# Weekly report of lessons

**Name:** Mayank Kumar
**Roll No:** 19CS30029
**The week:** 18-10-2021 to 24-10-2021

## The topics covered:

Linear Discrimination: Discriminant Functions, Perceptron Classifier, Augmented Input and Linear Form, Linearly Separable Classes, Error function, Gradient descent method for iterative optimization, Other error functions, More stringent criteria of linear separability, Batch Relaxation Algorithm, Single sample relaxation with margin, Support Vector Machine: Two class problem, Optimization problem, Convex quadratic optimization problem, Solution, SVM Testing; Non-separable case, Slack variable to define soft margin, Constraints with soft margin hyperplanes, Projecting to higher dimensional space, Kernel machines, Vectorial kernel functions, Kernels, Parametric two class classification using discriminant, Logistic discrimination of two classes, Learning weights of logit functions, Algorithm (Learning weights)

ANN: Perceptron modelling a neuron, Artificial Neural Network, Feed-forward Network, Multilayered Feed-forward Network, Mathematical description, Input Output Relation, Optimization Problem, Chain rule of computing gradient of a single neuron, Computing gradient: Back propagation method, Back propagation: Concept, Back propagation: Delta Rule, ANN Training, Improving convergence

## Summary topic wise:

### Linear Discrimination

<u>Discriminant Functions:</u> Linear function: $g_i(x|w_i, w_{i0}) = \sum\limits_{j=1}^{d} w_{ij}x_j + w_{i0}$

Quadratic Function: $g_i(x|W_i, w_i, w_{i0}) = x^T W_i x + w_i^T x + w_{i0}$

<u>Perceptron Classifier:</u> It is a linear classifier with a different perspective. $y = f(z) = sign(z)$

<u>Augmented Input and Linear Form:</u> W is computed by minimizing classification error for given $\{y_i, X_i\}$.

$z = \Sigma_i w_i x_i + w_0 \Rightarrow z = W^T X \qquad o = f(z) = sign(z)$

<u>Linearly Separable Classes:</u> If a solution to find a hyperplane separating data points of two classes exists, then the classes are called linearly separable.

<u>Error function:</u> Error function (Perceptron Criterion) $J(W) = \sum\limits_{Y\ misclassified} - W^T Y$.

<u>Gradient descent method for iterative optimization:</u> Start with initial $W^{(0)}$ and update $W^{(i)}$ iteratively by computing $\nabla J(W^{(0)})$. $W^{(i)} = W^{(i-1)} - \eta(i)\nabla J(W)$, where $\eta$ is a positive scale factor (learning rate).

<u>Other error functions:</u> $J_q(W) = \sum\limits_{Y\ misclassified} (W^T Y)^2$, $J_r(W) = \sum\limits_{Y\ misclassified} \frac{(W^T Y - b)^2}{||Y||^2}$, $W^T Y \le b$

<u>More stringent criteria of linear separability:</u> The margin of separation between two linearly separable data points of classes is maximized using linear support vector machines.

<u>Batch Relaxation Algorithm:</u> W is initialized with $W^{(0)}$ and iterated till convergence. M is computed with margin b such that $M = \{Y|W^T Y <= b\}$, where M is the set of misclassified samples. Gradient

$\nabla J_r(W) = \sum\limits_{W^T Y \le b} \frac{Y(W^T Y - b)}{||Y||^2}$ and W is updated as $W^{(i)} = W^{(i-1)} - \eta\nabla J_r(W^{(i-1)})$

<u>Single sample relaxation with margin:</u> W is initialized with W(0) and is updated by considering samples one by one in every iteration. For sample $Y_i$ at iteration k, $W^{(k)} = W^{(k-1)} + \eta(k)\frac{b - W^T Y_i}{||Y_i||^2} Y_i$ if $W^T Y_i <= b$, and the algorithm stops when change in update at the end of an iteration is very small.

<u>Support Vector Machine (SVM):</u> It is a linear discriminant classifier that uses Vapnik's principle and computes class boundaries without computing class distributions.

<u>Two class problem:</u> For $X = \{x^t, r^t\}$, t = 1 to N, $x^t$ in $R^d$, $r^t$ in $\{+1, -1\}$, w and $w_0$ are computed such that, $r^t(w^T x^t + w_0) \ge 1 \ \forall\ t$

<u>Optimization problem:</u> Constrained Problem: To minimize $||w||^2/2$ subject to $r^t(w^T x^t + w_0) \ge 1 \ \forall\ t$

Unconstrained Problem: $L_p = \frac{1}{2}||w||^2 - \sum\limits_{t=1}^{N} \alpha^t [r^t(w^T x^t + w_0) - 1]$, where $\alpha^t$ is the Lagrange multipliers.

<u>Convex quadratic optimization problem:</u> Dual Problem: To be maximized wrt to $\alpha^t$ (>0) subject to zero gradients of $L_p$ wrt w and $w_0$. From this, we get $w = \Sigma_t \alpha^t r^t x^t$ and $\Sigma_t \alpha^t r^t = 0$

<u>Dual optimization problem:</u> To maximize $L_d$ wrt $\alpha^t$, $L_d = -\frac{1}{2}\Sigma_t \Sigma_s \alpha^t \alpha^s r^t r^s (x^t)^T x^s + \Sigma_t \alpha^t$

<u>Solution:</u> On applying quadratic optimization technique, we get $w_0 = r^t - w^T x^t$

<u>SVM Testing:</u> Here, margin is not enforced and only support vectors decide class boundaries.

<u>Non-separable case:</u> If classes are not linearly separable, use slack variable, $\{s^t\}$, t = 1 to N

<u>Slack variable to define soft margin:</u> For soft margin hyperplane, to maximize margin, minimize $\|w\|^2$ subject to $r^t(w^T x^t + w_0) \geq 1 - s^t \, \forall \, t$

<u>Constraints with soft margin hyperplanes:</u> Slack variable is used to define constraints.
- $0 < s^t < 1$, $x^t$ correctly classified; $s^t > 1$, $x^t$ misclassified
- Soft Error = $\Sigma_t s^t$

<u>Projecting to higher dimensional space:</u> It may make them linearly separable.
$$z = \varphi(x) \qquad g(z) = w^T z \qquad g(x) = \Sigma_j w_j \varphi_j(x), \text{ where } z_j = \varphi_j(x), j = 1 \text{ to } k$$

<u>Kernel machines:</u> Discriminant function $g(x) = w^T \varphi(x) = \Sigma_t \alpha^t r^t K(x^t, x)$

The matrix of kernel values K is called Gram matrix, where $K_{t,s} = K(x^t, x^s)$
<u>Vectorial kernel functions:</u> Polynomials of degree q, Radial basis functions, Mahalanobis kernel functions, Distance function based, Sigmoid function
<u>Kernels:</u> It may be defined between a pair of objects flexibly without using any closed functional form.

<u>Parametric two class classification using discriminant:</u> Parameters $\Sigma$, $\mu_1$, and $\mu_2$ are estimated to compute coefficients of g(x): $w$ and $w_0$.
<u>Logistic discrimination of two classes:</u> $y = P(C_1|x) = \dfrac{1}{1+e^{-(w^T x + w_0)}}$

<u>Learning weights of logit functions:</u> For $y = P(r^t=1|x)$, minimize E using gradient descent to iterate on weights, where $E = -\Sigma_t (r^t \log y^t + (1 - r^t)\log(1 - y^t))$
<u>Algorithm (Learning weights):</u> Assume initial w and $w_0$ and compute $y = sigmoid(w^Tx+w_0)$, compute gradients, and update w and $w_0$. Continue these computations till convergence.

## ANN

<u>Perceptron modelling a neuron:</u> $Input \rightarrow z = \Sigma_i w_i x_i + w_0, f(z) \rightarrow output$, $w_0$ is bias

<u>Artificial Neural Network:</u> It is a network of perceptrons that describe input/output relations
<u>Feed-forward Network:</u> Signals travel in one direction only, and there is no feed back or loop
<u>Multilayered Feed-forward Network:</u> Every unit in a layer is connected with all units in the previous layer, and the $i^{th}$ layer takes input from $(i-1)^{th}$ layer and forwards its output to the input of $(i+1)^{th}$ layer.
<u>Mathematical description:</u> $ne_j^{(i)} \rightarrow j^{th}$ neuron of $i^{th}$ layer, weights $W_j^{(i)} = (w_{j1}^{(i)}, w_{j2}^{(i)}, ..., w_{jn\_(i-1)}^{(i)})$, bias $w_{j0}^{(i)}$
$n\_(i-1)$: dimension of input to neuron, $n\_i$: dimension of output at $i^{th}$ layer

Output of $j^{th}$ neuron in $i^{th}$ layer: $y_j^{(i)} = f(W_j^{(i)^T} X^{(i-1)} + w_{j0}^{(i)})$
<u>Input Output Relation:</u> In $i^{th}$ layer, $Y^{(i)} = f(W^{(i)}X^{(i-1)} + b^{(i)})$

<u>Optimization Problem:</u> Given $X_i$ and $O_i$, i = 1 to N, to minimize $J_n(W) = \frac{1}{N}\sum_{i=1}^{N} \|O_i - F(X_i; W)\|^2$; $X_i$ is input that produces output $O_i$. Gradient descent procedure: Start with initial $W_0$ and update W iteratively.
$W_i = W_{i-1} + \eta(i) \Sigma_k (O_k - F(X_k; W_{i-1})) \nabla F(X_k; W_{i-1})$
<u>Stochastic gradient procedure:</u> $W_i = W_{i-1} + \eta(i) (O_k - F(X_k; W_{i-1})) \nabla F(X_k; W_{i-1})$

<u>Chain rule of computing gradient of a single neuron:</u> $\dfrac{\partial E}{\partial w_i} = \dfrac{\partial E}{\partial o}\dfrac{\partial o}{\partial z}\dfrac{\partial z}{\partial w_i}$ $\qquad \nabla(W) = (\dfrac{\partial E}{\partial w_0}, \dfrac{\partial E}{\partial w_1}, ..., \dfrac{\partial E}{\partial w_n})$

<u>Computing gradient: Back propagation method:</u> Chain rule is applied from output to input and from output layer to input layer, and partial derivatives of weights at $(i-1)^{th}$ layer from $i^{th}$ layer are computed.

<u>Back propagation: Concept:</u> $\Delta w_{31}^{(2)} = \dfrac{\partial E}{\partial w_{31}^{(2)}} = \dfrac{\partial E}{\partial o}(\dfrac{\partial o}{\partial y_1^{(3)}}\dfrac{\partial y_1^{(3)}}{\partial y_1^{(2)}} + \dfrac{\partial o}{\partial y_2^{(3)}}\dfrac{\partial y_2^{(3)}}{\partial y_1^{(2)}})\dfrac{\partial y_1^{(2)}}{\partial w_{31}^{(2)}} \qquad \dfrac{\partial y_1^{(2)}}{\partial w_{31}^{(2)}} = f'(z_1^{(2)}) y_4^{(1)}$

<u>Back propagation: Delta Rule:</u> $\delta_1^{(2)} = \delta_{11}^{(3)} w_{11}^{(3)} + \delta_{12}^{(3)}\delta_{12}^{(3)}$

<u>ANN Training:</u> Initialize $W^{(0)}$ and for each training sample $(x_i, o_i)$, functional values of each neuron is computed in forward pass, weights of each link starting from the output layer is updated using back propagation, and the process continues till convergence.

<u>Improving convergence:</u> To avoid abrupt change in gradient, we use momentum: $\Delta w_i = \alpha \Delta w_i^{(t-1)} - \eta \frac{\partial E^{(t)}}{\partial w_i}$

Adaptive learning rate is used according to change in error.

## Concepts challenging to comprehend:
None

## Interesting and exciting concepts:
I found Kernel Functions and Perceptrons somewhat interesting.

## Concepts not understood:
None

## Any novel idea of yours out of the lessons:
In the kernels, we work in the space of $\varphi(x)$ and the dimensionality d of this new space may be much larger than the dataset size N, and as we use the N×N matrix $XX^T$, it might get computationally expensive. So, to solve this issue, we can use principal component analysis to reduce its linear dimensionality in the $\varphi(x)$ space and compute the projected new k-dimensional values. This will optimize the algorithm and make it efficient.

---------------------------------------------------------------------------------------------------------------------