$$y = \{0, 1\}, \quad p+1 = 3, \quad \mathcal{H} = \{ \mathbb{I} \; \vec{w} \cdot \vec{x} \geq 0 ; \; \vec{w} \in \mathbb{R}^3 \}$$
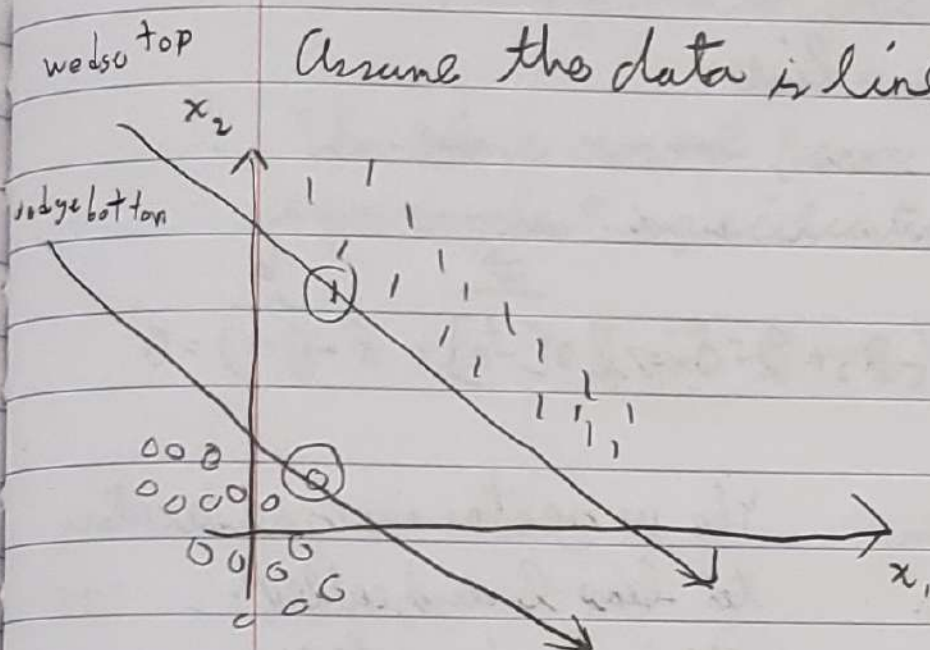
Assume the data is linearly separable so it looks like:

wedge top

wedge bottom



We need an algorithm that locates the middle of that wedge. Let the top of the wedge be the linearly separable model "closest" to the $y=1$'s and the bottom of the wedge be the linearly separable model "closest" to the $y=0$'s. The "max margin hyperplane" is the parallel line in the center of the top and bottom.
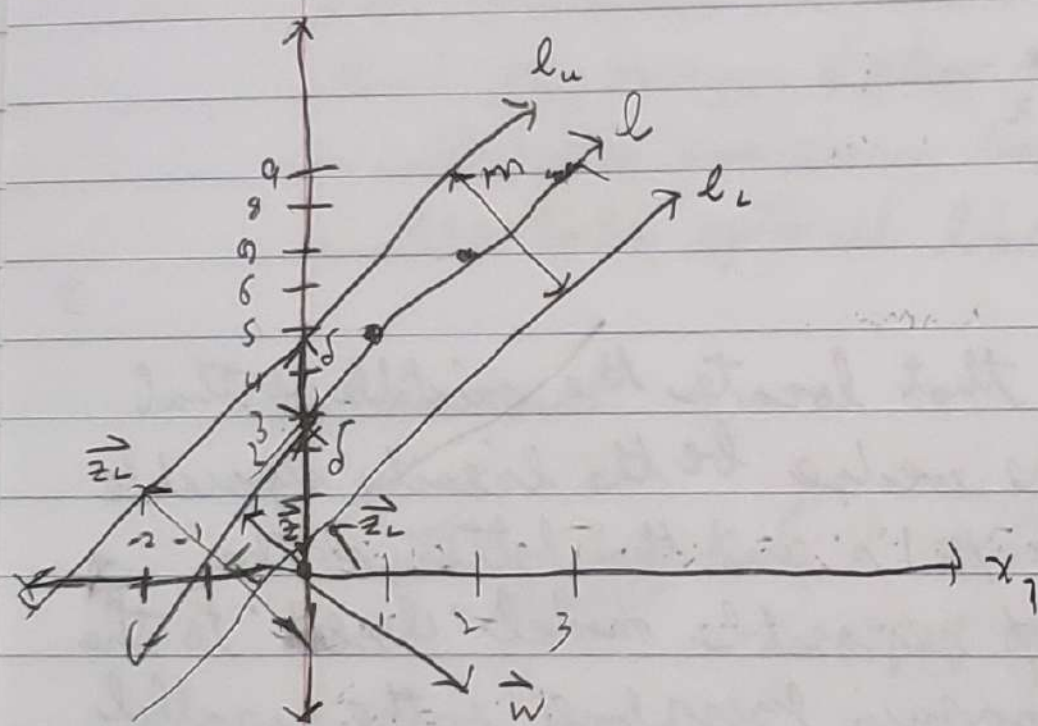
note: there are two critical observations (the circled points). Since observations are $x$-vectors, these critical observations are called "support vectors" and hence the final model is called a "support vector machine" (SVM). "machine" is a fancy word meaning "complex model"; So "machine learning" just means "learning complex models". To find the SVM...

First rewrite $\mathcal{H} = \{ \mathbb{I} \; \vec{w} \cdot \vec{x} - b = 0 \; ; \; \vec{w} \in \mathbb{R}^p, \; b \in \mathbb{R} \}$

Note $\underline{\vec{w} \cdot \vec{x} - b = 0}$ defines a line

Here normal form.

$l : x_2 = 2x_1 + 3 \implies l : 2x_1 - x_2 + 3 = 0 \implies l : \overset{\vec{w}}{\overbrace{\begin{bmatrix} 2 \\ -1 \end{bmatrix}}} \cdot \vec{x} - \overset{b}{\overbrace{(-3)}} = 0$



The $w$ vector is perpendicular to line $l$ and called the 'normal vector'

Let $\vec{w_0} := \dfrac{\vec{w}}{\|\vec{w}\|}$

The direction of the $w$ vector with unit length

$\vec{z} = \alpha \vec{v_0}, \quad \vec{z} \in l$

$\vec{w} \cdot \vec{z} - b = 0$

$\Downarrow$

$\vec{w} \cdot (\alpha \vec{v_0}) - b = 0$

$\implies \alpha \; \dfrac{\|\vec{w}\|^2}{\|\vec{w}\|} - b = 0$

$\implies \alpha = \dfrac{b}{\|\vec{w}\|} \implies \vec{z} = \dfrac{b}{\|\vec{w}\|} \vec{w_0}$

Let $m > 0$ be the perpendicular distance between $l_u$ and $l_L$ and let $\delta > 0$ be the distance between $l_u$ and $l$ (and $l_L$ and $l$) on the $x_2$ axis

$l_u : \vec{w} \cdot \vec{x} - (b + \delta) = 0, \quad \vec{z_u} = \dfrac{b + \delta}{\|\vec{w}\|} \vec{w_0}$

$l_L : \vec{v} \cdot \vec{x} - (b - \delta) = 0, \quad \vec{z_L} = \dfrac{b - \delta}{\|\vec{w}\|} \vec{w_0}$

$\Downarrow$

$M = \|\vec{z_u} - \vec{z_L}\| = \left\| \dfrac{b + \delta}{\|\vec{w}\|} \vec{w_0} - \dfrac{b - \delta}{\|\vec{w}\|} \vec{w_0} \right\|$

$= \dfrac{1}{\|\vec{w}\|} \, 2\delta \|\vec{w_0}\| = \dfrac{2\delta}{\|\vec{w}\|}$

Goal is to make $m$ as large as possible (maximum margin) $\iff$ making the $w$ vector as small as possible.

The Hesse normal form is not unique. There are infinite equivalent specification of a line:

$$\forall c \neq 0 \quad c(\vec{w} \cdot \vec{x} - b) = 0 \quad \text{Let } c = \frac{1}{g}$$

$$\Downarrow$$

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

Now we need two conditions

(I) All $y = 1$'s are above or equal to $l_1$:

$$\forall i \text{ st } y_i = 1 \quad \vec{w} \cdot \vec{x}_i - (b+1) \geq 0 \Rightarrow \vec{w} \cdot \vec{x}_i - b \geq 1 \Rightarrow \frac{1}{2}(\vec{w} \cdot \vec{x}_i - b) \geq \frac{1}{2}$$

$$\Downarrow$$

$$(y_i - \tfrac{1}{2})(\vec{w} \cdot \vec{x}_i - b) \geq \frac{1}{2}$$

(II) All $y = 0$'s are below or equal to $l_1$:

$$\forall i \text{ st } y_i = 0 \quad \vec{w} \cdot \vec{x}_i - (b-1) \leq 0 \Rightarrow \vec{w} \cdot \vec{x}_i - b \leq -1 \Rightarrow \frac{1}{2}(\vec{w} \cdot \vec{x}_i - b)$$

$$\leq -\frac{1}{2}$$

$$\Rightarrow -\frac{1}{2}(\vec{w} \cdot \vec{x}_i - b) \geq \frac{1}{2}$$

$$\Downarrow$$

$$(y_i - \tfrac{1}{2})(\vec{w} \cdot \vec{x}_i - b) \geq \frac{1}{2}$$

note how both inequalities are the same for both I and II. Thus this inequality satisfies both constraints. So all observation will be in their right places

$$\forall i \; (y_i - \tfrac{1}{2})(\vec{w} \cdot \vec{x}_i - b) \geq \frac{1}{2} \Rightarrow \text{line is linearly separable.}$$
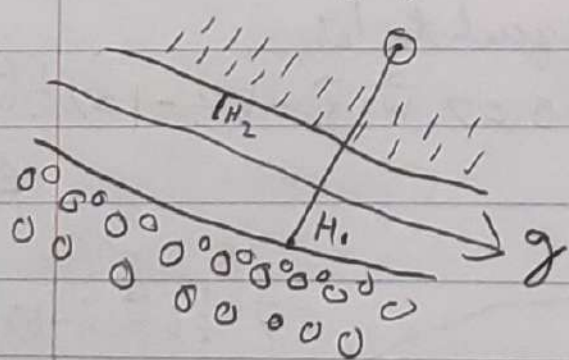
You compute the SVM by optimizing the following problem:

$$\min \|\vec{w}\| \quad \text{st} \quad \forall i \ \left(y_i - \tfrac{1}{2}\right)\left(\vec{w}\cdot\vec{x_i} - b\right) \geq \tfrac{1}{2} \text{ is true}$$

and return the resulting $w$ vector and $b$. There is no analytical solution. You need optimization algorithm. It can be solved with quadratic programming and other procedures as well.

note: everything we did above generalizes to $p > 2$. note: most textbooks have 1's in the place of our $\frac{1}{2}$'s that's because they assumed $y = \{-1, 1\}$ but we assumed binary.

What if the data is not linearly separable? You can never satisfy that constraint ... So this whole thing doesn't work. We will use a new objective function / loss function / error-tallying function called "hinge loss":

should be $\geq \frac{1}{2}$

$$H_i := \max\left\{0, \tfrac{1}{2} - \left(y_i - \tfrac{1}{2}\right)\left(\vec{w}\cdot\vec{x_i} - b\right)\right\}$$



Let's say a point is $d$ away from where it should be.

$$\left(z_i - \tfrac{1}{2}\right)\left(\vec{w}\cdot\vec{x_i} - b\right) = \tfrac{1}{2} - d$$

With this loss function, it is clear we wish to minimize the sum of the hinge errors:

$$H_i = \max\left\{0, \tfrac{1}{2} - \left(\tfrac{1}{2} - d\right)\right\} = \max\{0, d\} = d$$

$$SHE := \sum_{i=1}^{n} \max\left\{0, \tfrac{1}{2} - \left(y_i - \tfrac{1}{2}\right)\left(\vec{w}\cdot\vec{x_i} - b\right)\right\}$$

But we also want to maximize the margin. So we combine both consideration together into the objective function of Vapnik (1963):

$$\underset{\vec{w}, b}{\arg\min} \left\{ \frac{1}{n} \sum HE + \lambda \|\vec{w}\|^2 \right\}$$

minimizing distance errors

maximizing the width of the wedge.

Once $\lambda$ is set, the computer can do the optimization to find the resulting SVM even using out of the box R packages.

What is $\lambda$? It is a "hyperparameter", "tuning parameter". It is set by you! It controls the tradeoff between these two considerations.

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

What if you have the modeling setting where $y = \{1, 2, \ldots L\}$, a nominal categorical response with $L > 2$ levels. The model will still be a "classification model" but not a "binary classification model" and its sometimes called a "multinomial classification model". but not a "binary classification model" and its sometimes called a "multinomial classification model". What is the null model $g_0$? Again, $g_0 = $ Sample Mode $[Y]$.

Consider a model that predict on a new $x_*$ by looking through the training data and finding the "closest" $x_i$ vector and returning its $y_i$ as the predicted response value. This is called a "nearest neighbor" model. Further, you may also want to find the $K$ closest

observations and return the mode of these $K$ observations as the predicted response value (randomize ties). That's called "$K$ nearest neighbor" (KNN) model where $K$ is a natural number hyperparameter. There is another hyperparameter that must be specified the "distance function" $d: x^2 \to \mathbb{R} \geq 0$ The typical distance function is Euclidean distance squared:

$$d(\vec{x_*}, \vec{x_i}) := \sum_{j=1}^{p} (x_{ij} - x_{*ij})^2$$

what is $\mathcal{H}$? $\mathcal{A}$?