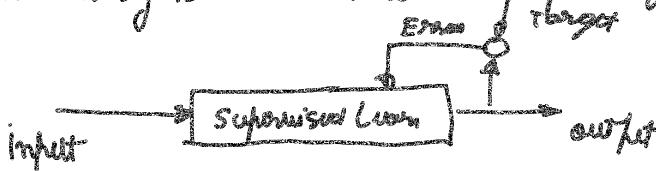


# Math of Intelligence - 1.

1

Machine learning is mathematical optimization.



input  $\rightarrow$  Unsupervised  $\rightarrow$  output.



Easiest way is to learn from labeled data.

Domain  $X$ ; Every point of  $X$  has features

label set  $Y$

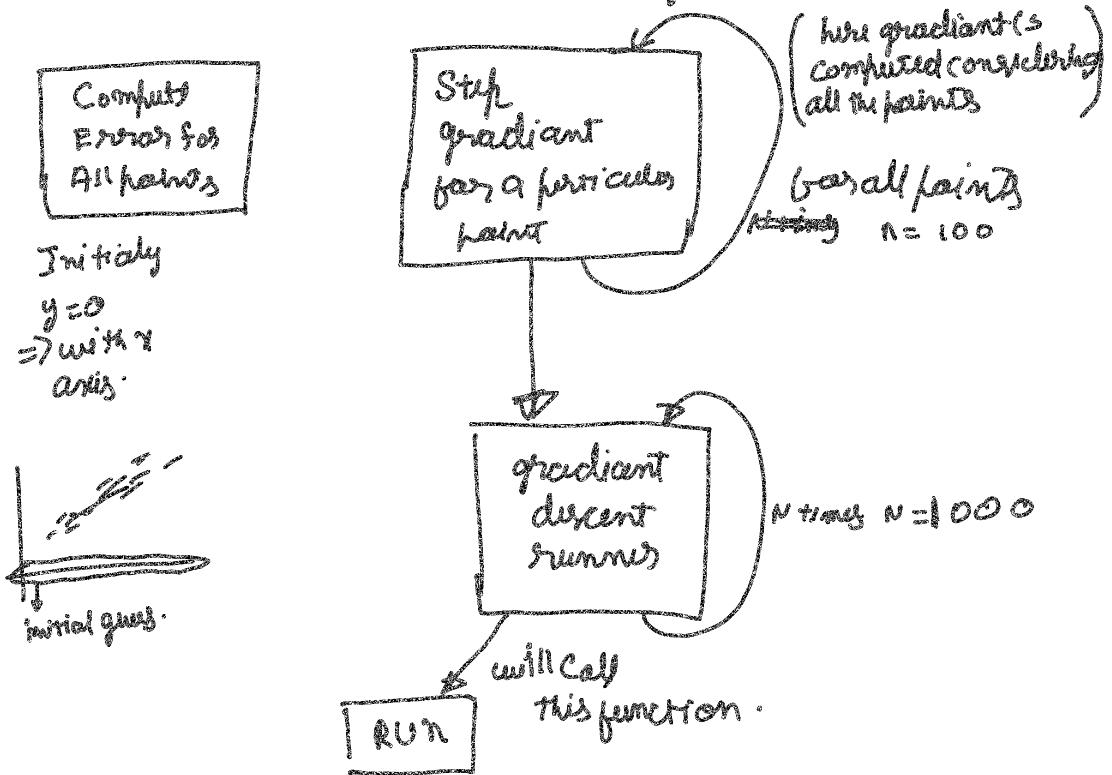
Data, a set of labeled examples  $S = \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq X \times Y$

Output  $\rightarrow$  a prediction rule / hypothesis.

Performance metric  $\rightarrow$  Error<sup>2</sup>

We keep minimizing Error.

## Architecture of code.



## Support Vector Machine

used mainly for classification

we will create a hyperplane to separate a data  
The dimension of a hyperplane is  $n-1$

input space  $\rightarrow$  Feature Space (kernel trick)

Loss function  $\rightarrow$  Hinge Loss function.

mainly comes to maximum-margin Classification.

$$c(x, y, f(x)) = (1 - y \cdot f(x))^+$$

$$c(x, y, f(x)) = \begin{cases} 0 & \text{if } y \cdot f(x) > 1 \\ 1 - y \cdot f(x) & \text{else.} \end{cases}$$

Objective function

$$\min_{\|w\|^2} \|w\|^2 + \sum_{i=1}^n (1 - y_i \cdot (x_i \cdot w))^+$$

~~regulariser term~~  $\rightarrow$  it's a tuning knob

If this is too large it might result in overfit if not underfit.

$$\frac{\delta}{\delta w_k} \lambda \|w\|^2 = 2 \lambda w_k$$

$$\frac{\delta}{\delta w_k} (1 - y_i (x_i \cdot w))^+ = \begin{cases} 0 & \text{if } y_i (x_i \cdot w) > 1 \\ -y_i x_i & \text{else.} \end{cases}$$

Updating weight  $w = w + \eta (-2 \lambda w)$

gradient is a vector which consists of change if contains direction of movement.

Gradient is a vector whose components contains  
partial derivatives all they represent the changes

## Second Order Optimization

function  $\rightarrow$  represents mapping

Machine learning  $\rightarrow$  process of optimization

Gradient Descent

Step 1  $\rightarrow$  Define a model

Step 2  $\rightarrow$  Define a error  $f^n = \frac{1}{N} \sum (y - \hat{y})^2$

Step 3  $\rightarrow$  compute the gradient

use the gradient to update the weights.

all these are first order optimization

next we have

Second order optimization.

\* They do not ignore Curvature of the surface.

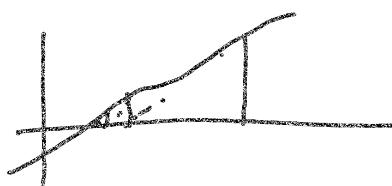
\* Better performance in a single optimization step.

## Newton's Method

\* Finding the roots

\* Optimization.

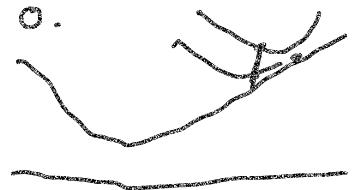
First lets code root finding.



soff iteratively find where  
 $f'' = 0$

In optimization method we will go to find where  
 the derivative of the function is 0.

$$x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$$



The minimizer of quadratic  $f''$  is used.

here we will replace first derivative with

\* gradient

\* Second derivative with Hessian. ( Matrix of second order derivatives. Describes a local curvature of the function ).

$\rightarrow$  Derivative  $\rightarrow \partial f / \partial x$

$\rightarrow$  Gradient  $\rightarrow [\partial f / \partial x_1, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3}, \dots]^T$

$\rightarrow$  Hessian  $\rightarrow \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \rightarrow$  Jacobia  $\rightarrow \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & \frac{\partial f}{\partial x_n} \end{bmatrix}$

first order  $\rightarrow$  uses Jacobia

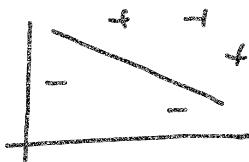
second order  $\rightarrow$  uses Hessian

Newton's method. is for second order method.

## Architecture of SVM Class

$x$ -val  
 $y$ -val  
bias

associated  
label  
 $y$



Loss function:

$$C(x, y, f(x)) = (1 - y \cdot f(x))^+$$

$C \rightarrow \text{Loss function}$      $x \rightarrow \text{sample}$      $y \rightarrow \text{true label}$   
 $f(x) \rightarrow \text{predicted label}$

$$C(x, y, f(x)) = \begin{cases} 0, & \text{if } y \cdot f(x) \geq 1 \\ 1 - y \cdot f(x) & \text{else} \end{cases}$$

Objective function

$$\min_w \underbrace{\lambda \|w\|^2}_{\text{regularization}} + \underbrace{\sum_{i=1}^n (1 - y_i \cdot \langle x_i; w \rangle)^+}_{\text{Loss}}$$

Mimimise the objective by getting its gradients  
for gradient descent.

$$\frac{\partial}{\partial w_k} \lambda \|w\|^2 = 2 \lambda w_k$$

$$\frac{\partial}{\partial w_k} (1 - y_i \cdot \langle x_i; w \rangle)^+ = \begin{cases} 0 & \text{if } y_i \cdot \langle x_i; w \rangle \geq 1 \\ -y_i x_i & \text{else} \end{cases}$$

- If we misclassified the sample
- Then we update weight vector  $w$  using gradients of both terms
- Else if classified correctly  
we will just update  $w$  by the gradient of the regularizer.

Miss Classification Condition

$$y_i < x_i, w > < 1$$

- update rule for our weights (miss classified)

$$w = w + \gamma (y_i x_i - 2 \lambda w)$$

$\gamma \rightarrow$  learning rate (This is the length of steps algorithm makes down the gradient on the error curve)  
 $\lambda \rightarrow$  regularizer.

$\gamma \rightarrow$  high  $\rightarrow$  algorithm might overshoot

$\gamma \rightarrow$  low  $\rightarrow$  algorithm takes long time to converge

---

Regularizer will control trade off between achieving low training error and low testing Error  
(ie for giving ability to generalise the classifier to unseen data) as regularising parameter  $w$  will choose 1/epoch such that this number will decrease as no. of epoch increases.

update rule for weights when classified correctly.

$$\omega = \omega + \eta (-2 \gamma \omega)$$