

BOOSTING

Interval function :-

$$h_{\theta_1, \theta_2, b}(x) = \begin{cases} -b & \text{if } \theta_1 \leq x \leq \theta_2 \\ +b & \text{otherwise} \end{cases}$$



Threshold function :-

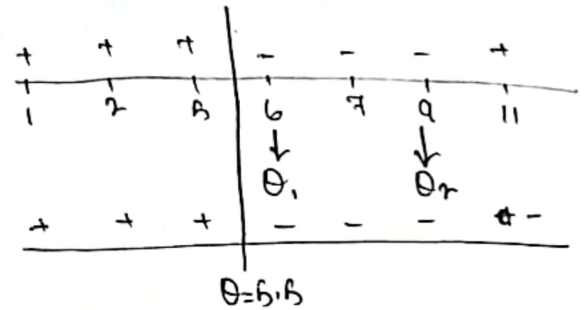
$$h_{\theta, b}(x) = \text{sign}(\theta - x) \cdot b$$

$$\theta = \frac{\theta_i + \theta_{i+1}}{2}$$

$$b \in \{-1, 1\}$$

Example :-

$(1, +), (2, +), (5, +), (6, -), (7, -), (9, -), (11, +)$ → arranged in asc



$$\theta = \frac{5+6}{2} = 5.5, b = 1 \rightarrow \text{for threshold function}$$

$$\text{min Error} \approx 1$$

Weak learners :-

A learning algorithm A is said to be λ -weak learner with $0 \leq \lambda \leq \frac{1}{2}$ for a hypothesis class H if there exists a sample size $m \geq \text{poly}(\frac{1}{\epsilon})$ such that for any $\delta \in (0, 1)$, for any distribution D over X , for any labelling function $f: X \rightarrow \{-1, 1\}$, if the realizability holds, then when the algorithm is run with a sample size of at least m , then the following holds:

$$P_n[R_n(h) \leq \frac{1}{n} - \lambda] \leq 1 - \delta$$

AdaBoost

input:

training set $S = (x_1, y_1), \dots, (x_m, y_m)$

weak learners WL

number of rounds T

$\xrightarrow{\text{Distribution}}$

initialize $D^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m}) \rightarrow$ initially, $D^{(1)}$ is importance same

for $t = 1, \dots, T$:

invoke weak learner $h_t = WL(D^{(t)}, S)$

generalisation error \leftarrow compute $\epsilon_t = \sum_{i=1}^m D_i^{(t)} |y_i \neq h_t(x_i)| \rightarrow$ $y_i \neq h_t(x_i)$ total $\epsilon_t = \sum_{i=1}^m D_i^{(t)} |y_i \neq h_t(x_i)|$ $\epsilon_t \leq \frac{1}{2} - \gamma$

let $w^t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right) \rightarrow w^t$ (as as $[w^t \propto \frac{1}{\epsilon_t}]$)

update $D_i^{(t+1)} = D_i^{(t)} \exp(-w_i y_i h_t(x_i))$ for all $i = 1, \dots, m$
 \downarrow
 D_i update $\sum_{j=1}^m D_j^{(t)} \exp(-w_j y_j h_t(x_i))$

Output:

the hypothesis $h_S(x) = \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right)$

\rightarrow weighted aggregate of the individual predictions of all the returned hypothesis.

* A data gets high probability if it had high probability in the previous round and also if the returned hypothesis makes a mistake on this data instance is identified as a problematic data instance.

* Empirical Risk bound of AdaBoost:-

$$R_S(h_{\text{agg}}) = \frac{1}{m} \sum_{i=1}^m |h_{\text{agg}}(x_i) \neq y_i| \leq e^{-2\gamma^2 T}$$

* Generalisation error:-

$$\sum_{i=1}^m D_i |h(x_i) \neq y_i| \rightarrow R_0(h)$$

$$R_{D, \theta}(h) = \sum_{i: y_i = 1} D_i |h(x_i) \neq 1| + \sum_{i: y_i = -1} D_i |h(x_i) \neq -1| \quad R_{D, \theta+1} = R_{D, \theta} - D_i y_i$$

SUPPORT VECTOR MACHINES

Convex function:-

$$f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v)$$

Condition:-

$$\text{At minimum } u, \nabla f(u) = 0$$

Concave function:-

$$f(\lambda u + (1-\lambda)v) \geq \lambda f(u) + (1-\lambda)f(v)$$

Tangent:-

$$f(x) \geq f(u) + \nabla f(u)^T (x - u)$$

Sub-Gradient:-

$$f(x) \geq f(u) + g^T (x - u)$$

* If u is such that $f(x)$ is differentiable at u , that is tangent & the remaining are sub-gradients.

SVM Objective function

$$\min_{w, b} \frac{1}{2} w^T w + C \sum \max(0, 1 - y_i (w^T x_i + b))$$

Here,

$w^T w$ is regularization term \rightarrow maximise margin

C is hyper-parameter \rightarrow controls trade-off betⁿ a large margin & small hinge loss

$C \sum \max(0, 1 - y_i (w^T x_i + b))$ is hinge/empirical loss \rightarrow penalizes weight vectors that make mistakes

modify,

$$w_0 \leftarrow w; w \leftarrow [w_0, b], x_i \leftarrow [x_i, 1]$$

$$\min_{w, b} \frac{1}{2} w_0^T w_0 + C \sum \max(0, 1 - y_i w^T x_i)$$

Gradient Descent for SVM

(1) Initialize w^0

(2) For $t = 0, 1, 2, \dots$

(i) Compute gradient of $J(w)$ at w^t . (Call $\nabla J(w^t)$)

(ii) Update w^t by taking w^{t+1} by taking a step in the opposite direction of the gradient.

$$w^{t+1} = w^t + \eta \nabla J(w^t)$$

Perceptron v/s SVM

$$L_{\text{perceptron}}(y, x, w) = \max(0, -y w^T x)$$

[No regularization]

Perceptron does not maximise margin width.

$$L_{\text{hinge}}(y, x, w) = \max(0, 1 - y w^T x)$$

[regularization]

SVM optimises hinge loss

Stochastic Gradient Descent for SVM

A training set $S = \{(x_i, y_i)\}$, $x_i \in \mathbb{R}^n$, $y_i \in \{-1, 1\}$ is given.

(1) Initialize $w^0 = 0 \in \mathbb{R}^n$

(2) For Epoch $= 1, \dots, T$

(i) Pick a random example (x_i, y_i) from S .

(ii) Repeat (x_i, y_i) to make a full dataset and take derivative of the SVM objective at w^{t-1} to be $\nabla J^t(w^{t-1})$

$$J^t(w) = \frac{1}{2} w^T w + C \sum \max(0, 1 - y_i w^T x_i)$$

(iii) Update:

$$w^t \leftarrow w^{t-1} - \eta_t \nabla J^t(w^{t-1})$$

(3) Return w

Stochastic Sub-Gradient Descent for SVM

A training set $S = \{(x_i, y_i)\}$, $x_i \in \mathbb{R}^n$, $y_i \in \{-1, 1\}$ is given.

(1) Initialize $w^0 = 0 \in \mathbb{R}^n$

(2) For Epoch $= 1, \dots, T$

(i) For each training example

$(x_i, y_i) \in S$:

if $y_i w^T x_i \leq 1$:

$$w \leftarrow (1 - \eta_t) w + \eta_t C y_i x_i$$

else:

$$w^t \leftarrow w + (1 - \eta_t) w^t$$

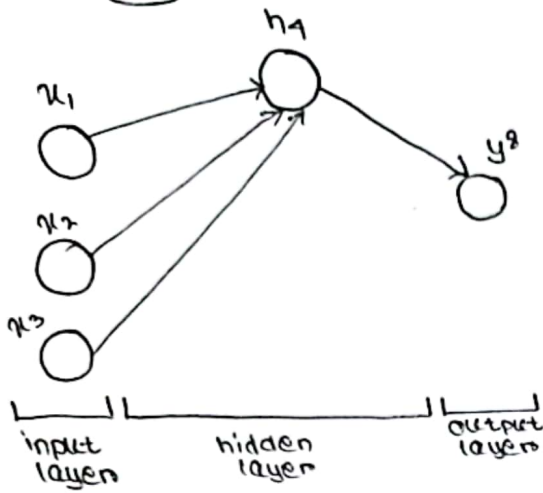
(3) Return w

* Sub-gradient method gives us an update as possible, but it takes time.

$$\eta_t = \frac{1}{1+t}, \quad t \uparrow, \eta_t \downarrow \quad \text{Here, } \eta_t \rightarrow \text{learning rate}$$

$$t \downarrow, \eta_t \uparrow \quad t \rightarrow \text{time}$$

= configuration



$$h_1 = w_{1,1} a_1 + w_{2,1} a_2 + w_{3,1} a_3 + b_1$$

\downarrow
 x_1

\downarrow
 x_2

\downarrow
 x_3

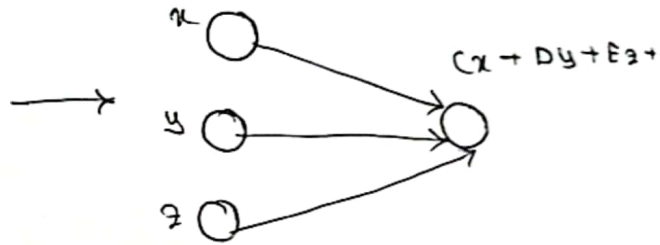
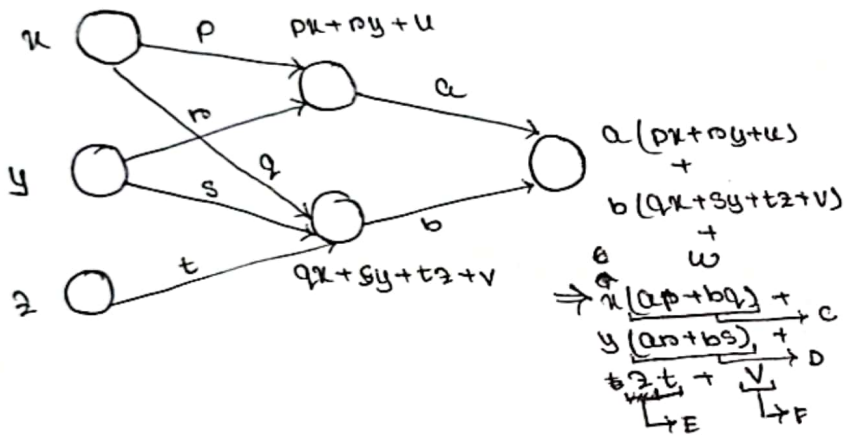
$$Q_1 = \frac{1}{1 + e^{-n_1}}$$

$$y^2 = \omega_{n,2} a_1$$

- identity function:-
 $f(x) = x$

- logistic function:-

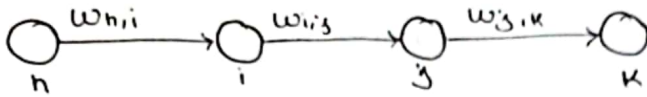
$$f(x) = \frac{1}{1 + e^{-x}}$$



: input layer (2) hidden layers, then hidden layers (2) output layer
layer 2 and 3 use activation function \tanh ,
layer 4 use activation function.

- Input layer has Identity function as the activation function.
- Output layers have logistic function as the activation function.

Back Propagation



$$E = (t_k - a_k)^2$$

$$\frac{\partial E}{\partial w_{j,k}} = \frac{\partial}{\partial w_{j,k}} (t_k - a_k)^2$$

$$= \frac{\partial}{\partial w_{j,k}} (-2(t_k - a_k))$$

$$\frac{\partial E}{\partial w_{j,k}} = a_j g'(in_k) (1 - g(in_k)) (-2(t_k - a_k))$$

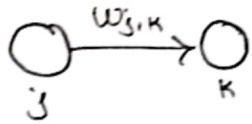
$$w_{j,k} = w_{j,k} - \lambda a_j \Delta_k$$

$$\frac{\partial}{\partial w_{j,k}} = \frac{\partial in_k}{\partial w_{j,k}} * \frac{\partial g(in_k)}{\partial in_k}$$

$$= a_j g(in_k) (1 - g(in_k))$$

$$g'(in_k) = g(in_k) (1 - g(in_k))$$

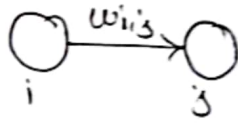
$$\Delta_k = g'(in_k) (-2(t_k - a_k))$$



$$\frac{\partial E}{\partial w_{i,j}} = -2 \sum_k \Delta_k a_i w_{j,k} g'(in_j)$$

$$w_{i,j} = w_{i,j} - \lambda a_i \Delta_j$$

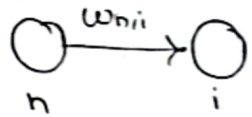
$$\Delta_j = -2 \sum_k \Delta_k w_{j,k} g'(in_j)$$



$$\frac{\partial E}{\partial w_{n,i}} = -2 \sum_j \Delta_j w_{i,j} a_n g'(in_i)$$

$$w_{n,i} = w_{n,i} - \lambda a_n \Delta_i$$

$$\Delta_i = -2 \sum_j \Delta_j w_{i,j} g'(in_i)$$



Adding Momentum to Back-Propagation

$$w_{i,j,n} = w_{i,j,n-1} + \alpha (w_{i,j,n-1} - w_{i,j,n-2}) + \lambda a_{i,n} \Delta_{j,n}$$

$$* a_{i,n} \Delta_{j,n} \neq 0 \rightarrow Gb \text{ doesn't stop}$$

1 K-NEAREST NEIGHBOURS

- * value of K is given
- * distance between testing data and training data is calculated. $d = \sqrt{(x_b - x_t)^2}$
- * Sorted.
- * K -number of shortest distance is checked for Yes/No +ve and No/01-ve.
- * If there are more Yes/No +ve than No/01-ve, then the testing data falls under Yes/No +ve class.
- * If the value of K is odd, then there is no possibility of no. of Yes to be equal to no. of No. However, if there's a tie, we can flip a coin to decide whether to give the testing data a class or not.
- * There's no correct way to determine best value for K .
A small value (such as $K=1$ and $K=2$) can be noisy and subject to outliers.
A large value can smooth over things so can avoid noise.
But, a large value can rule over results from small samples.