Null model Alg.
$$\mathcal{H} = \{0,1\}$$
$$g = A(\mathbb{D}, \mathcal{H}) = \text{Mode}[\vec{y}]$$

most likely category $\in \{0,1\}$
$\vec{x}$'s don't matter in null model    A

← Always keep
Null model to beat

Let's return to

$$\mathcal{H} = \left\{ \mathbb{1}_{\vec{w} \cdot \vec{x} > 0} : \vec{w} \in \mathbb{R}^{p+1} \right\}$$

$\vec{x}$ includes a 1 ⎯⎯⎯⎯⎯↑

Assume "linearly separable" i.e. $\exists \vec{w}$

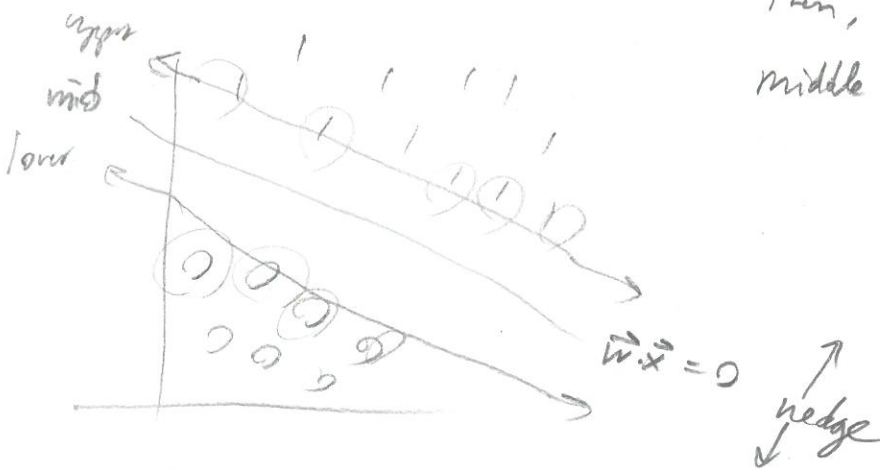s.t. there would be NO errors if
$g$ were used on obs's in $\mathbb{D}$.

Consider a

New $A$, different from perceptron learning algorithm.

Why not create a wedge, large as possible using parallel hyperplanes.



$\vec{w} \cdot \vec{x} = 0$

Which hyperplane (i.e. $\vec{w}$) is "best"?

Then, $g$ is built from the $\vec{w}$ in the middle of this wedge.

The "max-margin hyperplane"
(proven to be optimal linear classifier in 1998)



upper
mid
lower

$\vec{w} \cdot \vec{x} = 0$

↑ wedge

Which data pts most matter? Since data pts $\vec{x}_i$ are "vectors",
these are called the "support vectors". The midline model
is then called the "support vector machine"

weird name ... if I were naming this now...

"support vector" ⇒ "essential observation"
"machine" ⇒ "model" + "separation"

How to fit this?
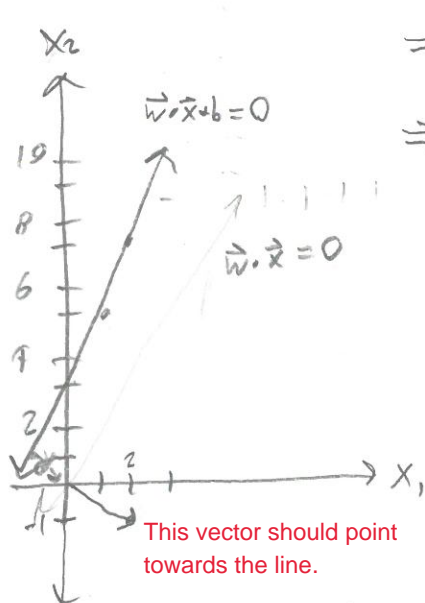
Unfortunately, it is convenient to revert back to the ugly slope-intercept form,

$$\mathcal{H}(= \left\{ \mathbb{1}_{\vec{w}\cdot\vec{x}+b > 0} : \vec{w} \in \mathbb{R}^p, b \in \mathbb{R} \right\}$$

this is equivalent to before, it's just a slight reparameterization...

Let's first review $8^{th}$ grade math.. Imagine the line $x_2 = 2x_1 + 3$



$$\Rightarrow 2x_1 - x_2 + 3 = 0$$

$$\Rightarrow \underbrace{\begin{bmatrix} 2 \\ -1 \end{bmatrix}}_{\vec{w}} \cdot \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_{\vec{x}} + \underbrace{3}_{b} = 0$$

$$\vec{w} \cdot \vec{x} + b = 0$$

where is $\vec{w}$ on this graph? It's the "normal vector"

i.e. perpendicular to the line

If $p > 2$, $\vec{w}\cdot\vec{x} = 0$ is a plane/hyperplane and $\vec{w}$ is $\perp$ to it.

Hesse Normal Form is wx - b = 0. Then the w vector will point towards the line. So the "b" here should be "-b".

This vector should point towards the line.

Note: $\|\vec{w}\|$ indicates length of the vector $:= \sqrt{\sum_{j=1}^{p} w_j^2}$

And the normalized $\vec{w}$ vector is defined as the vector in same direction with length 1.

$$\vec{w}_0 := \frac{\vec{w}}{\|\vec{w}\|} \qquad Proof: \quad \|\vec{w}_0\| = \sqrt{\sum_{j=1}^{p}\left(\frac{w_j}{\|\vec{w}\|}\right)^2} = \sqrt{\frac{1}{\|\vec{w}\|^2}\sum_{j=1}^{p} w_j^2} = \sqrt{\frac{1}{\|\vec{w}\|^2}\|\vec{w}\|^2} = 1$$

What is the length of the line $\vec{\ell} = \alpha\vec{w}_0$ where $\alpha \in \mathbb{R}$, constant?

$$\|\vec{x}\| = \sqrt{\sum_{j=1}^{p}(\alpha w_{0j})^2} = \sqrt{\alpha^2\sum_{j=1}^{p} w_{0j}^2} = \sqrt{\alpha^2\|\vec{w}_0\|^2} = |\alpha| \quad \text{Makes sense..} \quad \alpha \cdot 1 = 1 !$$

let $\vec{\ell}$ be the vector from the origin to the line $\vec{w}\cdot\vec{x}+b = 0$, perpendicular to it

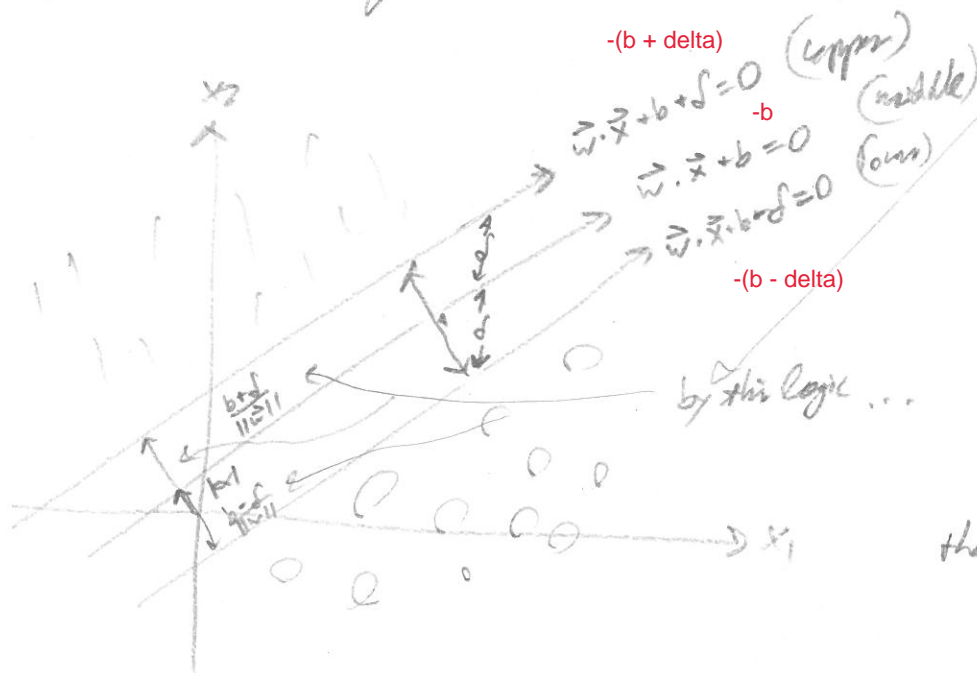Let's solve for $\alpha$.. $\vec{\ell}$ is on the line $\vec{w}\cdot\vec{x}+b = 0$   "-b"

$$\Rightarrow \quad \vec{w} \cdot \vec{l} + b = 0 \quad \Rightarrow \quad \vec{w} \cdot \alpha \vec{w}_0 + b = 0 \qquad \text{"-b"}$$

$$\Rightarrow \vec{w} \cdot \alpha \frac{\vec{w}}{\|\vec{w}\|} + b = 0 \quad \Rightarrow \quad \alpha \frac{\|\vec{w}\|^2}{\|\vec{w}\|} + b = 0 \quad \Rightarrow \quad \alpha = + \frac{b}{\|\vec{w}\|}$$

<span style="color:red">No negative sign here.</span>

$$\Rightarrow |\alpha| = \frac{b}{\|\vec{w}\|}$$

Now back to finding that "best" line:



$\vec{w} \cdot \vec{x} + b + \delta = 0$ (upper)  <span style="color:red">-(b + delta)</span>

$\vec{w} \cdot \vec{x} + b = 0$ (middle)  <span style="color:red">-b</span>

$\vec{w} \cdot \vec{x} + b - \delta = 0$ (lower)  <span style="color:red">-(b - delta)</span>

by this logic ...

distance between upper & lower lines is

$$\frac{b + \delta}{\|\vec{w}\|} - \frac{b - \delta}{\|\vec{w}\|} = \frac{2\delta}{\|\vec{w}\|}$$

<span style="color:red">this is still true now since alpha has correct sign</span>

therefore

Since $C(\vec{w} \cdot \vec{x} + b) = 0$, there are infinite solutions since $C \in \mathbb{R}$  <span style="color:red">-b</span>

Coerce $\delta = 1$ ... now there's a unique solution to the equation

$$\vec{w} \cdot \vec{x} + b + \delta = 0 \quad \Rightarrow \quad \vec{w} \cdot \vec{x} + b + 1 = 0$$

<span style="color:red">-(b + delta)</span> ... <span style="color:red">-b</span>

$$\Rightarrow \text{margin} = \frac{2}{\|\vec{w}\|}$$

Constrain all $y=1$'s to be $\geq$ upper; constrain all $y=0$'s to be $\leq$ lower

$$\vec{w} \cdot \vec{x} \,\,\text{<span style="color:red">-(b+1)</span>} \geq 0 \quad \Rightarrow \quad \vec{w} \cdot \vec{x} - b \geq {}_{+} 1 \quad \Rightarrow \forall i \text{ s.t. } y_i = 1 \quad \vec{w} \cdot \vec{x}_i - b \geq {}_{+} 1$$

multiply both sides by $\left(y_i - \frac{1}{2}\right)$

$$\left(y_i - \frac{1}{2}\right)\left(\vec{w} \cdot \vec{x}_i + b\right) \geq \left(y_i - \frac{1}{2}\right) \quad \Rightarrow \quad \left(y_i - \frac{1}{2}\right)\left(\vec{w} \cdot \vec{x}_i - b\right) \geq {}_{+}\left(1 - \frac{1}{2}\right) \quad \text{<span style="color:red">now positive</span>}$$

Now $y_i = 1$

$$\Rightarrow \left(y_i - \frac{1}{2}\right)\left(\vec{w} \cdot \vec{x}_i - b\right) \geq {}_{+} \frac{1}{2} \qquad \text{<span style="color:red">now positive</span>}$$

Now $\forall i$ s.t $y_i = 0$   $\vec{w} \cdot \vec{x}_i$ -(b-1) $\geq 0 \Rightarrow \vec{w} \cdot \vec{x}_i - b \geq 1$   <span style="color:red">now negative</span>

multiply both sides by $(y_i - \frac{1}{2})$

$(y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i + b) \geq -(y_i - \frac{1}{2})$    Now $y_i = 0 \Rightarrow$   <span style="color:red">now positive</span> $(y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i + b) \geq$ <span style="color:red">+</span>$\frac{1}{2}$

Same condition for $y \in \{0, 1\} \Rightarrow \forall i$

<span style="color:red">Condition of perfect separability. If the pt on the line, = +½. If not, > +½.</span>

Then we solve the following optimization problem:

Maximize $\frac{2}{\|\vec{w}\|} \Rightarrow$ | Minimize $\|\vec{w}\|$ subj. to. $\forall i$ $(y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i + b) \geq$ <span style="color:red">+</span>$\frac{1}{2}$
                 over $\vec{w} \in \mathbb{R}^p$, $b \in \mathbb{R}$

the $\{\vec{w}, b\}$ solution is the support vector machine

This approach assumes perfect linear separability. In the real world... who has that luxury? We need to upgrade it.

_____

We need a loss function. Previously we employed $SAE = \sum_{i=1}^{n} \mathbb{1}_{\hat{y}_i \neq y_i}$ and then allowed the computer to find $\vec{w}$. This is okay but we can do better. We can make the loss depodent on how bad the error is. Consider the following for the $i^{th}$ obs:

$H_i := \max \left\{ 0, \text{ } \color{red}{+}\color{black}\frac{1}{2} - (y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i - b) \right\}$    "hinge loss"

Let's see why this makes sense...

What are the params? $\vec{w} \in \mathbb{R}^p$, $b \in \mathbb{R}$. Still $p+1$ parameters.

What is $\lambda$? A predefined constant called a "hyperparameter."

$$g = A(D, \mathcal{H}, \lambda, ...)$$

It is considered a tuning knob on the $A$ itself. It is a meta idea. recall $g(\vec{x}) = \mathbb{1} \vec{w} \cdot \vec{x} + b$  there is no $\lambda$ here! Perceptron's meta idea could be a $\lambda$.

$\lambda$ only affects which $g \in \mathcal{H}$ will be selected.

We will discuss how the value of hyperparameters are selected later in the course. For now, who does $\lambda$ do?

- If $\lambda \stackrel{\sim}{=} 0$, we only care about errors and not a max. margin. One error for any can ruin our nice separation line.

- If $\lambda \stackrel{\sim}{=} \infty$, we only care about the best line of separation and not about errors... this makes no sense! I can just make little o

Considering $\lambda$ is picked "reasonably". How do we solve for $\{\vec{w}, b\}$ our params? Use modern numerical optimization methods:

  - quadratic programming
  - sub-gradient descent
  - coordinate descent

which we will likely not study. Lucky for us, R packages already implemented this for us.