

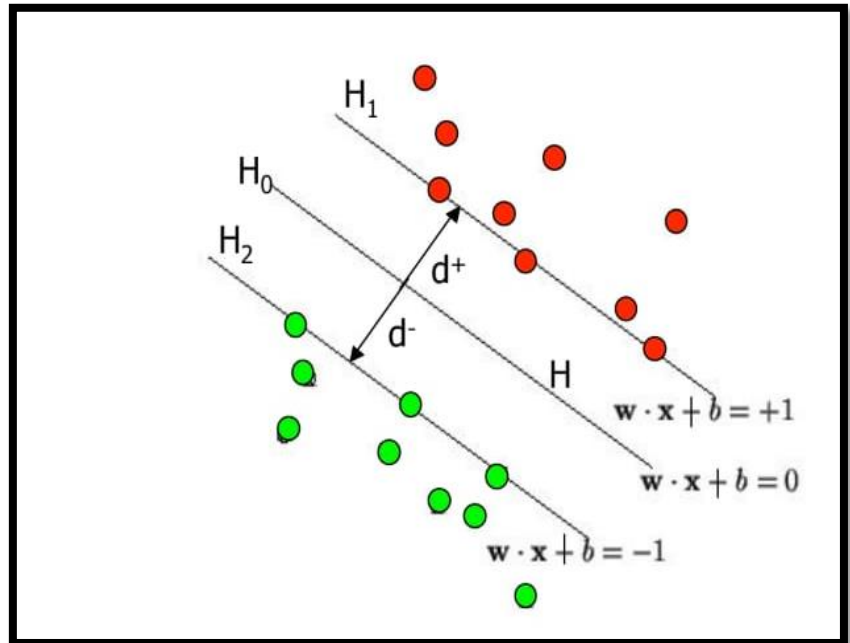
Decoding SVM

- Now we want to explore the maths and logic running behind the SVM which appears like a magic. We will decode it end to end. Let's get started:

First, we have 'm' number of feature vector in 'n' dimensional plane. So, we plot the m vectors in the n dimensional plane and we would be having m number of vector points.

- Now our objective is to find such a hyperplane that would:

- 1) Correctly classifies the classes of data points.
- 2) Ensures the maximum distance between the closest data point of each class and the boundary which is also known as 'margin'.



So, we want the equation of the hyperplane ensuring above two conditions.

As, equation of line in 2D is (Slope-Intercept form):

$$y = mx + c$$

Hyperplane is generalization is the analogy of line in 2D. So, its equation will be:

$$w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + \dots + b = 0$$
$$\Rightarrow wx + b = 0$$

Where,

$w = [w_1, w_2, w_3 \dots, w_n]$ (It determines the orientation of the hyperplane)
 $b =$ intercept which determines the position of hyperplane

The w is always perpendicular to the hyperplane. To prove it:

- We consider two feature points of the hyperplane which can't be traced in the figure (Just Imagine it). So:

$$w \cdot x_1 + b = 0$$

and

$$w \cdot x_2 + b = 0$$

$$\Rightarrow w \cdot (x_2 - x_1) = 0$$

Since the vector $(x_2 - x_1)$ will lie in the hyperplane and its dot product with the w is zero. This proves that they are perpendicular.

Here we are considering that we have Hard-Margin means we strictly considering that there is no point existing between the wall and the boundary and in the wrong side.

Now we assign labels to each class i.e.:

For red= $y_r = +1$

For green= $y_g = -1$

- So, as per the hard margin consideration the red points will have label= +1 and obey $w \cdot x + b \geq 1$ and for green points will have label= -1 and obey $w \cdot x + b \leq -1$.

Therefore, we can write the combined form of the conditions in a single inequality: $y_i * (w \cdot x_i + b) \geq 1$.

Now our interest is to find the hyperplane more specifically equation of the hyperplane achieving my objective.

Here the boundary we have considered are $w \cdot x + b = 1$ and $w \cdot x + b = -1$.

Let the margin to be 'd'.

Now the distance of a point (x_1, y_1) from a line $ax + by + c = 0$ is:

$$\text{Distance} = \frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

Similarly, the distance between the planes H_1 and H_0 is:

$$d = \frac{w \cdot x_1 + b}{\|w\|} = \frac{1}{\|w\|} \text{ as the point lies on } H_1 \text{ and satisfies } w \cdot x_1 + b = 1.$$

So, the total distance between the H_0 and H_1 is:

$$2d = \frac{2}{\|w\|} = \text{distance between the two hyperplanes.}$$

Now we need to calculate w and b such that the $2d$ distance is maximum as well as it is in between the hyperplanes to correctly classify the classes.

So, we have to maximize $\frac{2}{\|w\|}$.

\Rightarrow We need to minimize $\|w\|$ and for mathematical convenience we would be minimizing the $f(w) = \frac{1}{2} * \|w\|^2$ (as its derivative will be just w !)

But the w is subjected to the following constraint:

$$y_i * (w \cdot x_i + b) \geq 1$$

So, we can't consider any value of the w . It must be such a value which will minimize the $f(w)$ as well as satisfy the constraint.

It's a messy problem to solve with inequality. To remove inequality, we go with Lagrange Multiplier Method.

$$L = \frac{1}{2} w^2 - \sum_1^m a_i * [y_i * (w \cdot x_i + b) - 1]$$

Now we need to minimize the L wrt w and b and maximize wrt ai. As:

- 1) We want L to be minimum and the ai are positive so the L must be maximized wrt ai.
- 2) The L must be minimized wrt w and b as we want the maximum value of the margin.

So, implementing the minimization of the L wrt w and b:

$$\begin{aligned} \frac{\partial L}{\partial \vec{w}} &= \frac{1}{2} * \frac{\partial (\vec{w} \cdot \vec{w})}{\partial \vec{w}} - \sum_1^m a_i * y_i * \left[\frac{\partial (\vec{w} \cdot \vec{x}_i)}{\partial \vec{w}} \right] \\ &= \vec{w} - \sum_1^m a_i * y_i * \vec{x}_i \end{aligned}$$

Now computing the $\frac{\partial L}{\partial \vec{w}} = 0$, we get:

$$\vec{w} = \sum_1^m a_i * y_i * \vec{x}_i$$

As, for above step we can extend justification as:

$$\vec{w} = [w_1, w_2, \dots, w_n]$$

$$\vec{w} \cdot \vec{w} = [w_1^2, w_2^2, w_3^2, \dots, w_n^2]$$

$$\begin{aligned} \text{So, } \frac{\partial L}{\partial \vec{w}} &= \left[\frac{\partial (w_1)^2}{\partial w_1}, \frac{\partial (w_2)^2}{\partial w_2}, \frac{\partial (w_3)^2}{\partial w_3}, \dots, \frac{\partial (w_n)^2}{\partial w_n} \right] \\ &= \frac{1}{2} * [w_1, w_2, \dots, w_n] \\ &= \frac{1}{2} \vec{w} \end{aligned}$$

Similarly, for partial derivative of x_i dot product with w vector wrt w vector will give x_i .

Now, minimizing wrt b , we get:

$$\frac{\partial L}{\partial b} = 0 - \sum_1^m a_i y_i$$

On equating $\frac{\partial L}{\partial b} = 0$, then $\sum_1^m a_i y_i = 0$.

So, we got a constraint here that $\sum_1^m a_i y_i = 0$.

Reconstructing the original L :

$$L = \frac{1}{2} * \left(\sum_1^m a_i * y_i * \vec{x}_i \cdot \sum_1^m a_j * y_j * \vec{x}_j \right) - \sum_1^m a_i * \left[y_i * \left(\sum_1^m (a_j * y_j * \vec{x}_j) \cdot x_i + b \right) - 1 \right]$$
$$\Rightarrow L = \sum_1^m a_i - \sum_1^m a_i * y_i * b - \frac{1}{2} \sum_1^m a_i a_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

$$\Rightarrow L = \sum_1^m a_i - \frac{1}{2} \sum_1^m a_i a_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

Now, we need to maximize it wrt a_i and the equation is subjected to following two constraints:

$$1) \sum_1^m a_i y_i = 0$$

$$2) 0 \leq a_i$$

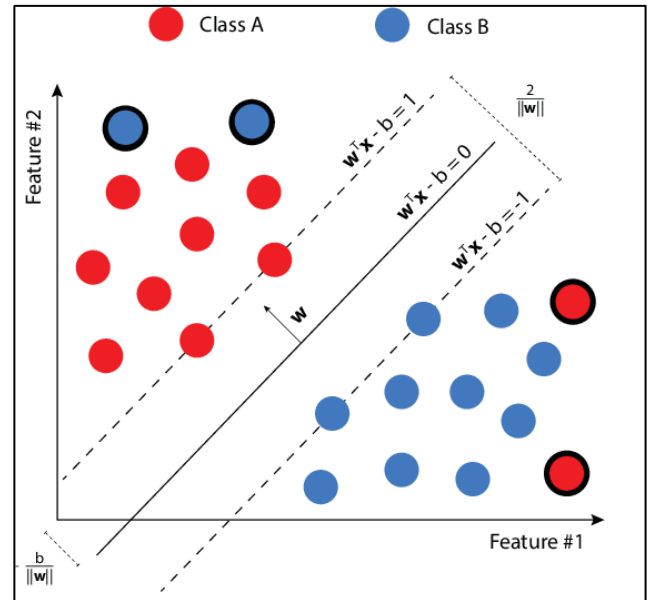
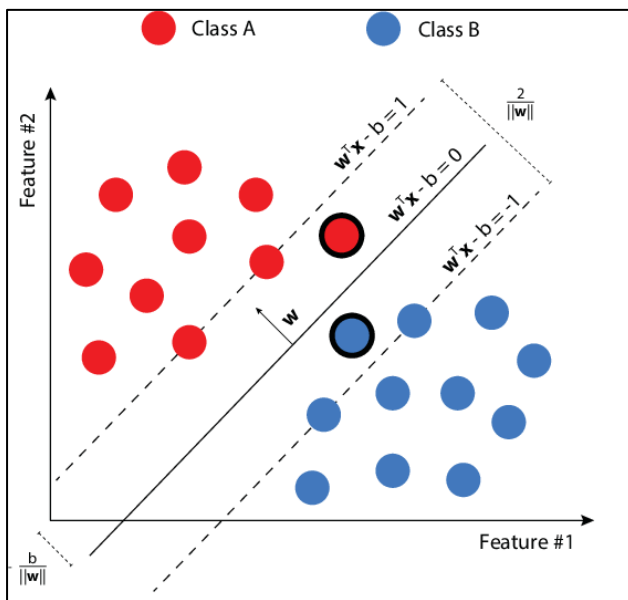
- The above dual optimization problem is solved by the optimizer to obtain the best possible values of a_i and a_j until it finds the best values that maximize the Lagrangian under constraints.

So, once the a_i and a_j are calculated for the dual optimization problem then w can be calculated very much and then the hyperplane equation can be easily obtained which we want very much.

- Now, the question arises is will we get the data points in our real-life dataset possessing a hard margin?

Certainly not as real-world datasets are very messy. Its very rear to have.

In real world we certainly have these situations:



Now, we need to accommodate this situation:

- The equation of the plane will remain the same very much.
- We need to identify all the points even which are within the margin from hyperplane and on the wrong side. For it, we construct the inequality:

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i$$

where, the values of:

$\xi=0$ for points present in the right location.

$0 < \xi < 1$ for points present within the margin and correctly classified.

$\xi > 1$ for the points on the wrong side of the plane.

Now, we want to punish all the points who are not present in the right location. So, our objective function becomes:

$$\text{Objective_Function} = \frac{1}{2}w^2 + C_i * \sum_1^m \xi_i$$

Now, we want two things as our objective:

- 1) To minimize the penalty for the losses.
- 2) To maximize the margin.

But the Objective_Function is subjected to the following constraints:

- 1) $y_i(w \cdot x_i + b) \geq 1 - \xi_i$
- 2) $0 \leq \xi \leq 1$

So, the Lagrangian will be:

$$L = \frac{1}{2}w^2 + C_i \sum_1^m \xi_i - \sum_1^m a_i * (y_i(w \cdot x_i + b) + \xi_i - 1) - \sum_1^m a_j * \xi_i$$

Now, we need to:

- 1) Minimize the Lagrangian wrt w , b , ξ_i .
- 2) Maximize the Lagrangian wrt a_i and a_j .

So,

$$\frac{\partial L}{\partial w} = \bar{w} - \sum_1^m a_i y_i \bar{x}_i$$

$$\Rightarrow \bar{w} = \sum_1^m a_i y_i \bar{x}_i \text{ when } \frac{\partial L}{\partial w} \text{ is equated with } 0$$

Also,

$$\frac{\partial L}{\partial b} = - \sum_1^m a_i y_i$$

$$\Rightarrow \sum_1^m a_i y_i = 0 \text{ when } \frac{\partial L}{\partial b} \text{ is equated with } 0$$

And,

$$\frac{\partial L}{\partial \xi} = \sum_1^m C_i - a_j - a_i$$

$$\Rightarrow \sum_1^m C_i = \sum_1^m a_i + a_j \text{ when } \frac{\partial L}{\partial \xi} \text{ is equated with } 0$$

As, a_i and a_j are greater than zero so we can easily conclude that:

$$0 \leq a_i \leq C_i$$

and

$$0 \leq a_j \leq C_j$$

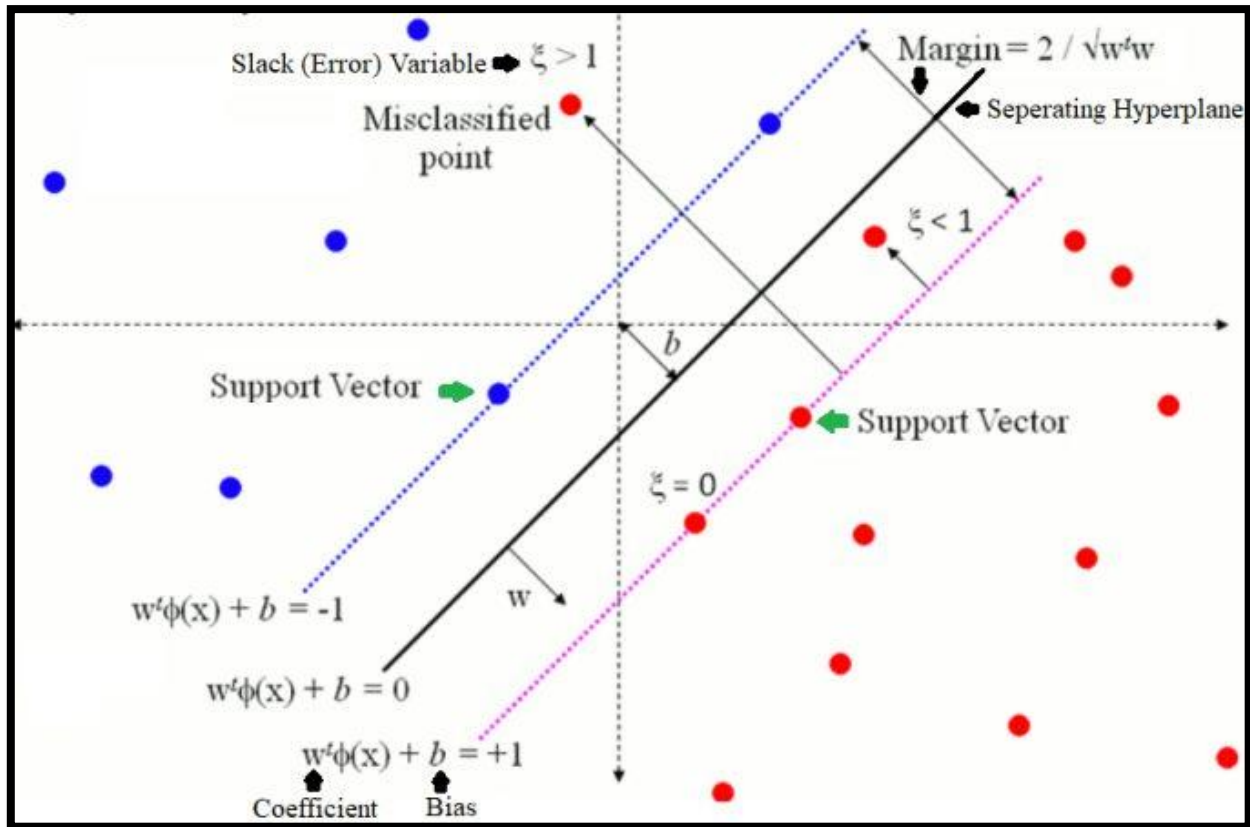
- Now, we are only remained with the multiplier finding which is too computed by the optimizer in the same way as done in Hard-Margin as the dual form is the very much same for them.
- And using support vectors we would be able to compute the value of b as datapoints are known to us so we would be easily able to do it.

So, this is how the SVM actually works.

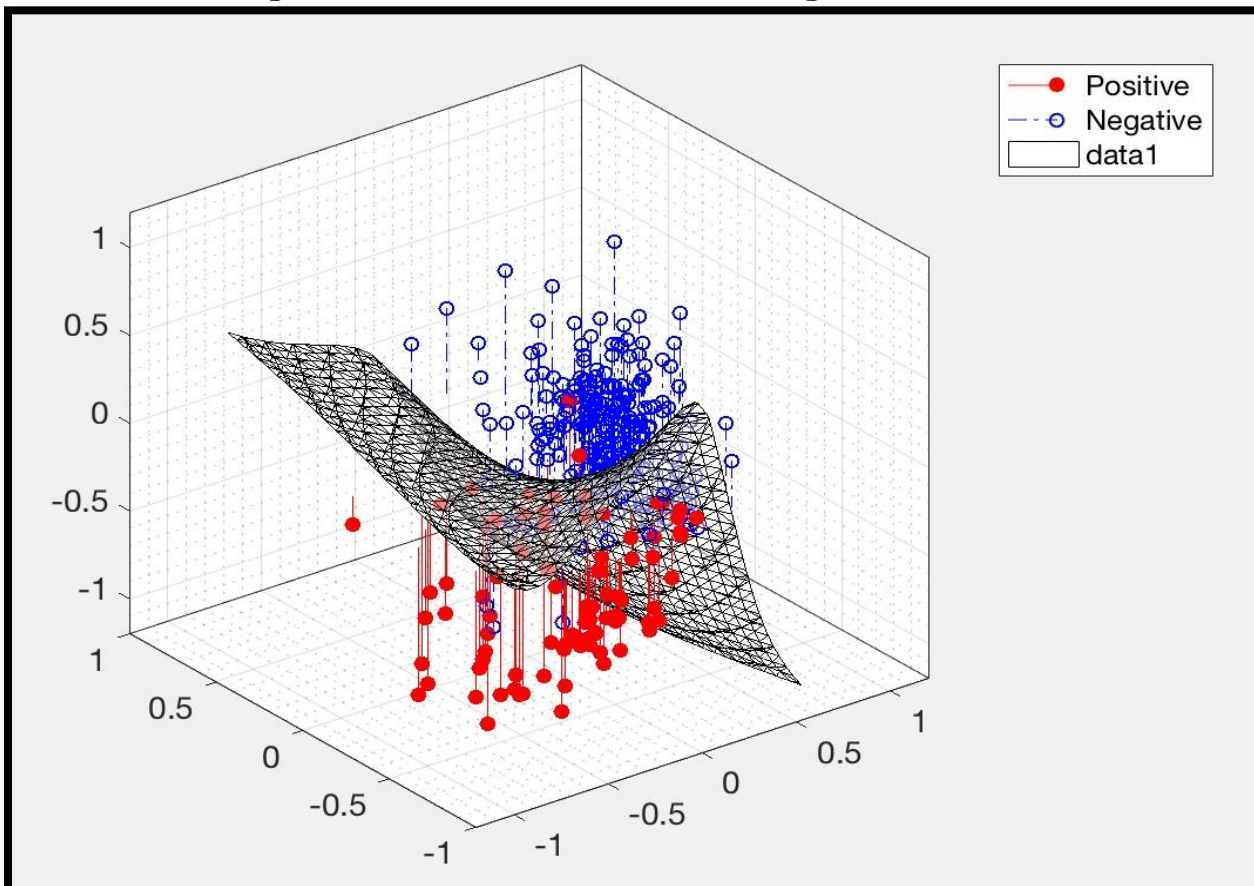
Further we will look for its implementation to solve real world problems and with time we will also explore the types of kernel and many more.

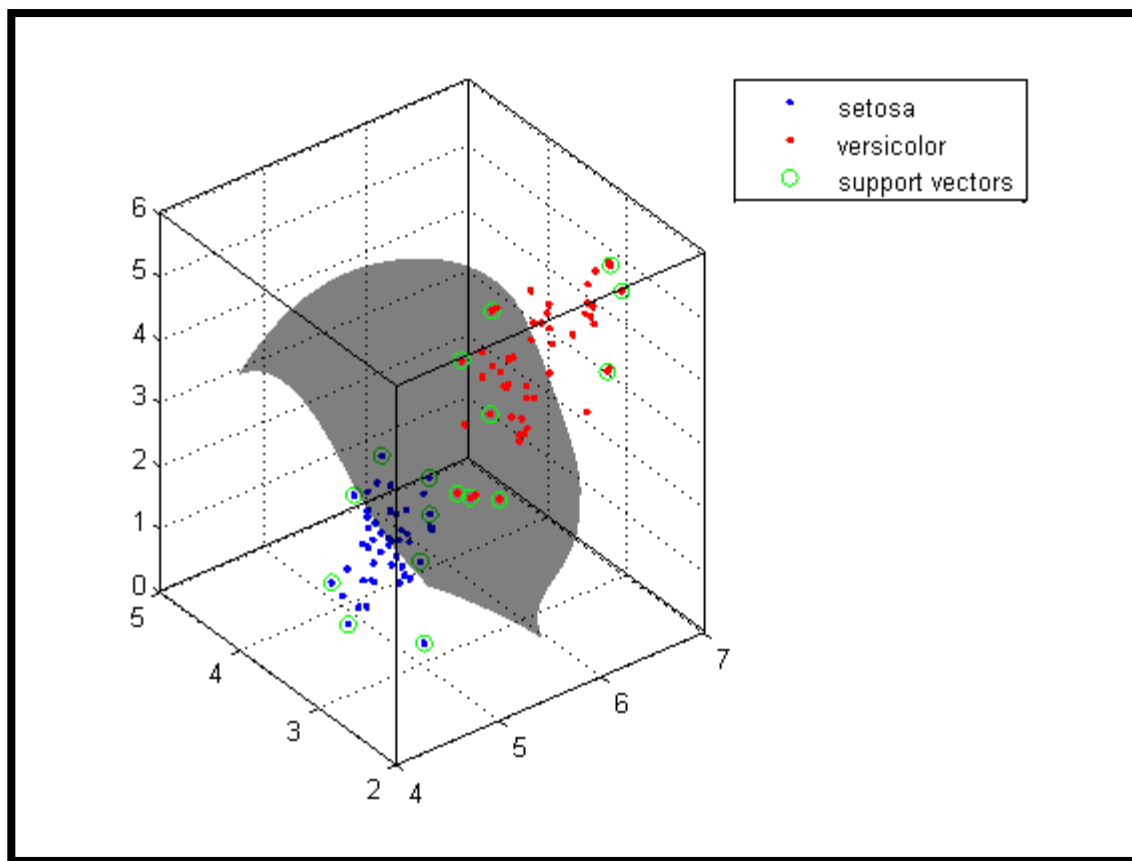
Visuals for summarizing

- What we have looked so far is:



- Some 3d insights to make stuffs interesting:





- Finally, the formulas:

$\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j x_i x_j$	$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i x_j$
Dot product distance or similarity measurement formula.	Kernel looks version
$\sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j f(x_i x_j)$	$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i x_j)$
Constraints from derivatives	
$w = \sum_{i=1}^n \alpha_i y_i x_i$ $\sum_{i=1}^n \alpha_i y_i = 0$	