

⇒ Margin: Distance of a data point from the decision boundary.

Minimum margin: Distance of the decision boundary to the single and closest data point in the training data.

★ One concern is with minimum margin. Sometimes, people use the word just "margin" instead of "minimum margin". We have to find that particular decision boundary whose minimum margin is maximum from both clusters of data. Such boundary lies in middle.

Those nearest data points are called Support Vectors.

In nutshell,

- Find all decision boundary.
- Find minimum margin corresponding to all boundaries.
- Most optimal boundary will be one with maximum value of minimum margin.

OK! Now, we see how to calculate this margin.

⇒ But before that... :-

Equation of a Hyper plane:

$$\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \dots + \omega_n x_n + b = 0$$

$$\Rightarrow \boxed{\vec{\omega}^T \vec{x} = -b}$$

$\vec{\omega}$: weight vector

\vec{x} : features vector

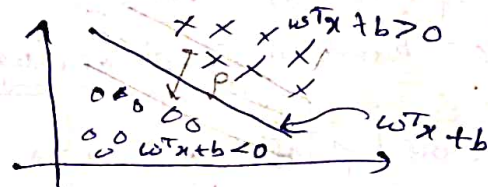
n = number of features

b : bias

$\vec{\omega}^T \vec{x} + b$ → decision hyperplane

$\omega^T x + b < 0$ for negative class

$\omega^T x + b > 0$ for +ve class



⇒ Margin 1.0 :-

Taking $\omega^T x + b$ as our discriminant score for margin.

≠ For any point, higher the value of $\omega^T x + b$, farther is the point from decision boundary.

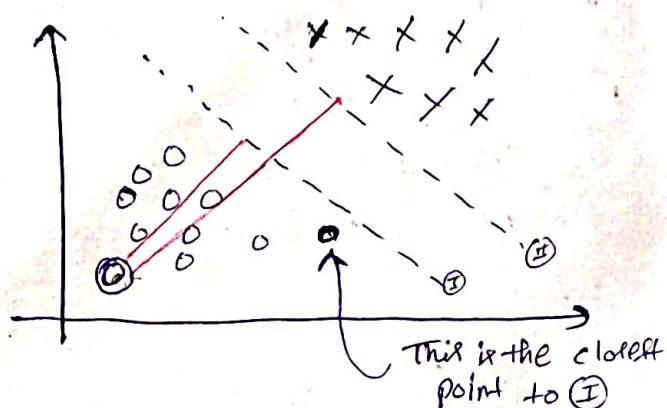
* High +ve value \Rightarrow point is very far and above the decision boundary

* High -ve value \Rightarrow point is very far and below the decision boundary

$$\text{minimum margin} = \min (w^T x + b)$$

Drawback of margin 1.0 :-

Minimum margin is always coming from one single point which is in negative class and deepest.



However, minimum margin is defined for the closest point.

⇒ Bringing constraints for optimal boundary / Need of change :-

Even closest points to decision boundary, should be as away as possible for confident classification.

⇒ For positive class:

$$w^T x + b \gg 0 \quad (\text{discriminant score})$$

⇒ For negative class:

$$w^T x + b \ll 0$$

From now on,

$$w^T x + b \geq 1 \quad \text{for +ve class}$$

$$\text{and, } w^T x + b \leq -1 \quad \text{for -ve class}$$

(constraints for decision boundary)

⇒ Margin 1.1 :-

We should take magnitude of distance even for negative class.

Taking discriminant score as $|w^T x + b|$

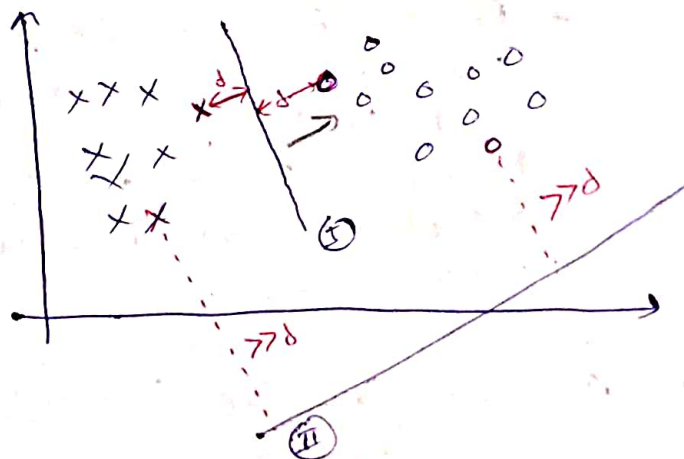
and constraint will become:-
 $|w^T x + b| \geq 1$ (For

$$\text{Min. margin} = \min (|w^T x + b|)$$

Now, the ^{min} margin distance will not be calculated w.r.t deepest point in -ve class but w.r.t to the closest point.

Drawback of margin 1.1 :-

If we go by this definition of margin the our algorithm can learn some poor maximum margin classifier.



Constraint is satisfied for (II) as well. Our algorithm will declare (II) as more better decision boundary than (I) because

min. margin is larger for it.

Here, Negative class is fully misclassified.

⇒ Margin ≥ 0 aka Functional Margin \Rightarrow

$$\text{margin} = y(\omega^T x + b)$$

y : output class label

$$y = \{-1, +1\}$$

Constraints:

$$\text{For all points, } y(\omega^T x + b) \geq 1 \\ (\text{i.e. margin} \geq 1)$$

$$\text{Minimum margin, } F = \min [y(\omega^T x + b)]$$

Peculiarity of this constraint:

If a negative class point is above the decision boundary:
 $\text{margin} = (-ve)(+ve) < 0$
 \Rightarrow Such decision boundary will be rejected

If a positive class point is below the decision boundary:
 $\text{margin} = (+1)(-ve) < 0$
 \Rightarrow Such decision boundary will be rejected.

3. Following things are successfully achieved:-

No misclassification

Minimum margin is calculated from the closest point only.

* \therefore margin value is +ve always.

⇒ Optimal Optimization problem to face \Rightarrow

Find decision boundary (i.e. $\vec{\omega}$, b) for which minimum margin is maximum

or

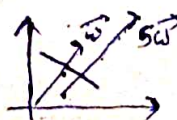
$$\max_{\vec{\omega}, b} [\min (y(\omega^T x + b))]$$

under the constraint:

$$y(\omega^T x + b) \geq 1 \quad \forall \text{ points}$$

⇒ Drawback with Functional Margin \Rightarrow

We can maximize the functional margin simply by scaling $\vec{\omega}$ (normal vector) and b (~~distance from origin or intercept~~ b). But by doing so our decision boundary is not changing.



For same decision boundary we can have many values of F .
So, functional margin is not a good metric.

➔ Margin 3.0 aka Geometric margin \Rightarrow

We constraint our margin with unit length normal vectors.

$$\text{margin} = \frac{y(\vec{w}^T \vec{x} + b)}{\|\vec{w}\|}$$

$$\text{Min. margin}_{(b)} = \frac{\min(y(\vec{w}^T \vec{x} + b))}{\|\vec{w}\|} = \frac{F}{\|\vec{w}\|}$$

For convenience, let us choose to require that functional margin of all data points is atleast 1 and that is equal to 1 for at least one data vector. which, means, for all items in the data:-

$$y_i(\vec{w}^T \vec{x}_i + b) \geq 1$$

and there exist support vectors for which the inequality is an equality.

$$\therefore F = \min(y(\vec{w}^T \vec{x} + b)) = 1$$

$$\text{So, } G = \frac{1}{\|\vec{w}\|}$$

~~Considering margin for both clusters, net margin will be:-~~

$$\rho = \frac{2}{\#}$$

Due to some reasons, a factor of 2 is also introduced

$$\rho = \frac{2}{\|\vec{w}\|}$$

Optimisation problem we got:-

$$\text{maximize } \frac{2}{\|\vec{w}\|} \text{ for } \vec{w}, b$$

$$\text{s.t. } \forall (\vec{x}_i, y_i) \in D$$

$$y_i(\vec{w}^T \vec{x}_i + b) \geq 1$$

Now, problem is that

$$\frac{1}{\|\vec{w}\|} \text{ i.e. } \frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}}$$

is a non-convex function
(such functions have many local optima)

Take reciprocal and remove square root to make it convex

\Rightarrow

$$\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \text{ or } \min_{\vec{w}, b} \frac{1}{2} \vec{w}^T \vec{w}$$

➔ Final Optimization Problem of SVM □

$$\begin{aligned} \min_{\vec{w}, b} & \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t. c:} & \\ \text{For all } (\vec{x}_i, y_i) & \\ y_i (\vec{w} \cdot \vec{x}_i + b) & \geq 1 \end{aligned}$$

➤ Method of Lagrange Multiplier + KKT conditions □

It is a method to optimise f(x) subjected to constraints.

$f(x)$: objective fcn which is to be optimised

constraints:-

$$g_i \leq 0 ; i = 1, 2, \dots, k$$

Lagrangian fcn :

$$L(x, \alpha) = f(x) + \sum_{j=1}^k \alpha_j g_j$$

α : Lagrange multipliers

• We work on ~~off~~ L instead of f then..

(we ≤ 0 in case of $g \geq 0$)

Primal and Dual Concept □

Mini. $f(x)$ (Primal Problem)

s.t. $g_i(x) \leq 0, i = 1, 2, \dots, k$

$$L(x, \alpha) = f(x) + \sum_{j=1}^k \alpha_j g_j \quad (5)$$

~~Mini~~ $\alpha_j \geq 0$ (from perceptron algorithm)

$$\text{Mini. } f(x) \equiv \min_x \left[\max_{\substack{\alpha \\ \text{s.t.} \\ \alpha \geq 0}} L(x, \alpha) \right] \quad (\text{Primal})$$

Reason : $\alpha_j g_j(w) = (+ve)(+ve) = +ve$

So, in order to maximise, we have to let $\alpha = 0$. That will yield original problem only.

↳ when α becomes 0, $\alpha_j g_j = 0$

In general :-

$$\max_{\alpha} \left[\min_x L(x, \alpha) \right] \leq \min_x \left[\max_{\alpha} L(x, \alpha) \right]$$

$$\text{Dual} \leq \text{Primal}$$

We solve dual problem :-

In most cases dual is easier to solve.

Since dual is less than or equal to primal, it is good for us as we are already in search of minimum value in primal problem. Here, we are getting even more less value, which favours us.

Primal = Dual
when KKT conditions are satisfied

KKT conditions :-

1) Convert to Lagrangian fcn.
Take partial partial derivatives
w.r.t variables and equate
them with 0.

2) $\alpha_i g_i = 0$

(Practical see that either $\alpha_i = 0$
or $g_i = 0$ is found, not
both)

3) $g_i \leq 0$

4) $\alpha_i \geq 0$

Take derivative of L w.r.t x and equate to 0.
Eliminate x from dual, then we will left only to
minimize w.r.t α .

Illustration :-

(i) Minimize $w_1^2 + w_2^2$
S.T.C $2w_1 + 3w_2 \geq 1$

Sol:-

$1 - 2w_1 - 3w_2 \leq 0$ { $g \leq 0$ }
form

$L(w_1, w_2, \alpha) = w_1^2 + w_2^2 + \alpha(1 - 2w_1 - 3w_2)$

dual: $\max_{\alpha} \left[\min_{w_1, w_2} L(w_1, w_2, \alpha) \right]$

S.T. $\alpha \geq 0$

$$\begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} 2w_1 - 2\alpha \\ 2w_2 - 3\alpha \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$\Rightarrow w_1 = \alpha$ and $w_2 = \frac{3}{2}\alpha$ (i)

Put these w_1 and w_2 in 2:-

$L = \alpha^2 + \frac{9\alpha^2}{4} + \alpha - 2\alpha^2 - \frac{9\alpha^2}{2}$

$L = -\frac{13\alpha^2}{4} + \alpha$

dual problem: $\max_{\alpha} L(\alpha)$

i.e. $\max -\frac{13}{4}\alpha^2 + \alpha$

S.T. KKT conditions

$\frac{\partial L}{\partial \alpha} = 0$

$\Rightarrow -\frac{13}{2}\alpha + 1 = 0$

$\Rightarrow \alpha = \frac{2}{13}$ { $\alpha \geq 0$ }
satisfied

$\therefore w_1 = \frac{2}{13}$ and $w_2 = \frac{3}{13}$

Minima: $\left(\frac{2}{13}, \frac{3}{13} \right)$

checking if KKT conditions are
satisfied:-

$1 - 2w_1 - 3w_2 = 1 - 2 \times \frac{2}{13} - 3 \times \frac{3}{13}$
 $= 0$

$\therefore g \leq 0$ satisfied

$\alpha_i g_i = 0$ { $\because g = 0$ }
satisfied

(6)

$$\text{i)} \quad \text{Minimize } \omega_1^2 + \omega_2^2$$

$$\text{S.T.C } 2\omega_1 + 3\omega_2 \leq 1$$

sol:-

$$2\omega_1 + 3\omega_2 - 1 \leq 0 \quad (g \leq 0)$$

$$L(\omega_1, \omega_2, \alpha) = \omega_1^2 + \omega_2^2 + \alpha(2\omega_1 + 3\omega_2 - 1)$$

$$\max_{\alpha} \left[\min_{\omega_1, \omega_2} L(\omega_1, \omega_2, \alpha) \right]$$

$$\text{S.T. } \alpha \geq 0$$

$$\frac{\partial L}{\partial \omega_1} = 0 \Rightarrow 2\omega_1 + \alpha = 0$$

$$\Rightarrow \omega_1 = -\frac{\alpha}{2}$$

$$\frac{\partial L}{\partial \omega_2} = 0 \Rightarrow 2\omega_2 + 3\alpha = 0$$

$$\Rightarrow \omega_2 = -\frac{3\alpha}{2}$$

$$\text{So } L = -\frac{13}{4}\alpha^2 - \alpha$$

$$\text{dual:- } \max_{\alpha} -\frac{13\alpha^2}{4} - \alpha$$

$$\text{S.T. } \alpha \geq 0$$

$$\frac{\partial L}{\partial \alpha} = -\frac{13\alpha}{2} - 1 = 0$$

$$\Rightarrow \alpha = -\frac{2}{13}$$

$$\text{But } \alpha \geq 0$$

$$\text{So, take } \alpha = 0$$

$$\Rightarrow \omega_1 = 0, \omega_2 = 0$$

$$\text{Minima:- } (0, 0)$$

$$g \leq 0$$

$$-1 \leq 0$$

✓

$$\alpha g = 0 \quad (\because \alpha = 0)$$

✓

(7)

⇒ Dealing with our SVM optimisation problem

$$\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \quad (\text{Primal})$$

$$\text{S.T.C } y_i (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \forall (x_i, y_i)$$

Lagrange fn:-

$$L(\vec{w}, b, \alpha) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\vec{w}^T \vec{x}_i + b) - 1]$$

$$\min_{\vec{w}, b} \left[\max_{\alpha_1, \alpha_2, \dots, \alpha_n} L(\vec{w}, b, \alpha_1, \alpha_2, \dots, \alpha_n) \right]$$

$$\text{S.T. } \alpha_i \geq 0$$

(Primal)

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_n} \left[\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\vec{w}^T \vec{x}_i + b) - 1) \right]$$

Taking partial derivative w.r.t \vec{w} , b and equating them to 0.

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \vec{w}} = 0 \Rightarrow \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$$

features vector

Now, plug in the value of \vec{w} and $\sum \alpha_i y_i = 0$ in (1.)

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j X_i^T X_j$$

S.T.C. $\forall (x_i, y_i)$

$$\alpha_i \geq 0$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

Dual Problem of SVM

All KKT conditions under which
primal = dual in SVM :-

$$1) \forall (x_i, y_i), y_i (\bar{w}^T \bar{x}_i + b) - 1 \geq 0$$

[Primal problem constraint]

$$\left. \begin{aligned} 2) \frac{\partial L(\bar{w}^*, b^*, \alpha^*)}{\partial \bar{w}} &= 0 \\ 3) \frac{\partial L(\bar{w}^*, b^*, \alpha^*)}{\partial b} &= 0 \end{aligned} \right\} \begin{array}{l} \text{Sufficient conditions} \\ \text{which are sufficient} \\ \text{to find maxima} \\ \text{or minima} \end{array}$$

$$4) \forall_i \alpha_i \geq 0 \quad (\text{from perceptron algorithm})$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

$\alpha = 0$ for almost all data points.
 $\alpha \neq 0$ for only support vectors
 α are weights
 $\alpha = 0 \Rightarrow$ no penalty weight
 $\alpha \neq 0 \Rightarrow$ weight is increased to avoid misclassification

[4) is basically the Dual problem constraint]

$$5) \forall_i \alpha_i [y_i (\bar{w}^T \bar{x}_i + b) - 1] = 0$$

Solution from SVM dual problem

$$\bar{w} = \sum \alpha_i y_i X_i$$

$$b = y_k - \bar{w}^T \bar{x}_k \text{ for any } \bar{x}_k \text{ such that } \alpha_k \neq 0 \Rightarrow \text{support vectors}$$

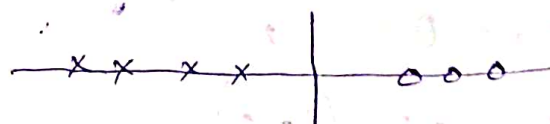
Each non-zero α_i indicates that corresponding \bar{x}_i is support vector

⇒ Most Important Advantage of Dual Problem in SVM :-

we are getting the dot product of feature vectors X_i and X_j

It is known that linear classifier fails to provide good classification. We do feature engineering to map our data to higher dimensions. Then we are able to draw linearly separating hyperplane.

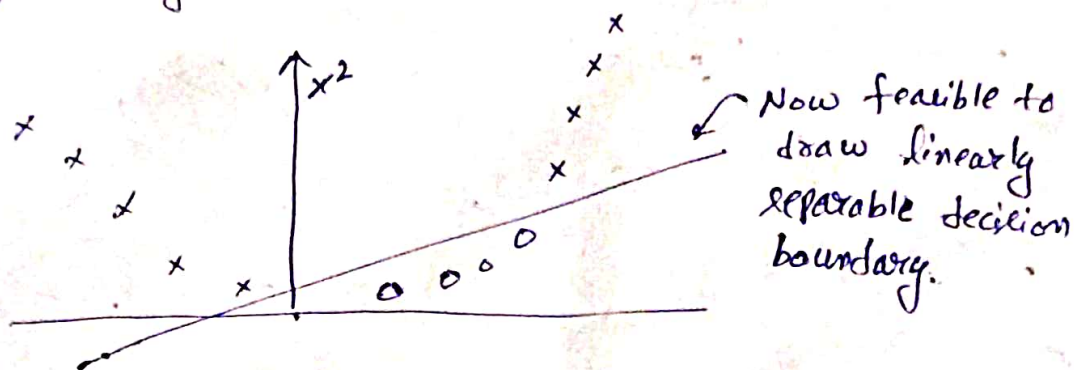
* Linearly Separable:



* Non-linearly separable:-



Projecting into higher dimension



This feature engineering involves hand coding which takes lot of efforts and it increases computational complexity for finding discriminant score.

However, in SVM, $X_i^T X_j$ naturally yields us engineered features without having the need to hand code them just with slight modification while providing fast calculation of dot product at the same time.

Kernel function: It is a function of vectors which projects the two vectors into original or higher dimensional space and takes their dot product (like $X_i^T X_j$) after projecting them and return the result of dot product.

Suppose $X_i = \begin{bmatrix} x_1^i \\ x_2^i \end{bmatrix}$ & $X_j = \begin{bmatrix} x_1^j \\ x_2^j \end{bmatrix}$

and modify $X_i^T X_j$ to Quadratic kernel $(X_i^T X_j + c)^2$

$$K(X_i, X_j) = (X_i^T X_j + c)^2$$

$$= (x_1^i \cdot x_1^j + x_2^i \cdot x_2^j + c)^2$$

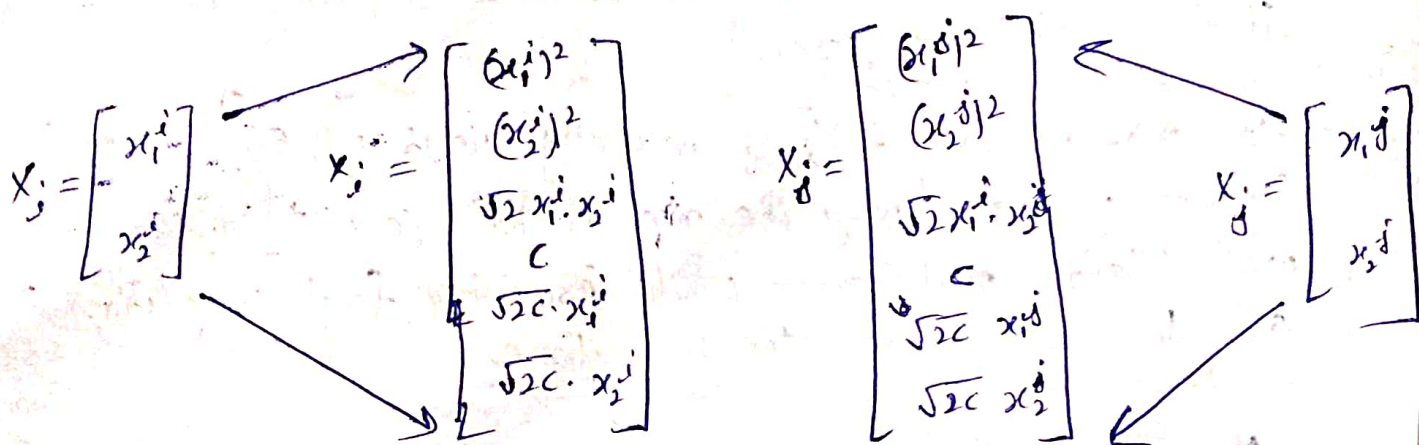
$$= (x_1^i \cdot x_1^j + x_2^i \cdot x_2^j)^2 + c^2 + 2 \cdot c \cdot (x_1^i \cdot x_1^j + x_2^i \cdot x_2^j)$$

(On rearranging further)

$$= (x_1^i)^2 \cdot (x_1^j)^2 + (x_2^i)^2 \cdot (x_2^j)^2 + \sqrt{2} \cdot x_1^i \cdot x_2^i \cdot \sqrt{2} \cdot x_1^j \cdot x_2^j + c \cdot c$$

$$+ \sqrt{2c} \cdot x_1^i \cdot \sqrt{2c} \cdot x_1^j + \sqrt{2c} \cdot x_2^i \cdot \sqrt{2c} \cdot x_2^j$$

Now, we can see feature engineering.



This implicit projection of the two feature vectors into a higher dimensional space while computing different slight modification with constant time complexity, is what we call the **KERNEL TRICK**.

Now, if we can add kernel^{trick} to our optimization dual problem

$$\max_{\alpha_1, \alpha_2, \dots, \alpha_n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{S.T.C } \forall_i \quad \alpha_i \geq 0$$

$$\sum_{j=1}^n \alpha_j y_j = 0$$