huge

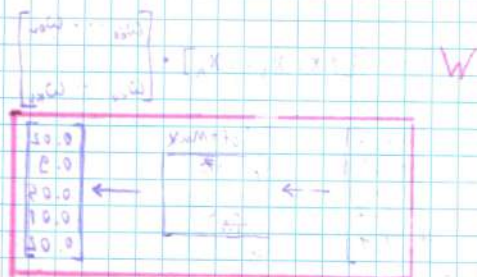
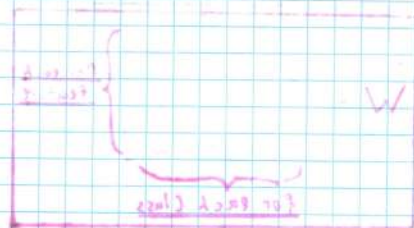
$$L(W) = -\frac{1}{m} \sum_{i=1}^m (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

Value if Prediction  $\frac{y_i}{\hat{y}_i}$  is complete  
Mismatch the correct class answer

## Gradient descent:

$$\frac{\partial L}{\partial W} = \frac{1}{n} X_{\text{individual}}^T (\hat{Y} - Y)$$





## Idea:

it is a logical algorithm that based on a binary tree. Each node contains a question (condition) connections - all possible results of a question, leafs - terminal node ready to predict.

After training (tree is complete) we could predict the result by simple tree traverse.

## Training:

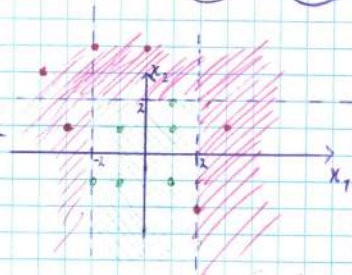
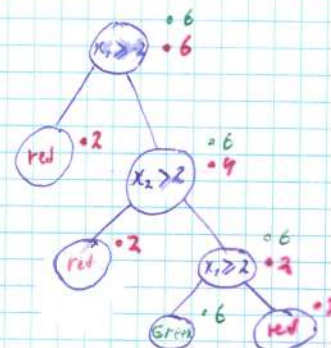
- Check all the leafs and calculate the Entropy or Gini Index for every possible combination of every feature value split ( $x_1 \leq 1, x_2 \geq 3, \dots$ )

- Calculate Information Gain for each possible split By:

$$IG = \text{entropy}(\text{parent}) - \sum_{i=0} W_i \cdot \text{entropy}(\text{child}_i)$$

where " $W_i$ " =  $\frac{\text{size}(\text{child}_i)}{\text{size}(\text{parent})}$

- Choose split with the highest IG value and repeat until "max-tree-depth" or every final node become leaf.



$$\text{entropy}(X) = - \sum_{i=1} P(x_i) \cdot \log_2 P(x_i)$$

# where " $p_i$ " - probability of a class.

$$P(\text{class}) = \frac{\text{Total objects of class amount}}{\text{Total objects}}$$

$$\text{gini-index}(X) = 1 - \sum_{i=1} P(x_i)^2$$

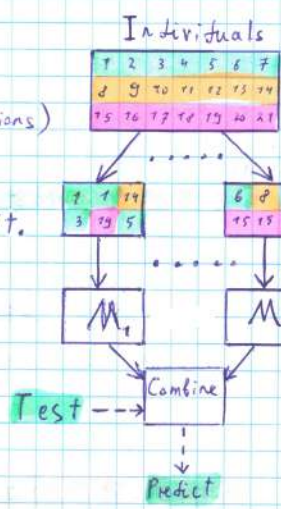


## Voting:

Use different algorithms with Prediction probabilities greater than  $\frac{1}{2}$  (50%) and simply get the mean value of all the results for Regression.  
For Classification problem the idea is the same, but we choose the most frequent class.

## Bagging (parallel):

Firstly, get some random samples (with repetitions) and train some models to get "Weak learners" so they would not overfit.  
Then, pass the test data through these models and mean the result.





## Random Forest:

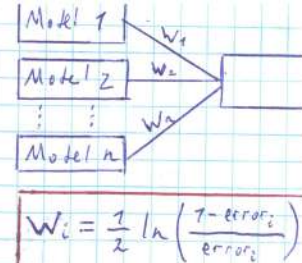
This is a variation of Bagging technique where we randomly sample not only individuals but their features too.

Also it is used Decision Tree as a model.

# Could learn very deep and hard dependencies.

Increase weight at those individuals in Training dataset.

- Repeat this steps and combine all weak learners into one strong learner using weights, which are depends on weak learner accuracy.



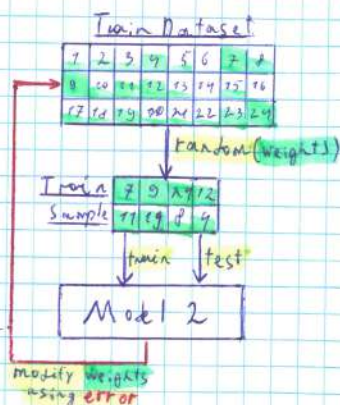
## Gradient Boosting:

### Boosting (sequential):

It is a process that uses a set of models (Weak Learners) combined into single Strong Learner model in order to increase the accuracy.

### Adaptive Boosting (AdaBoost):

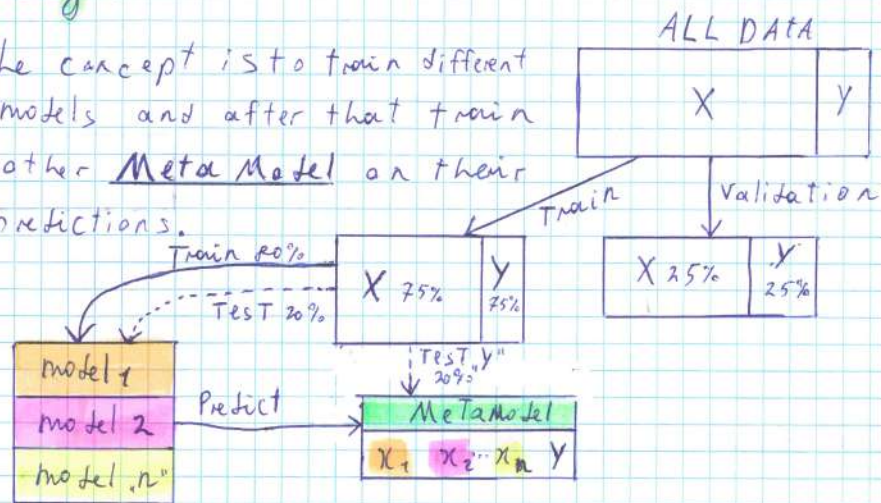
- Firstly add " $\frac{1}{n}$ " weight to each individ (or sample) in Training Dataset.
- Randomly (depending on weights) choose data to create training sample and train clear model (weak learner).
- Test this model on the same sample and mark up individuals were predicted with error.





## Blending:

The concept is to train different models and after that train other Meta Model on their predictions.

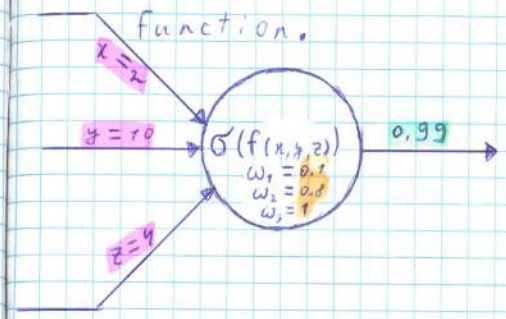


## Stacking:



## Neuron:

The main structure element of ANN's. The main idea behind it is to linearly transform some input data (constant, vector, matrix, tensor...) using neuron function and then apply non-linear activation function.



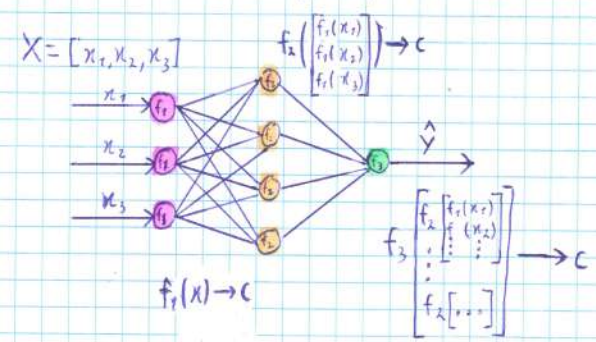
$$f(x, y, z) = w_1 \cdot x + w_2 \cdot y + w_3 \cdot z$$

$$\sigma(x) = \text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma(f(2, 10, 4)) = \text{Sigmoid}(0.1 \cdot 2 + 0.8 \cdot 10 + 1 \cdot 4) = \text{Sigmoid}(8.2) = 0.99$$

## Layer:

The way to group neurons of same type into one object.



- Input layer
- Output layer
- Other



### Idea:

The goal of any activation function is to prevent collapse of all linear function into one. Mostly, it also normalize data and should be differentiable.

# Consider sequential linear functions without activation

$$\boxed{w_1(k_1x + b_1)} \xrightarrow{f_1(x)} \boxed{w_2(k_2x + b_2)} \xrightarrow{f_2(x)} \boxed{w_3(k_3x + b_3)} \xrightarrow{f_3(x)}$$

$$f_3(f_2(f_1(x))) = w_3(k_3(w_2(k_2(w_1(k_1x + b_1)) + b_2)) + b_3) \sim \boxed{w(kx + b)}$$

could be simplified into this expression (linear)

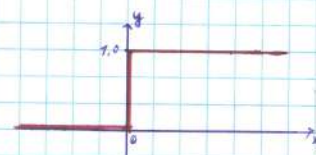
# Same but using activation:

$$\sigma f_3(\sigma f_2(\dots)) = \sigma(w_3(k_3(\sigma w_2(\dots) + b_3))) \rightarrow \text{could not be simplified because } \sigma \text{ -function!}$$

### Binary Step:

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

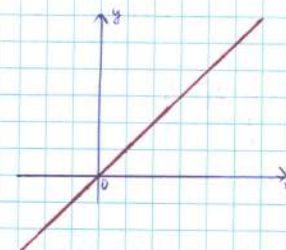
$$f'(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ \text{none} & \text{if } x = 0 \end{cases}$$



### Linear:

$$f(x) = x$$

$$f'(x) = 1$$

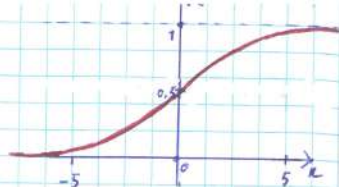




Sigmoid (logistic):

$$f(x) = \frac{1}{1+e^{-x}}$$

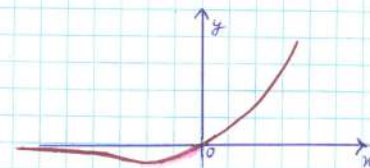
$$f'(x) = f(x) \cdot (1-f(x))$$



Swish:

$$f(x) = x \cdot \text{sigmoid}(x) \sim \frac{x}{1+e^{-x}}$$

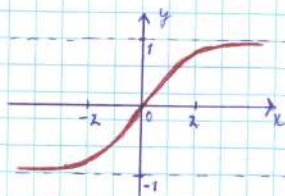
$$f'(x) = f(x) + \text{sigmoid}(x) \cdot (1-f(x))$$



Tanh:

$$f(x) = \frac{2}{1+e^{-2x}} - 1$$

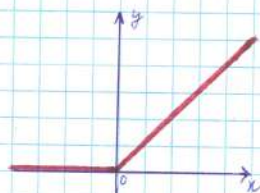
$$f'(x) = 1 - f(x)^2$$



ReLU (Rectified Linear Unit):

$$f(x) = \max(0, x) \sim \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

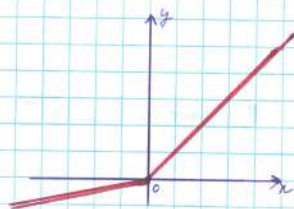
$$f'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Leaky ReLU:

$$f(x) = \max(\alpha x, x) \sim \begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

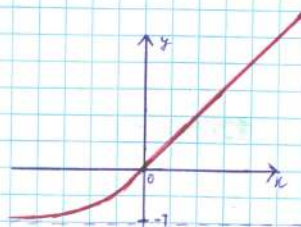
$$f'(x) = \begin{cases} \alpha & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



ELU (Exponential Linear Unit):

$$f(x) = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} f(x) + \alpha & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



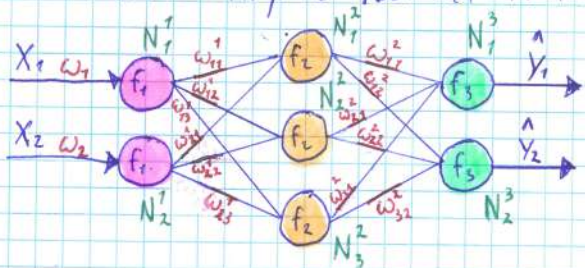


## Idea:

The main idea is to update all the parameters (weights) using gradient descent to minimize Loss function.

## Example:

Consider Simple Neural Network:



$$f_1(x) = \tanh(\omega x)$$

$$f_2(x_1, x_2) = \sigma(\omega_1 x_1 + \omega_2 x_2)$$

$$f_3(x_1, x_2, x_3) = 1 \cdot (\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3)$$

## Forward Pass:

$$\hat{y} = f_3(f_2[f_1[x], \dots], \dots, \dots)$$

$$\hat{y} = 1 \cdot (\omega_1^3 (\sigma[\omega_1^1 [\tanh(\omega_1 x)] + \omega_2^1 \dots]) + \omega_2^2 \dots + \omega_3^2 \dots)$$

## Backward Pass:

Consider Loss function as MSE

$$L = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

# Use Local Gradients to find Weight correction factors and further

Local Gradients for each Neuron and its connection.

To simplify, our goal is to find every partial derivative by each model parameter (weight).

## For $f_3$ :

$$\frac{\partial L}{\partial \omega_{11}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_3} \cdot \frac{\partial f_3}{\partial \omega_{11}^2} =$$

$$= \underbrace{-\frac{2}{n}(y - \hat{y})}_{L'} \cdot \underbrace{1}_{f'_3 \text{ (Activation)}} \cdot \underbrace{f_2(f_1[x], \dots)}_{\text{unweight previous result from } \omega_{11}^2}$$

For output Layer:

$$\text{Local Grad}(N_n^m) = L' \cdot f'_n$$

$$\omega_{\text{new}} = \omega_{\text{old}} - \eta \cdot \text{Local Grad}(N_n^m) \cdot \text{res}(w)$$

## For $f_2$ :

$$\frac{\partial L}{\partial \omega_{11}^2} = \sum_{i=1}^n (\omega_i \cdot \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_3}) \cdot f'_2 \cdot \frac{\partial f_2}{\partial \omega_{11}^2} =$$

$$= \underbrace{\sum_{i=1}^n (\omega_i \cdot \text{Local Grad}(N_i^3))}_{\text{Sum of weighted previous Local gradients}} \cdot \underbrace{\sigma(\dots)(1 - \sigma(\dots))}_{f'_2 \text{ (Activation)}} \cdot \underbrace{f_1(x)}_{\text{unweight prev. result from } \omega_{11}^2}$$

For Hidden Layers:

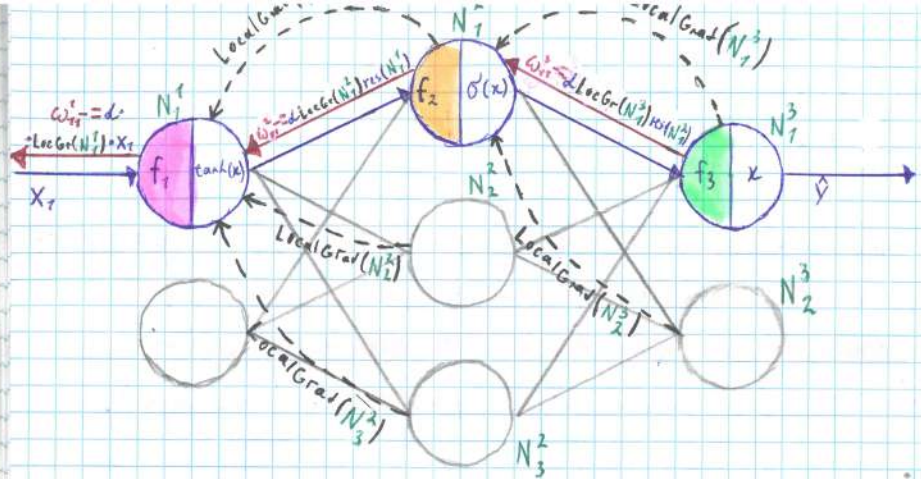
$$\text{Local Grad}(N_n^{m-1}) = \sum_{i=1}^n (\omega_i \cdot \text{Prev Local Grad}(N_i^m)) \cdot f'_{n-1}$$

## For $f_1$ :

$$\frac{\partial L}{\partial \omega_1} = \sum_{i=1}^n (\omega_i \cdot \sum_{j=1}^n (\omega_j \cdot \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_3}) \cdot f'_{2j}) \cdot f'_1 \cdot \frac{\partial f_1}{\partial \omega_1} =$$

$$= \underbrace{\sum_{i=1}^n (\omega_i \cdot \text{Local Grad}(N_i^2))}_{\text{Sum of weighted previous Local gradients}} \cdot \underbrace{(1 - \tanh^2(\dots))}_{f'_1 \text{ (Activation)}} \cdot \underbrace{X_1}_{\text{unweight prev. res from } \omega_1} - \text{Global Gradient}$$







# CNN

## Convolutional Neural Networks

### Idea:

It is a NN that contains special convolution layers.  
Mostly it is used for processing images.

### Convolution 2D:

An operation that is used to extract some locational information into Feature maps using kernels (filters).

**Kernel:** square matrix  $n \times n$  that represent a filter to "Detect" some patterns in image. Mostly, it consists of parameters (weights) which automatically change in train loop.

### Convolution:

$$\text{Stride} = 2; \text{Kernel} = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}_{3 \times 3}$$

Image  $7 \times 7$

1	9	8	4	9	6	7
4	8	6	7	5	1	7
4	0	5	9	3	8	9
7	3	6	9	0	5	4
7	4	1	1	8	1	2
7	6	6	9	8	7	6
3	6	3	5	4	2	7

$$\text{Out}_{11} = \text{Sum} \left( \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \odot \begin{bmatrix} 1 & 9 & 8 \\ 4 & 8 & 6 \\ 4 & 0 & 5 \end{bmatrix} \right)$$

$$\text{Out}_{12} = \text{Sum} \left( \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \odot \begin{bmatrix} 9 & 6 & 7 \\ 8 & 5 & 1 \\ 9 & 3 & 8 \end{bmatrix} \right)$$

**Stride:** amount of pixels on which kernel would be shifted during Convolution.

Output  
↓

21	7	-24
-7	30	12
7	25	18

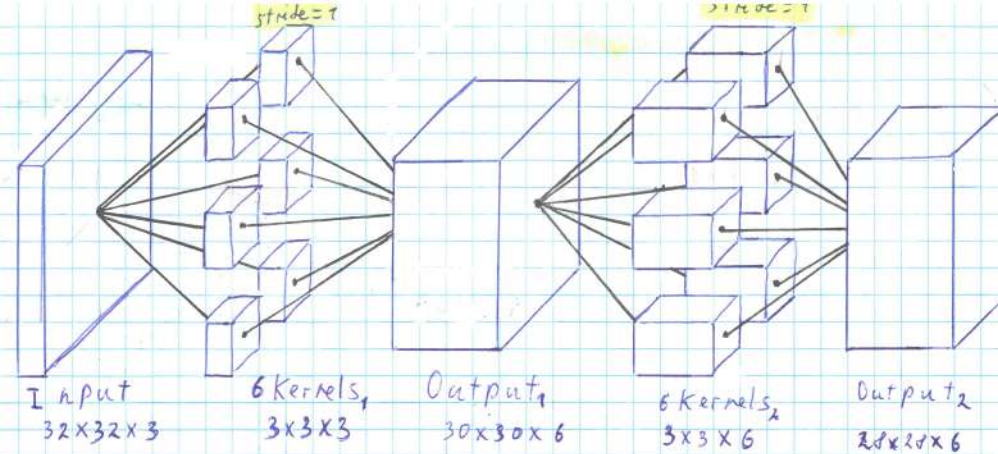
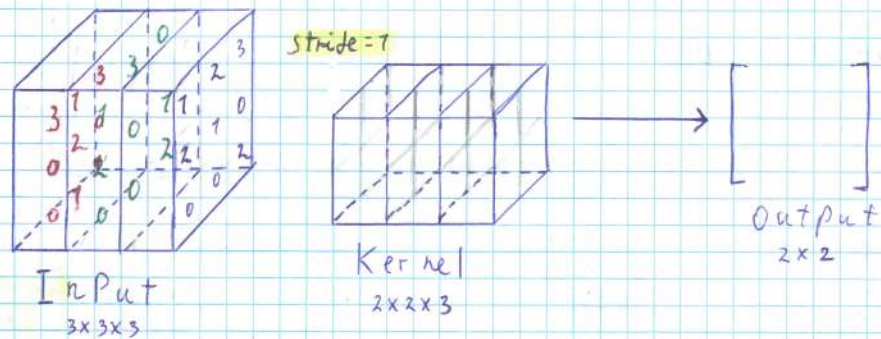
**Padding:** add frame of mean value of pixels to the image. Increases output image size after convolution.

**Feature Map:** Simply Shows if Kernel was activated (When Feature Map numbers are high with  $\text{Abs}()$ ) or not (When FM,  $\text{Abs}(\text{mins})$  are low). The one represents if there is a kernel pattern on the image.



## Convolution 3D:

This one is the same as 2D but the Kernel and input are now Tensors and feature map is still matrix.



## Pooling:

A technique of size reduction

of feature map By splitting matrix into pieces and

getting „Max“ or „Average“ or „Sum“ value. Used to

decrease sensitivity of next layers. Helps to recognize objects in different positions and sizes.

$$\begin{bmatrix} 1 & 7 & 2 & 6 \\ 2 & 4 & 0 & 2 \\ 1 & 4 & 0 & 8 \\ 0 & 0 & 4 & 9 \end{bmatrix} \xrightarrow{\text{Max Pooling}} \begin{bmatrix} 7 & 6 \\ 4 & 9 \end{bmatrix}$$

## Multiple Convolutions:

Consider convolutional layer as multiple kernels applied to single input and then stacking the results into one tensor. So this tensor could be the input to other convolutional layer with kernels of depth = stacked results.



# Alex Net as example :

