

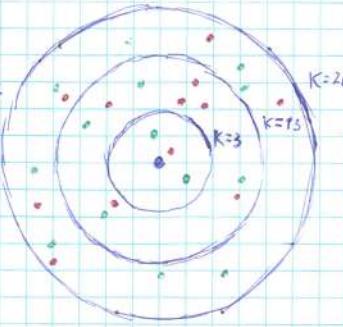
# Machine Learning

## Types of ML Problems:

- **Unsupervised Learning** - the machine did not use marked down factor. Could find Hidden patterns.
- **Supervised Learning** - the machine use marked down factor as training one. Could predict output due to previous training factor 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 Neighbours)

Ideas: Consider input Data as point   
(vector) 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 Dispersity 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:  
 $[0 \dots 2000] \text{ (m)}, [0 \dots 2] \text{ (km)} \rightarrow 2 \text{ m} \sim 2 \text{ km Wrong!}$

### Distance Metrics:

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

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

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

# Learning

## Accuracy Metrics:

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

### Classifier:

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

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

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

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

### Regressor:

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

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

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

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

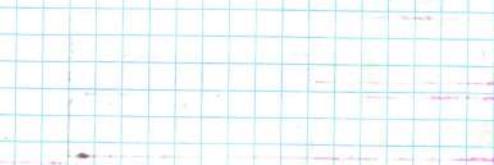
It is a function that shows how did prediction "y" differs from correct answer " $\hat{y}$ ". The lesser it is, the better predicted result you have.

## Overfitting:

Situation when model loss on training factor  $\rightarrow 0$  and loss on real factor  $\rightarrow \infty$ .

## Split Data:

all available data				
Training	Validation	Testing		



all available data					
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Testing
1	2	3	4	5	
split 1	1	2	3	4	5
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

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

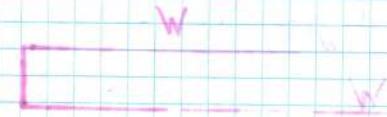
A method of adding special constraints to the loss function to decrease model complexity and prevent overfitting.

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

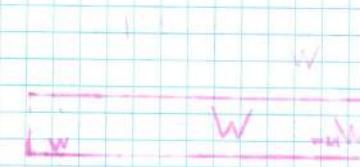
$$L_2 (\text{Ridge}) = \dots + \beta \sum_{i=1}^m w_i^2$$

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

$w$   $w$



$w$   $w$



$w$

$w$

$w$

$w$

$w$

## INPUT:

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & x_{21} & \dots & x_{2n} \\ 1 & x_{31} & \dots & x_{3n} \\ 1 & x_{41} & \dots & x_{4n} \\ \vdots & \vdots & & \vdots \\ 1 & x_{m1}, \dots & \dots & x_{mn} \end{bmatrix}$$

features

## Linear Regression

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

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  $\frac{\partial L}{\partial W} = 0$

## Gradient descent:

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

$$\frac{\partial L}{\partial W} = \frac{1}{m} X^T (XW - Y) = 0$$

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

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

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

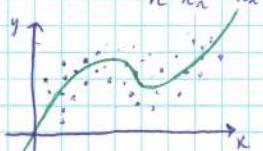
$\alpha$  - Learning Rate

## If dependencies are Non-Linear:

Simply add more Polynomial features and normalize them.

$$X = \begin{bmatrix} 1 & x_{11} \\ \vdots & \vdots \\ 1 & x_{n1} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$X' = \begin{bmatrix} 1 & x_{11} & x_{11}^2 & x_{11}^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n1}^2 & x_{n1}^3 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$



## 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.

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

Individuals

$$W = \begin{bmatrix} w_{00} & w_{01} & \dots & w_{0n} \\ w_{10} & w_{11} & \dots & w_{1n} \\ \vdots & \vdots & & \vdots \\ w_{K0} & w_{K1} & \dots & w_{Kn} \end{bmatrix}$$

for each feature  
for each class

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

Prob. of Belong to each class

## Train:

$$\text{logits}(X_{\text{individual}}) = X_{\text{individual}} \cdot W \sim [1 \ x_1 \ x_2 \dots x_n] \cdot \begin{bmatrix} w_{00} & \dots & w_{0n} \\ w_{10} & \dots & w_{1n} \\ \vdots & \vdots & \vdots \\ w_{K0} & \dots & w_{Kn} \end{bmatrix}$$

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

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

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 7.7 \end{bmatrix} \rightarrow \text{SoftMax} \rightarrow \begin{bmatrix} 0.02 \\ 0.9 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

## Loss:

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

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

# Log Loss - return huge value if prediction  $\hat{y}_i$  complies mismatch the correct class answer  $y_i$

## Gradient descent:

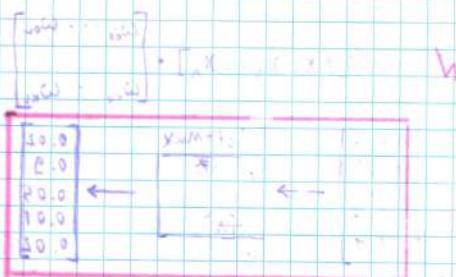
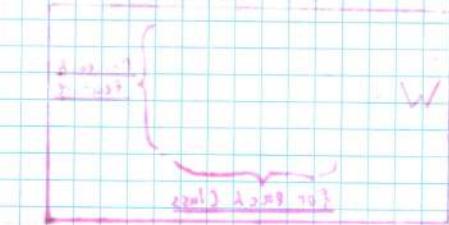
$$\frac{\partial L}{\partial W} = \frac{1}{n} \cdot X_{\text{individual}}^T (\hat{Y} - Y)$$

## DECISION TREES.

Idea:

it's 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 can't predict the result by simple tree traverse.



AT

CT

?

?

?

Training:

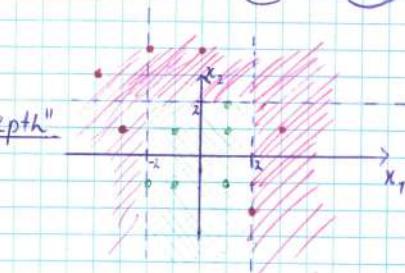
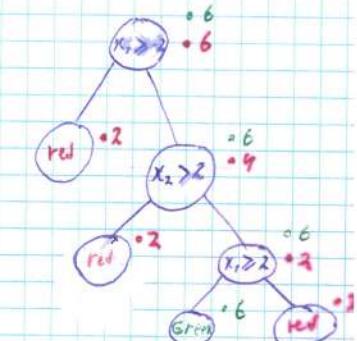
- Check all the leafs and calculate the Entropy or Gini Index for every possible combination of every feature value split ( $x_1 \leq 1, x_1 \geq 3 \dots$ )

- Calculate Information Gain for each possible split By :

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

Where " $w_i$ " =  $\frac{\text{size(child}_i)}{\text{size(parent)}}$

- Choose split with the highest IG value and repeat until "max\_tree\_depth" or every final note become leaf.



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

# where " $P(x_i)$ " = probability of  
 $P(\text{class}) = \frac{\text{Total objects}}{\text{class amount}}$  a class.

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

## Ensembles.

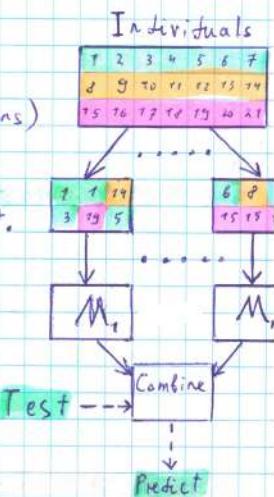
### Voting:

Use different algorithms with Prediction probability 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:

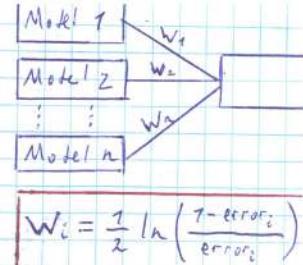
This is a variation of bagging technique where we randomly sample not only individuals but their features too.

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.



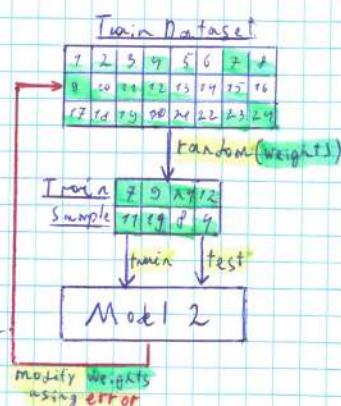
## Gradient Boosting:

## 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.

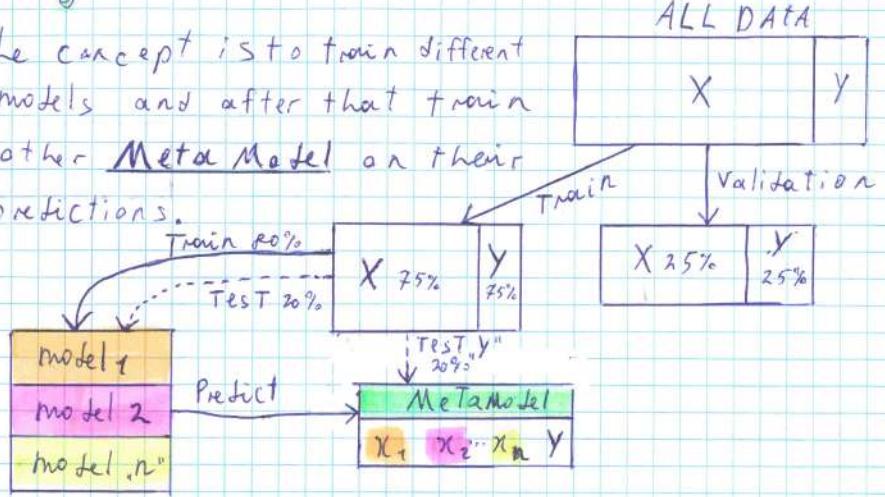
## Adaptive Boosting (AdaBoost):

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



## Blending:

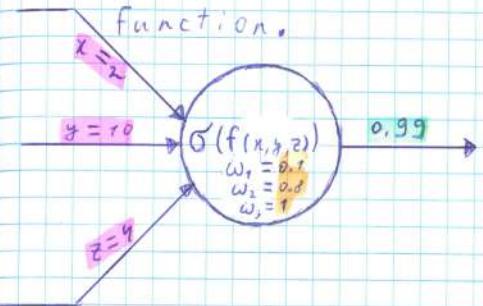
The concept is to train different models and after that train other Meta Model on their predictions.



## Stacking:

Neuron:

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



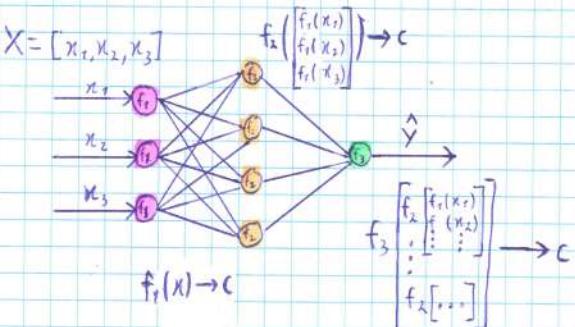
$$f(x, y, z) = w_1 \cdot x + w_2 \cdot \frac{y}{2} + w_3 \cdot z$$

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

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

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 functions into one. Mostly, it also normalize data and should be differentiable.

# Consider sequential linear functions without activation

$$\begin{array}{c} \boxed{\omega_1(k_1x + b_1)} \rightarrow \boxed{\omega_2(k_2x + b_2)} \rightarrow \boxed{\omega_3(k_3x + b_3)} \\ f_1(x) \qquad \qquad f_2(x) \qquad \qquad f_3(x) \end{array}$$

$$f_3(f_2(f_1(x))) = \omega_3(k_3(\omega_2(k_2(\omega_1(k_1x + b_1)) + b_2)) + b_3) \sim \boxed{\omega(kx + b)}$$

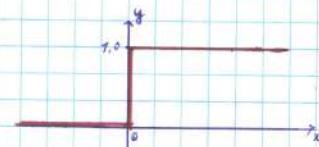
# Same but using activation:

$$\sigma(f_3(\sigma(f_2(\dots)))) = \sigma(\omega_3(k_3(\sigma(\omega_2(k_2(\dots + b_2)) + b_3))) \rightarrow \text{could not be simplified because "}\sigma\text{"-function!}$$

Binary Step:

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

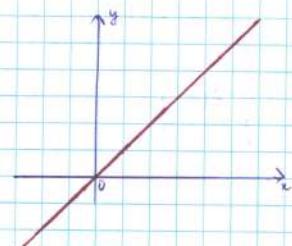
$$f'(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ \text{none} & \text{if } x = 0 \end{cases}$$



Linear:

$$f(x) = x$$

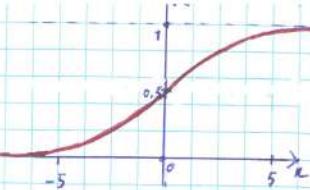
$$f'(x) = 1$$



Sigmoid (Logistic).

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

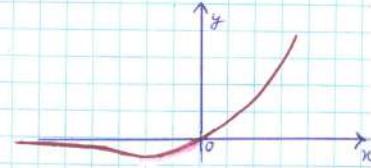
$$f'(x) = f(x) \cdot (1 - f(x))$$



Swish:

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

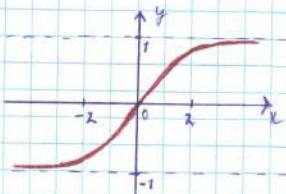
$$f'(x) = f(x) + \text{sigmoid}(x) \cdot (1 - f(x))$$



Tanh:

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

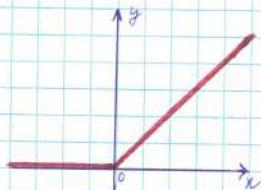
$$f'(x) = 1 - f(x)^2$$



ReLU (Rectified Linear Unit):

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

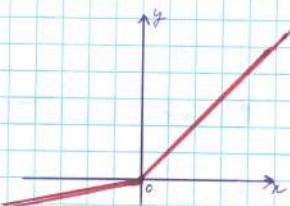
$$f'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Leaky ReLU:

$$f(x) = \max(0.01x; x) \sim \begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

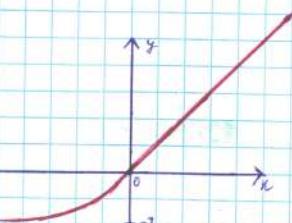
$$f'(x) = \begin{cases} 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



ELU (Exponential Linear Unit):

$$f(x) = \begin{cases} 0.01e^x - 1 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} f(x) + 0.01 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

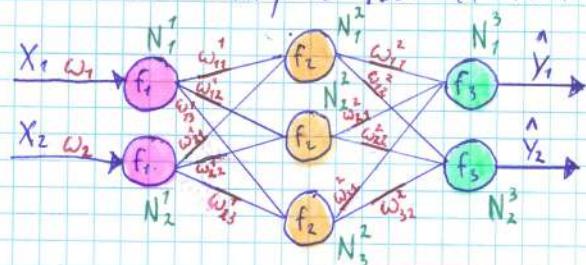


## Idea:

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

## Example:

Consider Simple Neural Network:



$$f_1(x) = \tanh(w_1 x)$$

$$f_2(x_1, x_2) = \sigma(w_1 x_1 + w_2 x_2)$$

$$f_3(x_1, x_2, x_3) = 1 \cdot (w_1 x_1 + w_2 x_2 + w_3 x_3)$$

## Forward Pass:

$$\hat{Y} = f_3(f_2[f_1[x_1, \dots], \dots, \dots])$$

$$\hat{Y} = 1 \cdot (w_1^2 (\sigma[w_1 [\tanh(w_1 x)] + w_2^1 \dots] + w_2^2 \dots + w_3^2 \dots])$$

## Backward Pass:

Consider Loss function as MSE

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

# Use LocalGradients to find weight correction factors and further

LocalGradients for each Neuron and its connection.

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

for  $f_3$ :

$$\begin{aligned} \frac{\partial L}{\partial w_{11}^2} &= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_3} \cdot \frac{\partial f_3}{\partial w_{11}^2} = \\ &= \underbrace{-\frac{2}{m}(y - \hat{y})}_{L'} \cdot \underbrace{1}_{f_3'} \cdot \underbrace{f_2(f_1[x], \dots)}_{\text{unweight previous result from } w_{11}^2} \end{aligned}$$

For output Layer:

$$\text{LocalGrad}(N_m^n) = L' \cdot f_n'$$

$$w_{\text{new}} = w_{\text{old}} - d \cdot \text{LocalGrad}(N_m^n) \cdot \text{res}(w)$$

for  $f_2$ :

$$\begin{aligned} \frac{\partial L}{\partial w_{11}^1} &= \sum_{i=1}^m \left( w_i \cdot \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_{3i}} \right) \cdot f_2' \cdot \frac{\partial f_2}{\partial w_{11}^1} = \end{aligned}$$

$$= \sum_{i=1}^m \left( w_i \cdot \text{LocalGrad}(N_i^3) \cdot \underbrace{\sigma(\dots)(1 - \sigma(\dots))}_{f_2'} \cdot f_1(x) \right) \cdot \underbrace{\text{sum of weighted previous local gradients, result from } w_{11}^1}_{\text{Activation}}$$

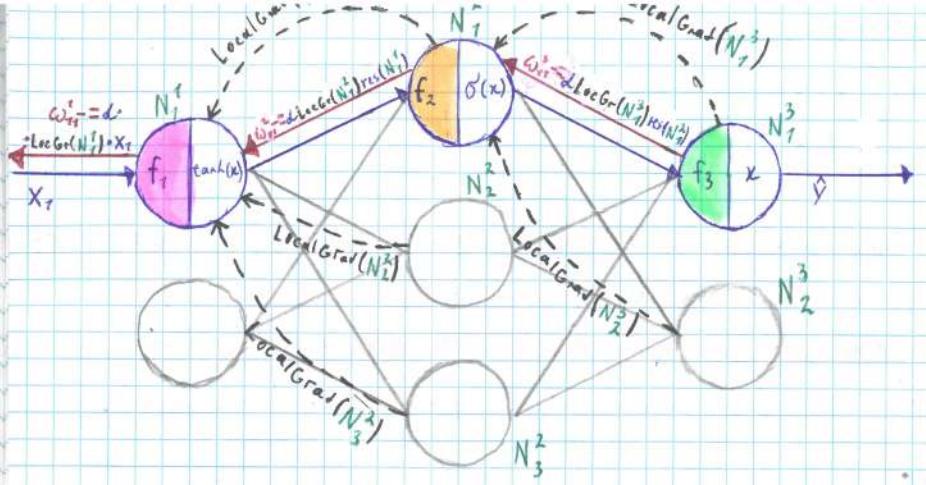
For Hidden Layers:

$$\text{LocalGrad}(N_m^{n-1}) = \sum_{i=1}^m (w_i \cdot \text{PrevLocalGrad}(N_i^n) \cdot f_{n-1}')$$

for  $f_1$ :

$$\frac{\partial L}{\partial x_1} = \sum_{i=1}^m \left( w_i \cdot \sum_{j=1}^2 \left( w_j \cdot \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial f_{3i}} \right) \cdot f_{1i}' \right) \cdot f_1' \cdot \frac{\partial f_1}{\partial w_1} =$$

$$= \sum_{i=1}^m \left( w_i \cdot \text{LocalGrad}(N_i^2) \cdot \underbrace{1 - \tanh^2(\dots)}_{f_1'} \cdot x_1 \right) \cdot \underbrace{\text{sum of weighted previous local gradients, result from } w_1}_{\text{Activation}} - \text{Global Gradient}$$



$$\left( \frac{x_1}{x_2}, \frac{x_2}{x_1} \right)$$

$$\left( \frac{x_1}{x_2}, \frac{x_2}{x_1} \right) \text{ hat wert} =$$

$$(\frac{x_1}{x_2})^2$$

$$\frac{x_1}{x_2}$$

$$\frac{x_2}{x_1}$$

# CNN

## Convolutional Neural Networks

Idea:

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

Convolution 2D:

An operation that 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:

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

Image  $2 \times 7$

$$\text{Out}_{11} = \text{Sum} \left( \begin{bmatrix} 1 & 9 & 8 & 4 & 9 & 6 & 7 \\ 4 & 8 & 6 & 7 & 9 & 1 & 7 \\ 4 & 0 & 5 & 9 & 3 & 8 & 9 \end{bmatrix} \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \right)$$

$$\text{Out}_{12} = \text{Sum} \left( \begin{bmatrix} 7 & 3 & 5 & 9 & 0 & 5 & 4 \\ 7 & 4 & 7 & 1 & 8 & 1 & 2 \\ 7 & 6 & 6 & 9 & 8 & 7 & 6 \\ 3 & 6 & 3 & 5 & 4 & 2 & 7 \end{bmatrix} \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} \right)$$

Output

$$\begin{bmatrix} 21 & 7 & -29 \\ -7 & 30 & 12 \\ 7 & 25 & 78 \end{bmatrix}$$

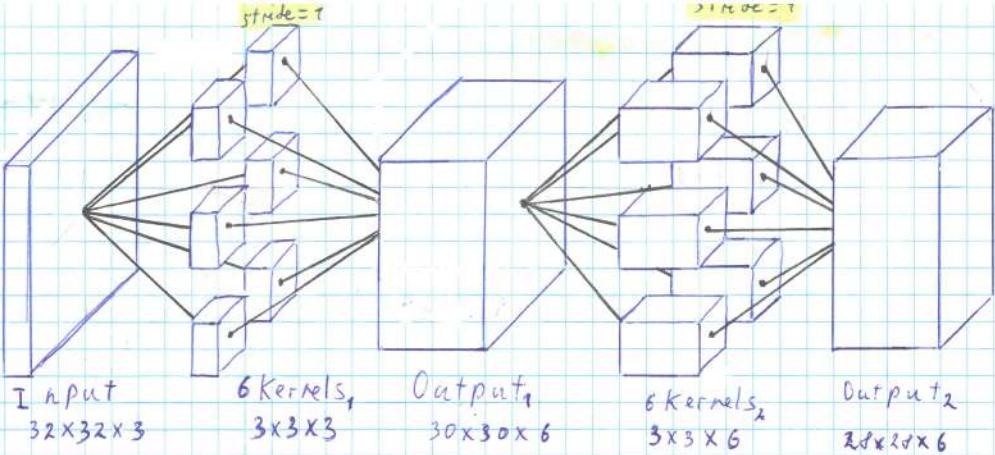
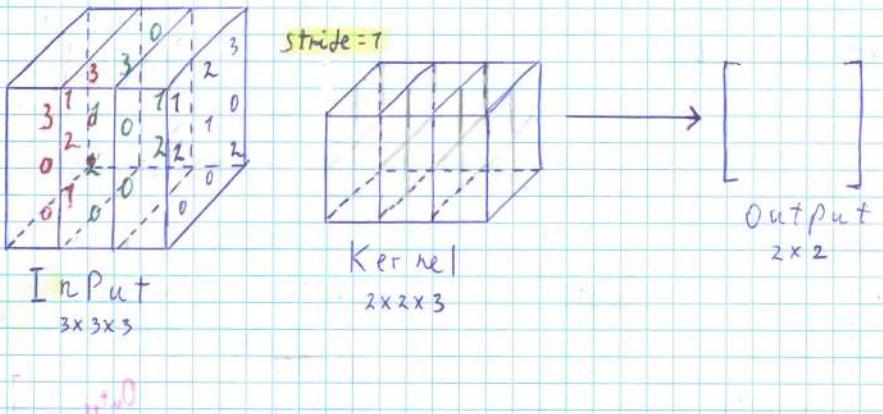
Stride: amount of pixels on which kernel would be shifted during Convolution.

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.  $\text{Abs}(\text{num})$  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:

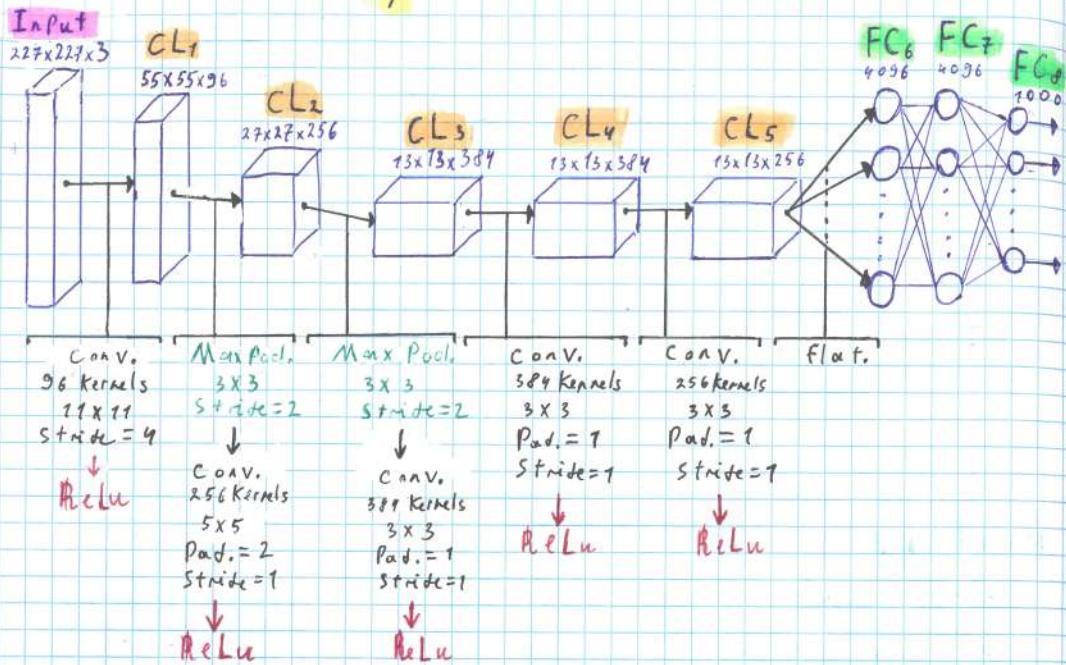
A technique of size reduction of feature map By splitting matrix into pieces and getting "Max" or "Average" or "Sum" value. Used to decrease sensitivity of next layers. Helps to recognize objects in different positions and sizes.

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

## 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.

# AlexNet as example:



## Vanishing Gradient

### Problem:

Consider NN have huge amount of layers ("Deep") and it uses Back propagation to train. So when we want to calculate weights:

$$w_0 = w_0 - \alpha \frac{\partial L}{\partial w_0} \dots w_n = w_n - \alpha \frac{\partial L}{\partial w_n}$$

$$\frac{\partial L}{\partial w_0} = \underbrace{\frac{\partial L}{\partial f_n} \cdot \frac{\partial f_n}{\partial w_0}}_{n \text{ multipliers less than 1.0}}$$

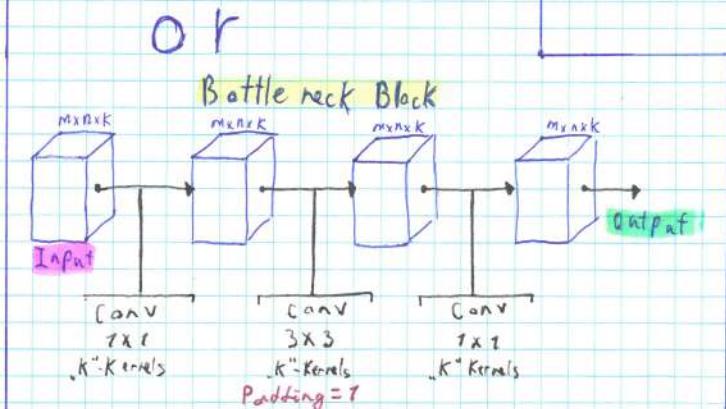
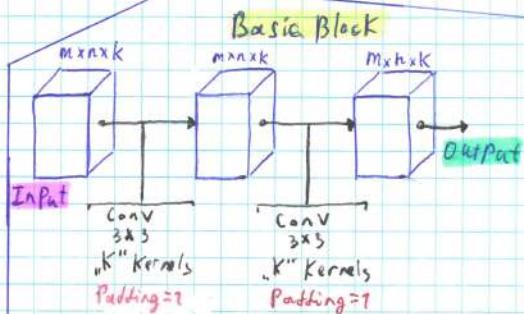
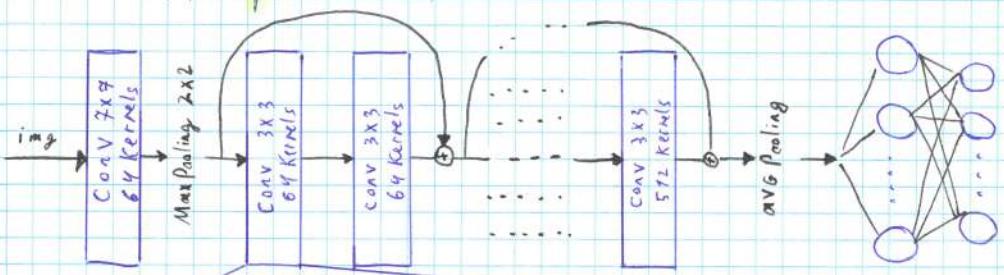
$$\frac{\partial L}{\partial w_n} = \underbrace{\frac{\partial L}{\partial f_n} \cdot \frac{\partial f_n}{\partial w_n}}_{2 \text{ multipliers less than 1.0}}$$

So  $\alpha \frac{\partial L}{\partial w_0} \rightarrow 0$ , That means that weight "w<sub>0</sub>" will not change so much. So the first layers won't learn.

have the same Dimensionality as an output tensor. To do it, output of a "Block" should have width and height same as "X".

# use Conv("Kernel=[1x1]; amount=n) to increase or decrease Depth of a tensor "Identity."

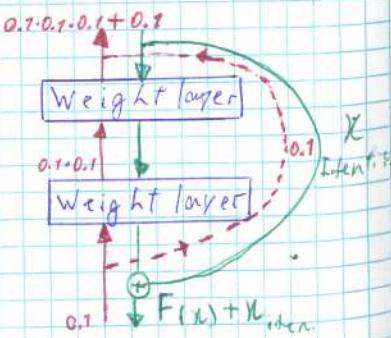
### Res Net Example:



### Skip Connection

In this method we simply connect and push signal from upper layers to lower ones.

Also It pushes high level information to deeper layers.



- In CNN: To perform such operation X\_identity should

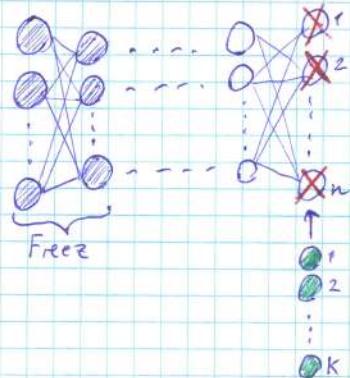
## Transfer Learning

The Idea is to somehow transfer "knowledge" from one network to another.

### Fine tuning:

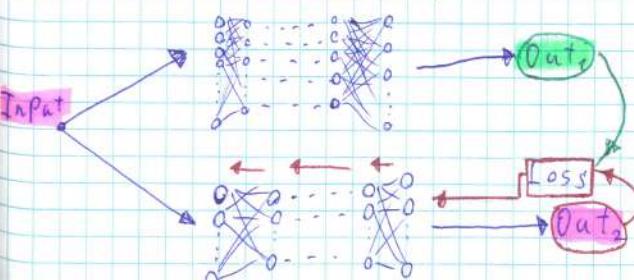
This method is used when you have an already trained model which solves fundamental problem. But you want to concretize or narrow down the problem.

- Copy trained model
- Drop down Output layer and replace with correct one.
- Freeze some (or all) first layers.
- Re-train it on new Dataset.



### Black Box Learning:

Uses if there exist huge already trained Neural Net and you want to get similar result or approximate it by smaller Network.



# SPIKING Neural Networks (SNN)

It is special type of ANN which principle of work is based on neuro morphic concept.

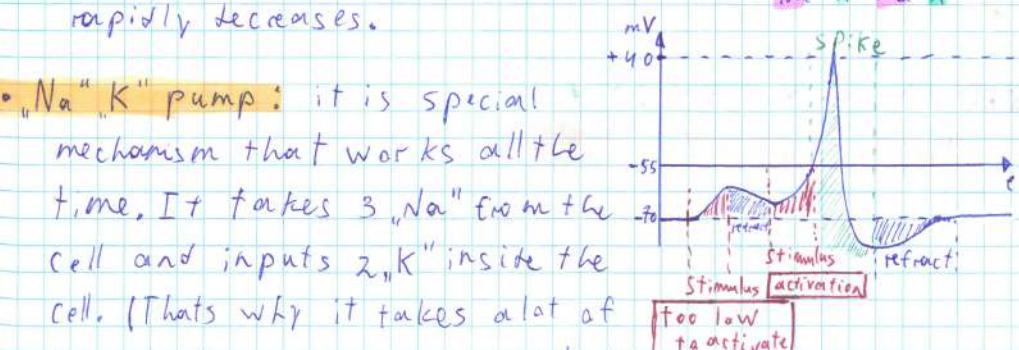
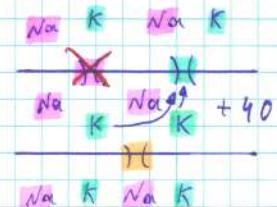
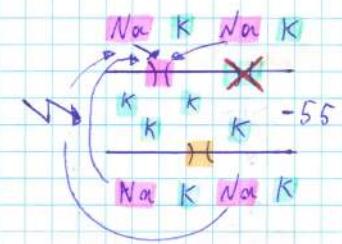
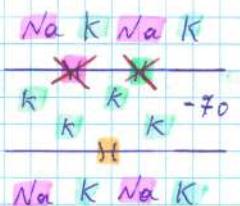
## Biological Concept:

• Firstly consider that every neuron cell is charged negatively due to small amount of " $K^+$ " ions and absence of " $Na^+$ " ions ( $-70\text{mV}$ ).

• Stimulus: when some comes, it increases the potential of cell. So when the potential reaches the threshold of ( $-55\text{mV}$ ) it activates " $Na^+$ " ions channels (Voltage Gates)

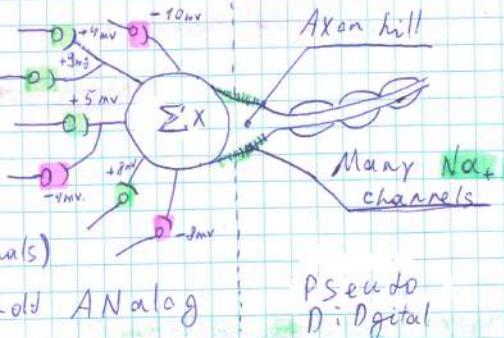
• Spike: when potential reaches ( $+40\text{mV}$ ) " $Na^+$ " ions channels closes and " $K^+$ " ions channels opens so the potential rapidly decreases.

• " $Na^+$ " " $K^+$ " pump: it is special mechanism that works all the time. It takes 3 " $Na^+$ " from the cell and inputs 2 " $K^+$ " inside the cell. (That's why it takes a lot of time to decrease or increase potential to the stable ( $-70\text{mV}$ )).



## Signal Processing:

When signals (current from synapses) comes to the soma, it changes the membrane potential (Actually soma simply sums all the signals) and when it reaches the threshold Analog Spike comes from Axon hillock.

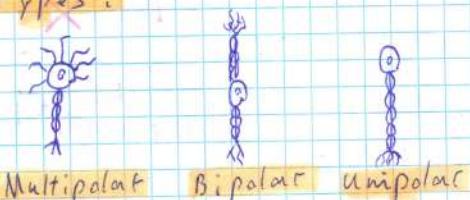


## Neuron Cells:

Dendrites - Inputs

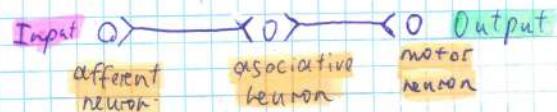
Axons - Output

### Types:



Pseudounipolar

No axon.



## SNN Data Encoding:

The input and output format should be encoded as a sequence of spikes.

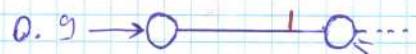
### • Rate Coding:

Corresponding input Neuron return spikes more frequently if the input is larger.



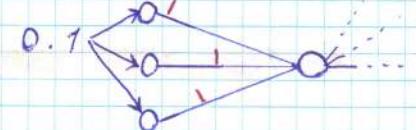
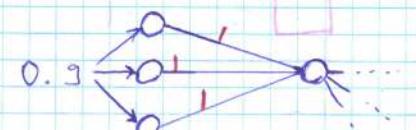
### • Temporal Coding:

Neuron with the largest input value fires first.



### • Population Coding:

Several input neurons (on one feature) to create spikes with different delay. So the second layer gets the sum of signals.



## • Leaky Integrate and Fire (LIF).

The idea is the same as Integrate and Fire model  
But now the potential ( $V$ - voltage) always decreases exponentially.

$$C \cdot \frac{dV}{dt} = -\text{leak} + I$$

exponential decrease  
of Voltage (always).

Where:  $\text{leak} = \frac{C}{T} (V_{(t)} - V_0)$

If spike comes  
from other neurons  
increase voltage on  
some "spike" constant.

Increase or slow  
down leakage.  
Used to synchronize  
system frequency.

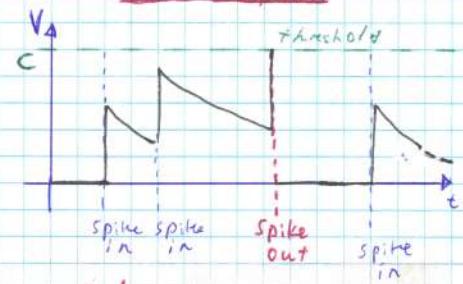
By substituting:

$$\frac{dV}{dt} = -\frac{1}{C} (V - V_0) + I$$

In other Interpretation:

$$\left. \begin{array}{l} \text{Always: } \frac{dV}{dt} = -\frac{1}{C} (V - V_0) \\ \text{If received } I : V = V + I \end{array} \right\} \begin{array}{l} \leftarrow \text{Leak} \\ \leftarrow \text{Integrate} \end{array}$$

$$\left. \begin{array}{l} \text{If } V > \text{Threshold} : V = 0 \end{array} \right\} \begin{array}{l} \leftarrow \text{Fire (non linearity)} \end{array}$$



## Izhikevich Model:

The concept is to "u(t)" function, that adaptively Decrease the Voltage per spike if there were multiple spikes in one before.

$$C \frac{dV}{dt} = 0.04V^2 + 5V + 140 + I - u$$

If input spike comes, the last voltage.

Decrease the Current per input spike.  
(recovery Variable)

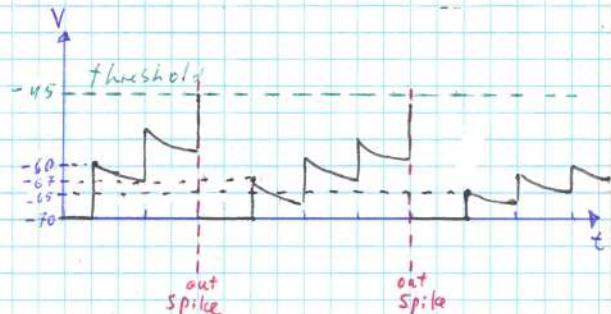
$$\frac{du}{dt} = \alpha(bV - u)$$

If  $V \geq \text{Threshold}$  :  $V = \dots$ ,  $u += d$

$\alpha$  - timescale of "u".

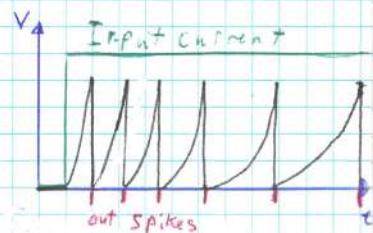
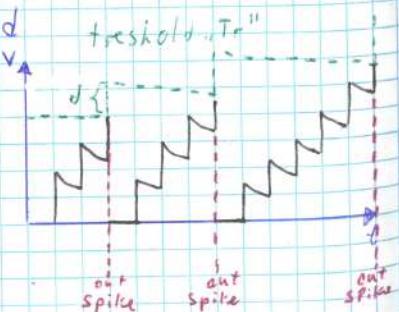
Smaller Values leads to slower recovery (default  $\alpha=0.02$ )

$b$  - Sensitivity of "u"



$d$  - After Spike reset of "u".

(By Default  $d=2$ )



General View.

## Adaptive (2D) LIF:

The same as basic LIF But the threshold is now dynamical " $T_r$ " (Homeostasis analog)

$$\frac{dV}{dt} = -\frac{1}{C}(V - V_0) + I$$

$$\frac{dT_r}{dt} = \frac{1}{\alpha} (1 - T_r)$$

If  $V \geq T_r$  then:  $V=0$ ,  $T_r += d$

## Hodgkin - Huxley Model:

This is one of the most precise and heavily computable model of neuron. It simulates the Ion channel process.

$$C \frac{dV}{dt} = I_K + I_{Na} + I_L$$

"K" Ion Current  
 "Na" Ion Current  
 "Na" "K"  
 channel current channel current pump

$$C \frac{dV}{dt} = g_K \cdot (V_K - V) + g_{Na} \cdot (V_{Na} - V) + g_L \cdot (V_L - V)$$

"K" Gating Variable  
 "Na" Gating Variable  
 leak Gating Variable

# Gating variables "g(V, t)" are functions that regulates when to open or close corresponding Ion channel.

$$C \frac{dV}{dt} = \bar{g}_K \cdot f_K (V_K - V) + \bar{g}_{Na} \cdot f_{Na} \cdot h_{Na} (V_{Na} - V) + \bar{g}_L (V_L - V)$$

4 Gates in K channel  
 Max. conductance for "K"  
 3 Gates in Na channel  
 Maximum activation conductivity (open) for "Na"  
 Deactivation (close) function  
 Constant leak

$$f_K = \frac{df_K}{dt} = \alpha_K (1 - f_K) - \beta_K f_K \quad \text{where: } \alpha_K = \frac{0.01(10 - V)}{\exp(\frac{10 - V}{10}) - 1}$$

$$\beta_K = 0.125 \cdot \exp(-\frac{V}{10})$$

$$f_{Na} = \frac{df_{Na}}{dt} = \alpha_{Na} (1 - f_{Na}) - \beta_{Na} f_{Na} \quad \text{where: } \alpha_{Na} = \frac{0.1 (25 - V)}{\exp(\frac{25 - V}{10}) - 1}$$

$$\beta_{Na} = 4 \cdot \exp(-\frac{V}{10})$$

$$h_{Na} = \frac{dh_{Na}}{dt} = \alpha_h (1 - h_{Na}) - \beta_h h_{Na} \quad \text{where: } \alpha_h = 0.075 \cdot \exp(-\frac{V}{20})$$

# Where "α" - Open Coefficient and "β" - Close Coefficient.

$$\beta_h = \frac{1}{\exp(\frac{30 - V}{10}) + 1}$$

## Training:

SNN's could not be trained by using basic methods such as backpropagation because they operate with spikes

## Spike Timing dependent Plasticity (STDP):

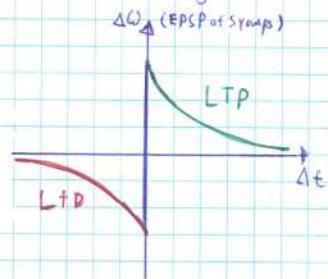
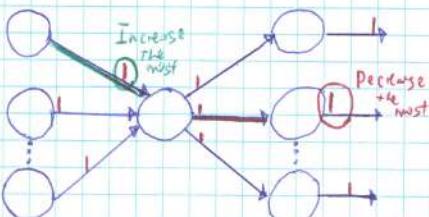
It is an unsupervised learning algorithm that operates with only two basic rules:

- If PreSynaptic Neuron Spikes and after that post synaptic; then increase this weight.
- If PostSynaptic Neuron Spikes and then presynaptic; then Decrease this weight.

#The one is probably increasing synchronous activity.

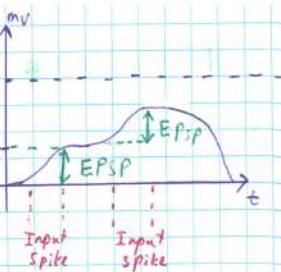
## Detailed Realization:

When the neuron Spikes, we check all its input and output connections, get " $\Delta t$ " time interval in which other neurons spike. Now using our function Decrease all output weights and increase input

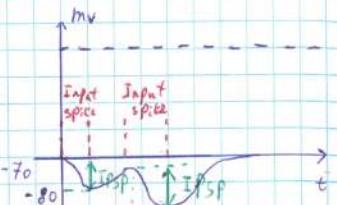


$$w_{\text{new}} = w + \Delta w \begin{cases} \text{if presynaptic: } \Delta w = A_p e^{-\frac{\Delta t}{\tau}} \\ \text{if postsynaptic: } \Delta w = A_{ps} e^{-\frac{\Delta t}{\tau}} \\ \text{else: } \Delta w = 0 \end{cases}$$

**EPSP:** How much synapse would amplify the input signal. The more it is, the faster neuron would fire. (Excitatory Postsynaptic Potential.)



**IPSP:** How much synapse would decrease (Inhibit) signal. The higher it is, the more unreachable is threshold. (Inhibitory postsynaptic Potential)



### • Homeostatic STDP:

The idea is to make STDP more stable in recurrent cases (all loops would infinitely increase their weights - "Hebbian rule").

### Synaptic Normalization:

Each neuron have Synaptical resource:

$$w_i = w_{\min} + \frac{(w_{\max} - w_{\min})}{n}$$

## Biological Concepts:

### Types of synapses:

- Electrical
- Chemical
- Axo Dendritic
- Axo Somatic
- Axo Axonic.

- Excitatory
- Inhibitory
- Modulatory

### How Does This Work:

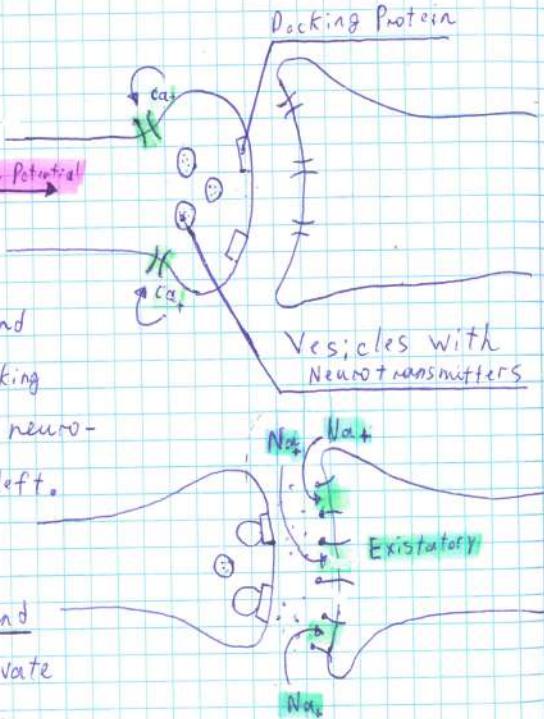
#### When the Action Potential

comes to synapse if Action Potential opens " $\text{Ca}^{2+}$ " Ion channels.

" $\alpha$ " connects to Vesicles and make them connect to the docking protein, which results the neurotransmitters are in Synaptic Cleft.

#### Depending on Neurotransmitters

Corresponding receptors (Ligand Gated Channels) would activate and increase (Excitatory) or decrease (Inhibitory) the membrane potential.



# Mostly, Neurotransmitters are synthesized in the Soma and transmitted by Motor Proteins to the synapses.