The distance of a training point $\underline{x_i}$ to a hyperplane $(\underline{w}, b)$ is

$$\frac{|\langle \underline{w}, \underline{x_i}\rangle + b|}{\|\underline{w}\|}$$

Assuming that $\|\underline{w}\| = 1$, the distance of the closest point to the hyperplane is

$$\min_{i \in [m]} |\langle \underline{w}, \underline{x_i}\rangle + b|$$

The hard SVM learning rule says that pick up the hyperplane which maximizes the margin.

$$\underset{\substack{(\underline{w}, b) \\ \|\underline{w}\|=1}}{\arg\max} \left[ \underbrace{\min_{\underset{\text{Min}}{i \in [m]}} \underbrace{|\langle \underline{w}, \underline{x_i}\rangle + b|}_{\text{distance}}}_{\text{margin}} \right]$$

— abs value.

such that $\forall i \quad y_i (\langle \underline{w}, \underline{x_i}\rangle + b) > 0$

for the separable case when the hyperplane can correctly classify all the examples.

For the separable case we are sure that

$$y_i (\langle \underline{w}, \underline{x_i}\rangle + b) > 0$$

∴ $\underbrace{|\langle \underline{w}, \underline{x_i}\rangle + b|}_{} = y_i \underbrace{(\langle \underline{w}, \underline{x_i}\rangle + b)}$

ground truth prediction.

1$^{st}$ formulation

∴ the equivalent problem is $\underset{\substack{(\underline{w}, b) \\ \|\underline{w}\|=1}}{\arg\max} \underbrace{\min_{i \in [m]} y_i (\langle \underline{w}, \underline{x_i}\rangle + b)}_{\text{margin}}$

Hard SVM learning rule.

In another formulation of Hard SVM, we assume that

$\|\underline{w}\| \neq 1$ ←

Therefore the distance of a point $\underline{x_i}$ to the hyperplane is

$$\frac{|\langle \underline{w}, \underline{x_i}\rangle + b|}{\|\underline{w}\|}$$

The __margin__ for a given training set is therefore

$$\min_{i \in [m]} \left[ \frac{|\langle \underline{w}, \underline{x_i}\rangle + b|}{\|\underline{w}\|} \right]$$

We now assume that the __smallest value for__ achieved through scaling of $\underline{w}$ and $b$.

$$\min_i |\langle \underline{w}, \underline{x_i}\rangle + b| = \min_i y_i(\langle \underline{w}, \underline{x_i}\rangle + b) = 1.$$

i.e. $\forall i \quad y_i(\langle \underline{w}, \underline{x_i}\rangle + b) \geq 1$ ←

Again consider

$$\operatorname*{arg\,max}_{(\underline{w}, b)} \left[ \min_{i \in [m]} \left( \frac{y_i(\langle \underline{w}, \underline{x_i}\rangle + b)}{\|\underline{w}\|} \right) \right]$$

$$= \operatorname*{arg\,max}_{(\underline{w}, b)} \left( \frac{1}{\|\underline{w}\|} \right) \underbrace{\min_{i \in [m]} y_i(\langle \underline{w}, \underline{x_i}\rangle + b)}_{= 1}$$

$$= \operatorname*{arg\,max}_{(\underline{w}, b)} \left( \frac{1}{\|\underline{w}\|} \right)$$
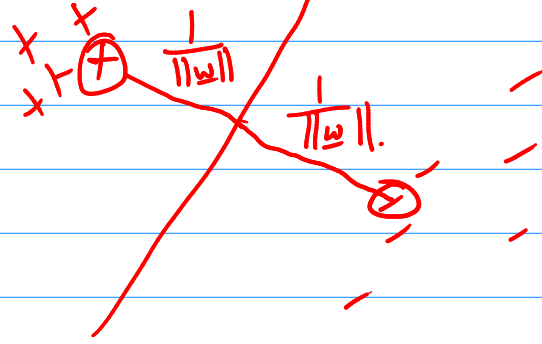
∴ the __2nd formulation of Hard margin learning rule__

$$(\underline{w_0}, b) = \operatorname*{arg\,min}_{(\underline{w}, b)} \|\underline{w}\|^2 \quad \text{such that} \quad \forall i \left( y_i(\langle \underline{w}, \underline{x_i}\rangle + b) \geq 1 \right)$$

Not normalized.   Quadratic objective.

$1\left(\dfrac{1}{\|\underline{w}\|}\right)$ : Margin.

$\dfrac{1}{\|\underline{w}\|}$ : unit to measure the margin.

$$\underline{w} \in \mathbb{R}^d$$

$$\underline{x_i} \in \mathbb{R}^d.$$

In many applications $d$ is very large.

The 2$^{nd}$ formulation is called as the <u>quadratic formulation</u> of the Hard Margin SVM problem.

Final output

$$\hat{\underline{w}} = \frac{\{\underline{w_0}\}}{\|\underline{w_0}\|} \qquad \hat{b} = \frac{b_0}{\|\underline{w_0}\|}$$

Here we have enforced that the <u>margin is 1.</u>

Let's verify that the solution to the first formulation will also be a <u>valid formulation</u> to the 2$^{nd}$ formulation and vice versa.

<u>solution</u>

Let $\underline{w^*, b^*}$ be a solution of the first formulation

$\underline{\|\underline{w^*}\| = 1.}$  Since $\forall i$ $\underline{y_i\left(\langle \underline{w^*}, \underline{x_i}\rangle + b^*\right) \geq 0}$ $\Rightarrow$ Correct classification

Let $\gamma^* = \boxed{\min_i} y_i\left(\langle \underline{w^*}, \underline{x_i}\rangle + b^*\right)$

$\therefore$ $\forall i$ we have $\underline{y_i\left(\langle \underline{w^*}, \underline{x_i}\rangle + b^*\right)} \geq \gamma^*$

Normalizing both sides by $\gamma^*$ we have

$$y_i\left(\langle \frac{w^*}{\gamma^*}, \underline{x_i}\rangle + \frac{b^*}{\gamma^*}\right) \geq 1.$$

These are the constraints of the 2$^{nd}$ formulation.

<u>Candidate</u>

$\therefore$ $\left(\frac{w^*}{\gamma^*}, \frac{b^*}{\gamma^*}\right)$ satisfies all the constraints of the 2$^{nd}$ formulation.

but we are not sure that whether this leads to the min value of the objective fn of the 2$^{nd}$ formulation

$\underline{w}, \hat{b}$ : normalized.

The unnormalized answer given by the 2nd formulation is $\underline{w_0}$.

∴ $\underline{w_0}$ has the minimum norm among all the $\underline{w}$ vectors that satisfy the constraints of the 2nd formulation.

∴ $\left\| \underline{w_0} \right\| \leq \left\| \dfrac{\underline{w}^*}{\gamma^*} \right\| = \dfrac{1}{\gamma^*}$   ∵ $\left\| \underline{w}^* \right\| = 1$

→ min   Candidate

Now consider the final answer of the 2nd formulation

Normalize→ $\underline{\hat{w}}$ and $\hat{b}$   $\left\| \underline{\hat{w}} \right\| = 1$.

We should check whether it satisfies the constraints of the first formulation

$$y_i \left( \langle \underline{\hat{w}}, \underline{x_i} \rangle + \hat{b} \right) = \dfrac{1}{\left\| \underline{w_0} \right\|} \; y_i \left( \langle \underline{w_0}, \underline{x_i} \rangle + b_0 \right)$$

2nd formulation   $\geq 1$

As per the 2nd formulation,

$$y_i \left( \langle \underline{w_0}, \underline{x_i} \rangle + b_0 \right) \geq 1 \quad \forall i$$

$$\therefore \dfrac{1}{\left\| \underline{w_0} \right\|} \; y_i \left( \langle \underline{w_0}, \underline{x_i} \rangle + b_0 \right) \geq \dfrac{1}{\left\| \underline{w_0} \right\|} \geq \gamma^*$$

$\geq 1$   $\gamma^* > 0$ positive.   ∵ $\left\| \underline{w_0} \right\| \leq \dfrac{1}{\gamma}$

∴ $\left( \underline{\hat{w}}, \hat{b} \right)$ satisfy the constraints of the 1st formulation.

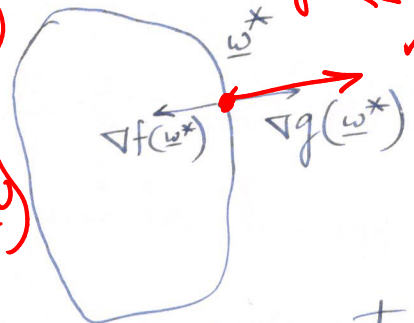further, $\left\| \underline{\hat{w}} \right\| = 1$

$\Rightarrow \left( \underline{\hat{w}}, \hat{b} \right)$ is an optimal solution to the first formulation.

# Dual formulation of SVM

$\min f(\underline{w})$ subject to $g(\underline{w}) \leq 0$

$h(\underline{w}) \geq \gamma.$

$g(\underline{w}) \leq 0$  $\underline{w}^+$  $\underline{w}^+$

$g(\underline{w}) = 0$

$g_1(\underline{w})$
$g_2(\underline{w})$
$g_3(\underline{w})$

$f(\underline{w})$
feasible region

$g(\underline{w}) \leq 0$

$g(\underline{w}) = 0$

$\underline{w}^*$
$\nabla f(\underline{w}^*) = 0$

$g(\underline{w}) = 0$

$\underline{w}^*$
$\nabla f(\underline{w}^*)$  $\nabla g(\underline{w}^*)$

constraint is active.
because $\underline{w}^*$ satisfies $g(\underline{w}^*) = 0$

$\alpha$ : scalar

$\nabla f(\underline{w}^*) = - \alpha \nabla g(\underline{w}^*)$

Hold at $\underline{w}^* \longrightarrow$ $\nabla f(\underline{w}^*) + \alpha \nabla g(\underline{w}^*) = 0$

$g(\underline{w}) = 0$

When there are several constraints $g_i(\underline{w}) \leq 0$, the feasible region is the intersection of all the regions $g_i(\underline{w}) \leq 0$. $\forall i$

Assuming that $f$ and $g_i$ are differentiable,

$$\nabla f(\underline{w}^*) = - \sum_{i \in I} \alpha_i g_i(\underline{w}^*)$$

where $I$ is the set of constraints which are active at $\underline{w}^*$

$$\nabla f(\underline{w}^*) + \sum_{i \in I} \alpha_i g_i(\underline{w}^*) = \underline{0}$$

The Hard SVM learning rule specifies
$$f(\underline{w}) = \frac{1}{2} \|\underline{w}\|^2 \quad \text{and} \quad g_i(\underline{w}) \leq 0 \quad g_i = 1 - y_i \langle \underline{w}, x_i \rangle$$