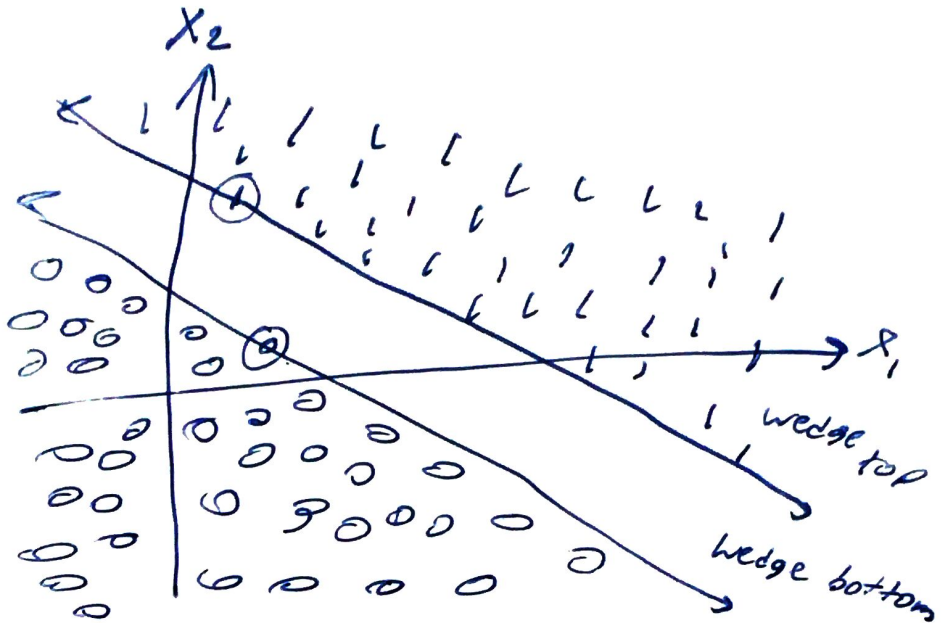# Lecture 05

$y = \{0, 1\}$, $p + 1 = 3$, $H = \{\mathbb{1}_{\vec{w} \cdot \vec{x} \geq 0} : \vec{w} \in \mathbb{R}^3\}$

Assume the data is linearly separable so
it looks like:



we need an algorithm that locates the
middle of the 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-factors, 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
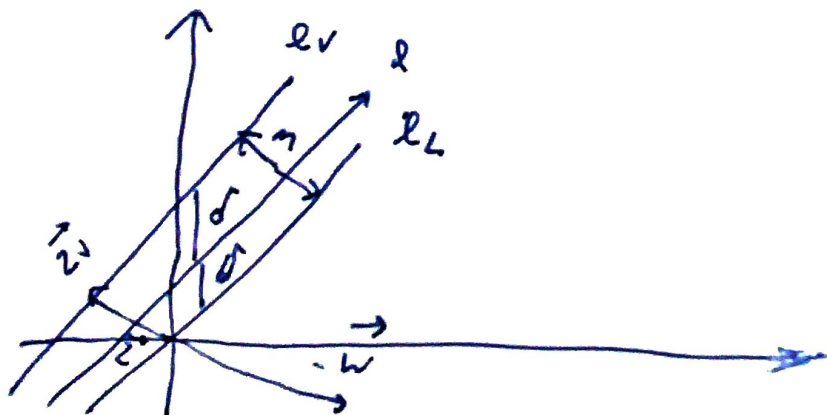means "learning complex models." To find
SVM...

First rewrite
$$H = \{ \mathbb{1}\, \vec{u} \cdot \vec{x} - b \geq 0 : \vec{w} \in \mathbb{R}^{\ell}, b \in \mathbb{R} \}$$

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

Hesse Normal
form

$$\ell = X_2 = 2x_1 + 3 \Rightarrow \ell : 2x_1 - x_2 + 3 = 0 \Rightarrow$$

$$\ell : \begin{bmatrix} 2 \\ -1 \end{bmatrix} \cdot \vec{x} - (-3) = 0$$

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

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

the direction of the $w$ vector with unit length

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.

$$\vec{Z} = \alpha \vec{w_0}, \quad \vec{Z} \in l$$

$$\boxed{\vec{w} \cdot \vec{Z} - b = 0}$$

$\Downarrow$

$$\vec{w} \cdot (\alpha \vec{w_0}) - b = 0 \Rightarrow \frac{\alpha}{\|\vec{w}\|} \|\vec{w}\| - b = 0$$

$$\Rightarrow \alpha = \frac{b}{\|\vec{w}\|} \Rightarrow \vec{Z} = \frac{b}{\|\vec{w}\|} \vec{w_0}$$

$$l_U: \vec{w} \cdot \vec{x} - (b + \delta) = 0, \quad \vec{z_U} = \frac{b + \delta}{\|\vec{w}\|} \vec{w_0}$$

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

$$m = \|\vec{z_U} - \vec{z_L}\| = \left\| \frac{b + \delta}{\|\vec{w}\|} \vec{w_0} - \frac{b - \delta}{\|\vec{w}\|} \vec{w_0} \right\|$$

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

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

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

$$\Downarrow$$

$$m \leftarrow \frac{z}{\|\vec{w}\|}$$

Now we need two conditions

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

$$\forall i \text{ s.t. } Y_i = 1$$

$$\vec{w} \cdot \vec{x}_i - (b+1) \geq 0 \implies \vec{w} \cdot \vec{x}_i - b \geq 1$$

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

$$\Downarrow$$

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

(II) All $y = 0$'s are bellow or equal to $l-L_i$:

$$\forall i \text{ s.t. } y_i = 0 \quad \vec{w}_i \vec{x}_i - (b-1) \leq 0$$

$$\implies \vec{w} \cdot \vec{x}_i - b \leq -1$$

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

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

Note how both inequalities are the same
for both I and II. Thus this inequality
satisfies both constraints, so all observations
will be in their right place

$$\forall_i \; (9_i \cdot -\tfrac{1}{2}) \; (\vec{w} \cdot \vec{x}_i - b) \ge \tfrac{1}{2}$$

⟹ line is linearly separable

You compute the sum by optimizing the
following problem:

min $\|\vec{w}\|$ s.t ⟵ true

and return the resulting w vector and b,
There is no analytical solution. You need
optimization algorithms. It can be solved with
quadratic programming and other procedures as well.
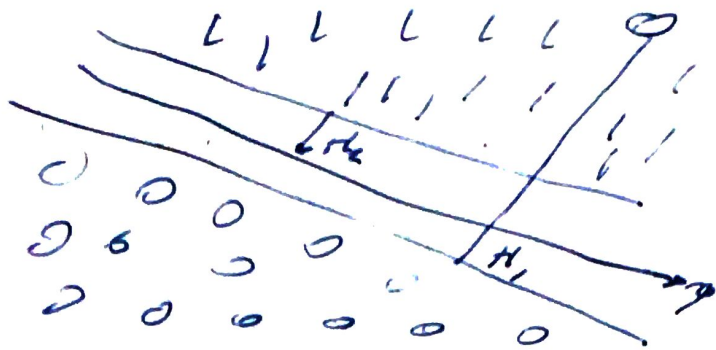
Note: everything we did above generalizes to $p \ge 2$.
Note: most textbooks have 1's in the place of
our 1/2's that's because they assumed
   $9 = \{-1, 1\}$ but we assumed binary

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

$$H_1 := \max\{0, \tfrac{1}{2} - (\varsigma_1 - \tfrac{1}{2})(\vec{u}\cdot\vec{x}_1 - b)\}$$

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

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

$$(\varsigma_1 - \tfrac{1}{2})(\vec{w}\cdot\vec{x}_1 - b_1) = \tfrac{1}{2} - d$$

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

$$SHE := \sum_{i=1}^{n} \max\{0, \tfrac{1}{2} - (\varsigma_1 - \tfrac{1}{2})(\vec{w}\cdot\vec{x}_1 - b)\}$$

But we also want to maximize the margin, so we combine both considerations together into the objective function of Vapnik (1963):

$$\underset{\vec{w}, b}{\text{argmin}}\left\{\tfrac{1}{n} SHE + \lambda \|\vec{w}\|^2\right\}$$

↑ minimizing distance errors

↑ maximize the width of the wedge

Once $\lambda$ Isset, the computer can do the optimization to find the resulting svm even using out of the box R packages

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

$$g = A(\{0, 1\}, \lambda)$$

What if you have the modelling setting where $y = \{1, 2, \ldots, L\}$, a nominal categorical response with $L \geq 2$ levels. The model will still be a "classification model" but not a "binary classification model" and it's sometimes called a "multinominal classification model".

What is the null model $g \sim 0$? Again $g \sim 0$ = SampleMode[$y$].

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

$$d(\vec{x}_x, \vec{x}_i) := \sum_{j=1}^{p} (x_{i,j} - x_{v,j})^2$$

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