

For solving NLP there is no unique method to get the solution — Lagrangian multiplier approach is one approach.
 → This can be solved using Lagrange multiplier approach.

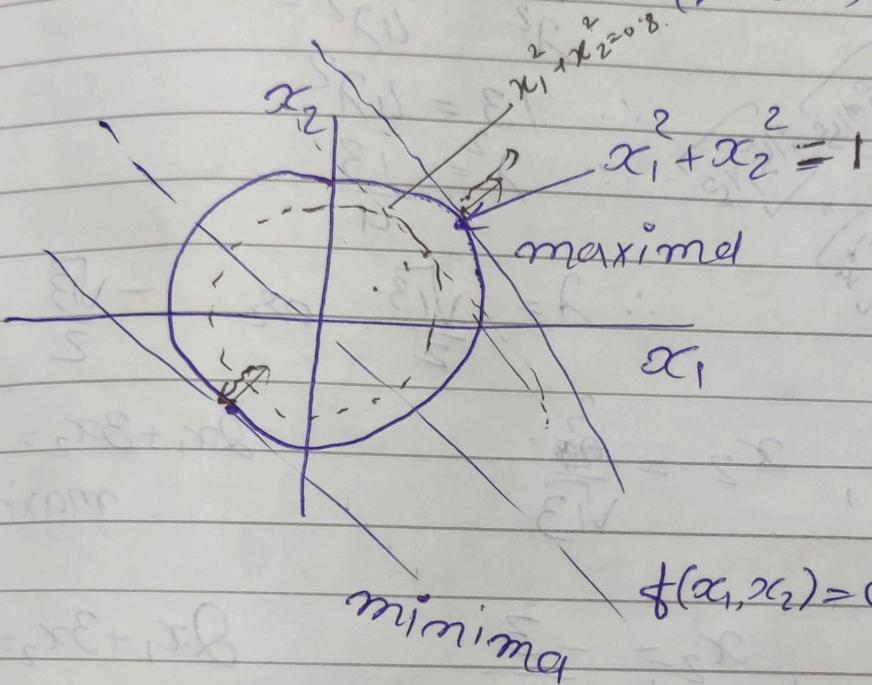
Ex. Optimize $\underbrace{2x_1 + 3x_2}$ such that $x_1^2 + x_2^2 = 1$

$f(x_1, x_2)$, objective

function (linear)

constraint
(non-linear)

$$g(x_1, x_2) = x_1^2 + x_2^2 - 1$$



$$f(x_1, x_2) = 0$$

minima

maxima

x_1

x_2

$$\begin{aligned} x_1 &= \sqrt{1-x_2^2} \\ x_2 &= \sqrt{1-x_1^2} \\ x_1 &= \sqrt{1-x_2^2} \\ x_2 &= \sqrt{1-x_1^2} \\ x_1 &= \sqrt{1-x_2^2} \\ x_2 &= \sqrt{1-x_1^2} \end{aligned}$$

→ Lagrangian $L(x_1, x_2, \lambda) = f(x_1, x_2) + \lambda g(x_1, x_2)$

$$\text{Lagrange} = (2x_1 + 3x_2) + \lambda(x_1^2 + x_2^2 - 1)$$

multiplier

→ Differentiate w.r.t. x_1, x_2, λ and equate to zero and solve.

$$\frac{\partial L(x_1, x_2, \lambda)}{\partial x_1} = 2 + 0 + 2x_1\lambda = 0$$

where λ is called Lagrange multiplier

$$x_1 = -\frac{1}{\lambda}$$

* Gradient of f , $\nabla f(x_1, x_2) = \begin{bmatrix} \frac{3}{\sqrt{13}} \\ \frac{-3}{\sqrt{13}} \end{bmatrix}$ & at $\lambda = \begin{bmatrix} -\frac{3}{\sqrt{13}} \\ \frac{-4}{\sqrt{13}} \end{bmatrix} = \begin{bmatrix} -\frac{3}{\sqrt{13}} \\ -\frac{6}{\sqrt{13}} \end{bmatrix}$ the same direction
for min gradients are in opposite direction

$$-\frac{\partial L(x_1, x_2, \lambda)}{\partial x_2} = 0 + 3 + 2x_2 \lambda = 0 \Rightarrow x_2 = -\frac{3}{2\lambda}$$

$$-\frac{\partial L(x_1, x_2, \lambda)}{\partial \lambda} = 0 + 0 + x_1^2 + x_2^2 - 1 = 0 \Rightarrow x_1^2 + x_2^2 = 1$$

$$\frac{1}{\lambda^2} + \frac{9}{4\lambda^2} = 1$$

Gradient vector $\nabla f = \begin{bmatrix} 3 \\ -3 \end{bmatrix}$
Directional derivative $\nabla f \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ -3 \end{bmatrix} \cdot \begin{bmatrix} \frac{x_1}{\sqrt{13}} \\ \frac{x_2}{\sqrt{13}} \end{bmatrix} = \frac{3}{\sqrt{13}}$

$$\therefore 13 = 4\lambda^2 \Rightarrow \lambda^2 = \frac{13}{4}$$

$$\therefore \lambda = \frac{\sqrt{13}}{2} \text{ or } -\frac{\sqrt{13}}{2}$$

$$\lambda = \frac{-\sqrt{13}}{2} \therefore x_1 = \frac{2}{\sqrt{13}}, x_2 = \frac{3}{\sqrt{13}}$$

$$2x_1 + 3x_2 = 3.6 \text{ maxima}$$

$$\lambda = +\frac{\sqrt{13}}{2} \quad x_1 = -\frac{2}{\sqrt{13}}, \quad x_2 = -\frac{3}{\sqrt{13}}$$

$$2x_1 + 3x_2 = -3.6 \text{ minima}$$

$\lambda > 0 \Rightarrow$ maxima and vice versa

* Consider the following problem which we call as primal optimization problem.

$$\min_{\theta} f(\theta)$$

$$\text{s.t. } g_i(\theta) \leq 0 \quad i=1 \text{ to } k$$

$$h_i(\theta) = 0 \quad i=1 \text{ to } l$$

θ are called primal variables



→ To solve this problem, define generalized lagrangian.

$$L(\underline{\theta}, \underline{\alpha}, \underline{\beta}) = f(\underline{\theta}) + \sum_{i=1}^K \alpha_i g_i(\underline{\theta}) + \sum_{i=1}^l \beta_i h_i(\underline{\theta})$$

- α_i, β_i are called lagrange multipliers.

→ Consider the quantity

$$\Theta_p(\underline{\theta}) = \max_{\underline{\alpha}, \underline{\beta}, \alpha_i \geq 0} L(\underline{\theta}, \underline{\alpha}, \underline{\beta})$$

Primal

→ let some $\underline{\theta}$ be given if $\underline{\theta}$ violates any of the primal constraints i.e. $g_i(\underline{\theta}) > 0$, $h_i(\underline{\theta}) \neq 0$
then $\Theta_p(\underline{\theta}) = \infty$

- If primal constraints are satisfied

$$\Theta_p(\underline{\theta}) = f(\underline{\theta}) \rightarrow \text{section will be that } f(\underline{\theta}) \text{ which satisfies constraint. [See Example } \frac{x_1^2}{2} + \frac{x_2^2}{2} \leq 1]$$

→ So we can write original primal problem as

$$\min_{\underline{\theta}, \underline{\alpha}, \underline{\beta}} \Theta_p(\underline{\theta})$$

maximize such that
the conditions are
satisfied.

Find the max g by choosing $\underline{\alpha}, \underline{\beta}$ with the constraints satisfied.

$$\therefore \min_{\underline{\theta}} \left(\max_{\underline{\alpha}, \underline{\beta}, \alpha_i \geq 0} L(\underline{\theta}, \underline{\alpha}, \underline{\beta}) \right)$$

↑ choose $\underline{\theta}$
max over α, β — (i)
then $f(\underline{\theta})$ comes to
do this for max $\underline{\theta}$

This will give optimum values of parameters $\underline{\theta}^*, \underline{\alpha}^*, \underline{\beta}^*$

→ The dual problem can be written as

$$\max_{\underline{\alpha}, \underline{\beta}, \alpha_i \geq 0} \left[\min_{\underline{\theta}} L(\underline{\theta}, \underline{\alpha}, \underline{\beta}) \right] \quad -\text{(ii)}$$

dual variables

→ let at $\underline{\theta}^*, \underline{\alpha}^*, \underline{\beta}^*$ value of eq. (i) be p^*
and value of eq. (ii) be d^* .

$\begin{array}{ll} \text{for some } \rightarrow & \underline{\theta}_1 \\ & \underline{\theta}_2 \\ & \underline{\theta}_3 \\ & \underline{\theta}_4 \end{array}$	$\begin{array}{l} \max = 25 \\ \max = 26 \\ \max = 27 \\ \max = 88 \end{array}$	{	for some	$\begin{array}{ll} \underline{\alpha}_1, \underline{\beta}_1 \\ \underline{\alpha}_2, \underline{\beta}_2 \\ \underline{\alpha}_3, \underline{\beta}_3 \\ \underline{\alpha}_4, \underline{\beta}_4 \end{array}$	min = 2 min = 3 min = 4 → min = 1 <small>max of (2, 3, 4)</small>

$$P^* = 25 \quad d^* = 4$$

for minimum of maximum (25, 26, 27, 88)

min (25, 26, 27, 88) maximum of minimum

$$d^* \leq P^*$$

→ In certain conditions $d^* = P^*$.
 What are the conditions under which $d^* = P^*$? So that one can get $\underline{\theta}$ either by using primal or dual formulation of the problem. These conditions are called as Kausch - Kuhn - Tucker (KKT) conditions listed below.

$$(i) \quad \textcircled{1} \quad \frac{\partial}{\partial \theta_i} [L(\underline{\theta}, \underline{\alpha}, \underline{\beta})] \Big|_{\substack{\downarrow \\ \text{all } \theta's}} \Bigg|_{\substack{\theta = \theta^*, \alpha = \alpha^*, \beta = \beta^*}} = 0 \quad i = 1 \dots n$$

Note that I am using $\underline{\theta}$ & $\underline{\beta}$ both for vector notation PAGE NO. _____ DATE _____
 ↗ Abuse of notation ↘ Parameters

$$\textcircled{2} \quad \frac{\partial}{\partial \beta_i} [L(\underline{\theta}, \underline{\alpha}, \underline{\beta})] \Big|_{\substack{\underline{\theta}=\underline{\theta}^*, \underline{\beta}=\underline{\beta}^*, \underline{\alpha}=\underline{\alpha}^*}} = 0 \quad i=1 \text{ to } m$$

$$\begin{array}{ll} \textcircled{3} & \alpha_i^* g_i(\underline{\theta})^* = 0 & i=1 \text{ to } k \\ \textcircled{4} & g_i(\underline{\theta})^* \leq 0 & > \text{ for } i=1 \text{ to } k \\ \textcircled{5} & \alpha_i^* \geq 0 & \text{maxima } i=1 \text{ to } k \end{array}$$

→ These conditions will be satisfied if ~~f, g, h~~ are convex and all ~~f, g, h~~ are affine.
 (all linear (affine) functions are convex)

* Convex :- $f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$

f is a real function defined on interval i $\alpha \in [0, 1]$

~~affine~~ :- $h(x) = ax + b$ is both convex and concave.

Let us come to SVM now

For SVM, we can write Lagrangian as

$$L(\underline{\beta}\underline{\theta}, \theta_0, \underline{\alpha}) = \frac{1}{2} \|\underline{\beta}\underline{\theta}\|^2 + \sum_{i=1}^m \alpha_i [1 - y^{(i)} (\underline{\theta}^\top \underline{x}^{(i)} + \theta_0)].$$

$\underbrace{\dots}_{m \text{ terms}} \quad \underbrace{\dots}_{f(\underline{\theta}, \theta_0)} \quad \underbrace{\dots}_{g_i(\underline{\theta}, \theta_0)}$

→ Using the dual optimization approach

$$\max_{\underline{\alpha}} \left[\min_{\underline{\theta}, \theta_0} L(\underline{\theta}, \theta_0, \underline{\alpha}) \right], \text{ there are } m \text{ number of } \alpha's$$

$$\alpha_i, i=1 \text{ to } m.$$

Our problem is $\min_{\underline{\theta}, \theta_0} \frac{1}{2} \|\underline{\theta}\|^2$ s.t. $\{1 - y^{(i)} [\underline{\theta}^\top \underline{x}^{(i)} + \theta_0]\} \leq 0$

for $\min_{\underline{\theta}, \theta_0} f(\underline{\theta}, \theta_0)$ s.t. $g_i(\underline{\theta}, \theta_0) \leq 0, i=1 \text{ to } m$ examples
 NO $h(\underline{\theta})$ here

→ Considered minimization
 - differentiate $L(\theta, \theta_0, \underline{\alpha})$ w.r.t. each of $\theta_i, i=0 \text{ to } n$ and equate it to 0.

$$\frac{\partial}{\partial \theta} [L(\theta, \theta_0, \underline{\alpha})] = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} \nabla_{\theta} (L(\theta, \theta_0, \underline{\alpha})) = 0$$

$$\frac{\partial}{\partial \theta_n} [L(\theta, \theta_0, \underline{\alpha})] = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\frac{\partial}{\partial \theta_0} [L(\theta, \theta_0, \underline{\alpha})] = 0$$

→ Hence we will get

$$\theta = - \sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)} = \underline{\theta} \quad \begin{array}{l} \text{vector} \\ \text{this is a vector} \end{array}$$

$$\theta_0 - \sum_{i=1}^m \alpha_i y^{(i)} x_0^{(i)} = 0 \quad \begin{array}{l} \text{scalar} \end{array}$$

$$\theta_m - \sum_{i=1}^m \alpha_i y^{(i)} x_n^{(i)} = 0$$

$$\underline{\theta}^* = \sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)}$$

where α_i 's are unknown

before we go for getting α_i^* 's, $\alpha_i \geq 0$

→ differentiate w.r.t. $\underline{\theta}_0$

$$\therefore \underline{\theta} = \sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)}$$

→ By substituting $\underline{\theta}$, into $L(\underline{\theta}, \underline{\theta}_0, \underline{\alpha})$

$$L(\underline{\theta}, \underline{\theta}_0, \underline{\alpha}) = \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)} \right]^T \left[\sum_{j=1}^m \alpha_j y^{(j)} \underline{x}^{(j)} \right]$$

$$+ \sum_{i=1}^m \alpha_i \left[1 - y^{(i)} \left[\left(\sum_{j=1}^m \alpha_j y^{(j)} \underline{x}^{(j)} \right)^T \underline{x}^{(i)} + \underline{\theta}_0 \right] \right]$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} \underline{x}^{(i)}^T \underline{x}^{(j)} + \underline{\theta}_0$$

because $\sum \alpha_i = 1$

→ This has to be maximized with constraints
as $\alpha_i \geq 0$ for $i=1$ to m , $\sum_{i=1}^m \alpha_i y^{(i)} = 0$

→ Form lagrangian $L(\underline{\alpha}, \underline{\mu}, \underline{\lambda})$ objective function
and can be solved. HW

→ Program is available for Quadratic optimization with linear constraints.

→ Once we get $\underline{\alpha}$

$$\underline{\theta} = \sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)}$$

* looking at KKT conditions

$$-\quad \alpha_i^* g_i(\underline{\theta}^*, \underline{\theta}_0^*) = 0 \quad i = 1 \dots m$$

$$-\quad g_i(\underline{\theta}^*, \underline{\theta}_0^*) = 1 - y^{(i)} [\underline{\theta}^{*T} \underline{x}^{(i)} + \underline{\theta}_0^*] \leq 0$$

$\alpha_i^* \geq 0$

→ If $\alpha_i > 0$

nonzero

Lagrange multiplier
are nonzero
only for
support
vec., data

$$\therefore y^{(i)} [\underline{\theta}^{*T} \underline{x}^{(i)} + \underline{\theta}_0^*] = 1$$

support vectors

→ If $\alpha_i = 0$ then $g_i(\underline{\theta}^*, \underline{\theta}_0^*) < 0$, non support vectors.

$$(i.e. \underline{\theta}^{*T} \underline{x}^{(i)} + \underline{\theta}_0^* > 1)$$

This implies
and

In getting contributing. ~~contributing~~ only support vectors are

$$\underline{\theta}^* = [\alpha_1, \alpha_2, \dots, \alpha_m]$$

which means if we know in advance which data points constitute

Support vectors

~~the opposite~~
 ~~$\underline{\theta}^*$ (nonzero) can be~~
~~computed only for~~
~~those data points~~

It is not possible to know
know which data points
constitute support vectors

* For positive support vectors:-

$$\underline{\theta}^{*T} \underline{x}^{(i)} + \underline{\theta}_0^* = 1$$

$$\therefore \underline{\theta}_0^* = 1 - \underline{\theta}^{*T} \underline{x}^{(i)}$$

- For negative support vectors:-

$$\underline{\theta}_0^* = -1 - \underline{\theta}^{*T} \underline{x}^{(i)}$$

- In general -

$$\boxed{\underline{\theta}_0^* = y^{(i)} - \underline{\theta}^{*T} \underline{x}^{(i)}}$$

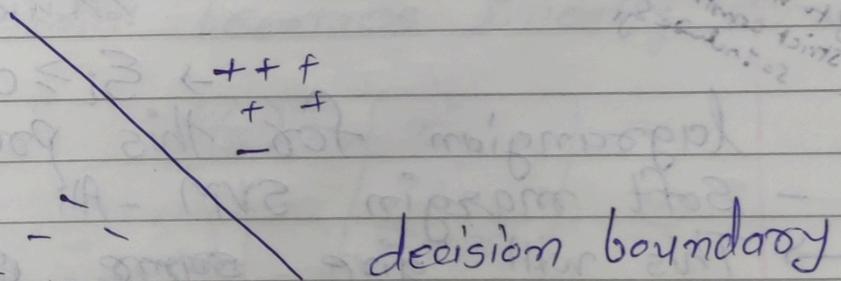
→ This is a convex quadratic optimization problem so unique solution for $\underline{\theta}, \underline{\theta}_0$
(An Advantage with SVM) → unlike neural nets.

- * Few points about Lagrangian that we get after substituting $\underline{\theta}^{(i)}$.
- Always $\underline{x}^{(i)}$ and $\underline{x}^{(i)T}$ appear as dot product.
- Unknowns are in which is irrespective of size of \underline{x} . (data vector) \rightarrow this

What we did:

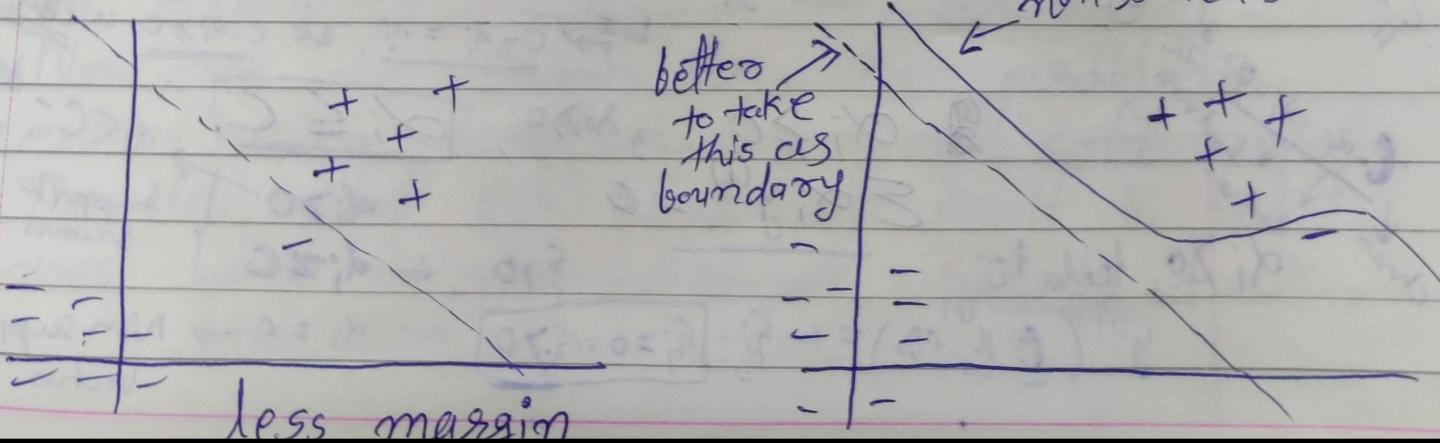
$$\min \frac{1}{2} \|\underline{\theta}\|^2 \quad \text{s.t. } y^{(i)} (\underline{\theta}^T \underline{x}^{(i)} + \theta_0) \geq 1 \quad i=1 \text{ to } m$$

- This is called as hard margin classifier.
- If the data is not separable by hyperplane then this method has no solution.



→ To solve such a problem we use soft margin classifier.

* Soft margin SVM (SVM with regularization):



* Problem formulation :-

- We still get hyperplane taking care of outliers & non-linear data.

* Objective function

$$\min_{\underline{\theta}, \xi_i} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^m \xi_i$$

data that is not separable by linear boundary (hyperplane is linear boundary)

If we don't include this we have no constraint since any point can violate it

Any point which violates the margin it adds an penalty in the objective fn & we try to minimize this objective fn.

such that $\rightarrow y^{(i)} (\underline{\theta}^T \underline{x} + \underline{\theta}_0) \geq 1 - \xi_i$

slack variable
not held tightly in position

slack parameter

$\rightarrow \xi_i \geq 0$

H.W.

lagrangian for this problem, dual form

- Soft margin SVM - Ali Ghodsi

\rightarrow This will give some

$$\underline{\theta} = \sum_{i=1}^m \alpha_i y^{(i)} \underline{x}^{(i)}$$

In soft margin
 $0 \leq \alpha_i \leq C$

Convex
Hence global minimum

$$\max_{\underline{\alpha}} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} \underline{x}^{(i)} \underline{x}^{(j)}$$

but $\rightarrow C - \alpha_i = \underline{\lambda}_i$, i.e. $C - \alpha_i \geq 0 \Leftrightarrow \underline{\lambda}_i \leq C$

$$\alpha_i \geq 0, \underline{\lambda}_i \geq 0, \underline{\lambda}_i \leq C$$

$\alpha_i > 0 \rightarrow$ support vector
 $\alpha_i = 0 \rightarrow$ non-support vector

$\alpha_i > 0$, leads to

$$y^{(i)} (\underline{\theta}^T \underline{x}^{(i)} + \underline{\theta}_0) = 1 - \xi_i, \quad \xi_i \geq 0 \text{ or } \xi_i < 0$$

$\alpha_i = 0 \rightarrow$ non-support vector

→ data is highly non-linear, then SVM with Kernel function. ~~for~~, for non-linearly separable data points Kernel trick works better.

* Kernel Trick:-

- Hard margin classifier = linearly separable
- Soft margin classifier = slack variable ξ_i are used
- takes care of non-linearly separable but still not the best classifier
- Use of Kernel function - Best SVM classifier with high accuracy
- When the data is not linearly separable, project it onto higher dimensional space where the data becomes linearly separable.

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\phi(\underline{x}) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{bmatrix}$$

\downarrow lower dimensional data \downarrow higher dimensional data

\rightarrow we saw that using polynomial regression we work on features gives us results

i.e. lower to higher dimensional mapping

m examples $\xrightarrow{\phi}$ m examples

- Now, we can formulate the problem either using hard margin or soft margin
- The data used $\phi(\underline{x})$. But ϕ itself is not used ~~due to the use of~~ what we use is KERNEL function.