Label 1    $\langle \underline{w}_1, \underline{x} \rangle + b_1$

Label 2    $\langle \underline{w}_2, \underline{x} \rangle + b_2$

Label 3    $\langle \underline{w}_3, \underline{x} \rangle + b_3$.

$\underline{x}$ [ ... ] $w \times h$

Linear discriminant function

prediction
$= \underset{i}{argmax} \langle \underline{w}_i, \underline{x} \rangle$

$W = np.random.randn(3, 3072)$

$\xrightarrow{3072}$

$\begin{matrix} 1 \\ 2 \\ 3 \end{matrix}$ [ - - - - ] [ ] $\underline{x}$ $+ \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix}$

$32 \times 32$

$1024 \times 3$
$= 3072$

$\begin{bmatrix} (r,g,b) & (r,g,b) & - & - & - & \cdots & 32 \\ (rgb) & (rgb) & - & - & \cdots & \cdot \\ & & \cdot \end{bmatrix}$

Image

Flatten    $\begin{bmatrix} r \\ g \\ b \end{bmatrix}$ 1st pixel  $\begin{bmatrix} r \\ g \\ b \end{bmatrix}$ 2nd pixel

$3072$ [ ⋮ ]   1-D vector
flattened vector

Writing the mapping $\underline{x} \longmapsto \langle \ddot{\underline{w}}, \ddot{\underline{x}} \rangle$ makes it
a <u>homogeneous linear</u> function.
  i.e. a homogeneous map.
What's the advantage of writing a homogeneous form?
We can <u>scale the vector $\underline{w}$</u> to impose a desired
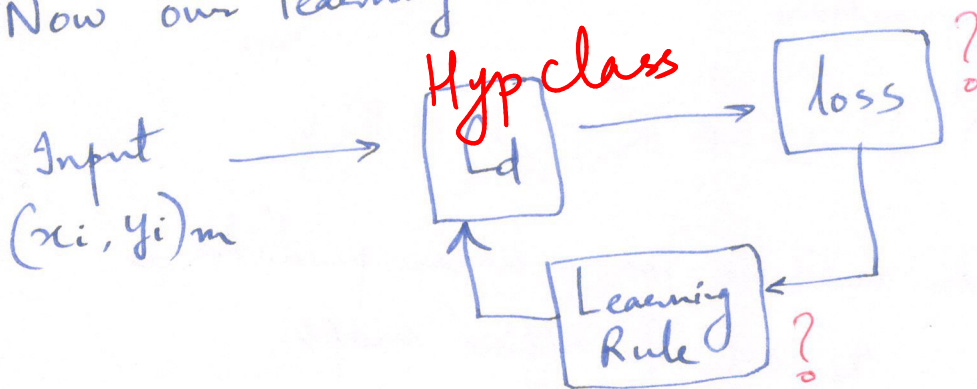minimum value

$$\left.\begin{array}{l} \langle \underline{w_1}, \underline{x_1} \rangle = 10 \\ \langle \underline{w_2}, \underline{x_2} \rangle = 7 \\ \langle \underline{w_1}, \underline{x_3} \rangle = \boxed{\textcircled{3}} \end{array}\right]$$

If we insist that
the minimum value
should <u>be 1</u>

$\textcircled{$\underline{w}$} \leftarrow \frac{1}{3} \underline{w_1}$

<span style="color:red">scaled</span>

then
$$\langle \underline{w}, x_1 \rangle = 10/3$$
$$\langle \underline{w}, x_2 \rangle = 7/3$$
$$\langle \underline{w}, x_3 \rangle = 3/3 = 1. \quad \longleftarrow \text{<span style="color:red">min value 1.</span>}$$

Now our learning model is

<span style="color:red">Hyp class</span>

Input $\longrightarrow$ [ $d$ ] $\longrightarrow$ [ loss ] <span style="color:red">?</span>
$(x_i, y_i)_m$ $\longrightarrow$ [ Learning Rule ] <span style="color:red">?</span>

The loss function to be used depends on the final
task.

<u>Task 1</u> : Binary Classification <span style="color:red">2 classes.</span>

$$h_{\underline{w}}(\underline{x}) \in \textcircled{$\mathbb{R}$} \longmapsto \{ +1, -1 \}$$
<span style="color:red">linear Hypothesis</span>  by checking  the sign.  $\{-,-,\}$

<span style="color:red">+1</span>
<span style="color:red">-1.</span>

(5)

i.e. if $\underline{h_w(x)} \gtreqless 0 \xmapsto{\text{map to}} +1$

$h_w(x) < 0 \xmapsto{\quad} -1$ $\Bigg\}$ $\underline{\text{sign}(h_w(x))}$

$h_{w_1}(x)$
$h_{w_2}(x)$
$h_{w_3}(x)$

i.e. $\underline{\textcircled{x}} \xmapsto{\quad} \text{sign}(h_w(x))$ $\to d$ parameters $\quad$ 2d

binary $\quad h_{w_1}(x)$

$h_{w_2}(x)$

Linear Hyperplane in d-dim space

$x_2$

$0^+ \langle w, x \rangle$

+ve half $\quad ax_1 + bx_2 + c = 0$

$0^+ \quad \langle \underline{w}, x \rangle = 0$

$\langle w, x \rangle > 0$

$0 - \quad \textcircled{o} \quad \textcircled{o}$

$\langle w, x \rangle < 0$

-ve half.

$x_1$

In general, a d-dim hyperplane divides the d-dim plane into 2 halves.

$x_2$

point/dot an example in d-dim feature space

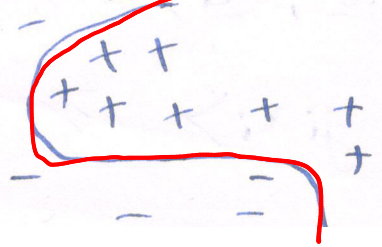Several solutions possible

Separable

$x_1$

This is a simple training dataset where the dataset is separable, i.e, the positive and negative points are well separated.

Data can also be non-separable

Not linearly separable

This requires a non-linear separating boundary

The separable case is also called as the
Realizable case

For this case, the best hypothesis $h \in \mathcal{L}_d$ has
$$error(h) = 0$$

The non separable data corresponds to the non realizable
case. Also called as Agnostic scenario.

The best hypothesis $h \in \mathcal{L}_d$ will have $error(h) > 0$

$S : (x_i, y_i)_m$

How to find the best hypothesis in $\mathcal{L}_d$ for the
Realizable case ?

$\underline{w}$ gives the best hypothesis

The (best) hypothesis will make sure that

100% correct classification $\quad sign(<\underline{w}, x_i>) = y_i \qquad \forall i = 1 \dots m.$
scalar $\qquad\qquad\qquad$ for all examples

condition for i.e. $\rightarrow \boxed{y_i <\underline{w}, x_i> \quad > 0} \quad \forall i = 1 \dots m \quad y_i \in \{+1, -1\}$
correct classification $\qquad\qquad$ LMS $\qquad\qquad\qquad$ set of constraints

For homogeneous linear functions we can scale the
parameter vector $\qquad \underline{w} \longrightarrow \underline{\breve{w}} \qquad \widehat{\underline{w}}$

such that $\qquad \boxed{y_i <\underline{\breve{w}}, x_i> \quad \geqslant 1} \, \forall_i$

Denoting $\underline{w} \equiv \underline{\breve{w}}$ we require for correct classification
of all data points $\qquad y_i <\underline{w}, x_i> \geqslant 1 \qquad \forall_i$

$y_i \Rightarrow \begin{bmatrix} x_i^1 \\ x_i^2 \\ x_i^3 \\ x_i \\ x_i \end{bmatrix} x_i \bigg\} d \qquad \boxed{<\underline{w}, y_i x_i> \quad \geqslant 1} \qquad \forall_i \quad \substack{m \\ constraints}$

100 examples

$$y_1 \langle \underline{w}, \underline{x_1} \rangle = 10 \qquad \langle \underline{\tilde{w}}, \underline{x_1} \rangle = 10/6$$

$$y_2 \langle \underline{w}, \underline{x_2} \rangle = 13 \qquad \underline{x_2} = 13/6$$

$$y_3 \langle \underline{w}, \underline{x_3} \rangle = \boxed{6} \leftarrow \text{min} \qquad \underline{x_3} = 6/6 =$$

$$\underset{\underset{\text{scaling}}{\uparrow}}{\underline{\tilde{w}}} = \frac{1}{6} \underline{w}$$

$$\frac{x_4}{\vdots} = 23.$$

$$100$$

$$\begin{bmatrix} y_1 x_{11} & y_1 x_{12} & y_1 x_{13} & \cdots & y_1 x_{1d} \\ y_2 x_{21} & y_2 x_{22} & & \cdots & y_2 x_{2d} \\ & & \ddots & \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \geq 1$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxx}}_{A} \qquad \underbrace{\phantom{xxx}}_{w}$$

$$A\,\underline{w} \geq \underline{v} \qquad \text{where} \quad \underline{v} \equiv \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{array}{l} \text{vector} \\ \text{of 1s.} \end{array}$$

**LP solvers.**

The standard <u>linear programming problem</u> is

<u>Input</u>
$\begin{array}{c} \underline{u} \\ A \\ \underline{v} \end{array}$

<u>Output</u>
$\underline{w}$

$$\underset{\underline{w}}{\text{maximize}} \quad \left( \underline{u}^T \underline{w} \right)$$

subject to linear inequality constraints

$$\boxed{A\,\underline{w} \geq \underline{v}}$$

$$\underline{u} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} \underline{w} \end{bmatrix} \begin{bmatrix} \phantom{x} \end{bmatrix}$$

So we design a linear program for our problem:

$$\begin{bmatrix} \underset{\underline{w}}{\max} & \boxed{\text{constant } 1} & \text{We have nothing to maximize} \\ \text{subject to} & \langle \underline{w}, y_i \underline{x}_i \rangle \geq 1 & \forall i \end{bmatrix}$$

solution given by the LP solver is the weight vector $\underline{w}$

Using the weight vector, any given point $\underline{x}$ can be classified using $\text{sign}(\langle \underline{w}, \underline{x} \rangle)$ ] prediction

LP solver gives one solution.
There may be many solutions

Generalization

Gap

Gap

S: training set

best
$$h^* = \arg\min_{h \in H} error_S(h)$$

ERM learning rule        M: minimization

Empirical Risk (expected loss)

training set