* Error Relaxation

# Machine learning

A What is ML?

→ A set of methods for making decisions from data.
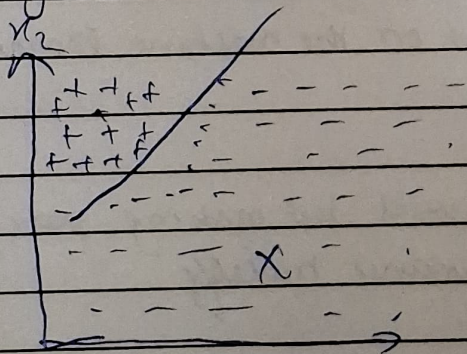
● Why study ML?

→ To apply; to understand; to evaluate.

& We have training data.

For data point $i \in \{1, \dots n\}$
- Feature vector
$$x^{(i)} = (x_1^{(i)}, \dots x_d^{(i)})^T \in \mathbb{R}^d$$

- label $y^i \in \{-1, +1\}$



Training Data $D_r = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

* What we want? A good way to ~~the~~ label new pts.
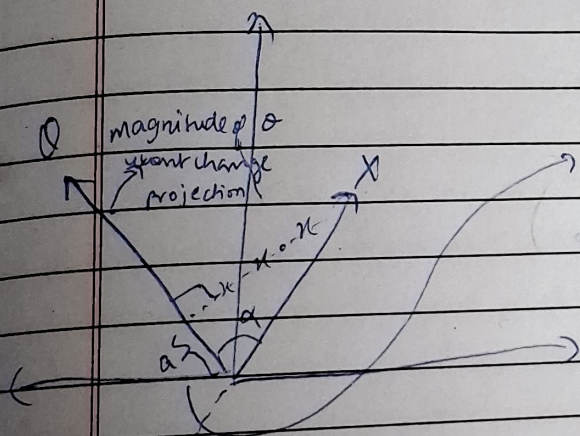
* Let $h: R^d \to \{-1, +1\}$

    $x \circlearrowright y - x \to h \to y$

* linear Classifiers

Example H: All hypotheses that label +1 on one side of a
line -1 on the other side.



Q magnitude of $\theta$     $\to$    $\theta^T x / \|\theta\|$ (Projection of $x$ onto $\theta$)
   won't change    $X$         $[1 \times d] \times [d \times 1]$ ; dimensions.
   projection

if $\alpha = 90°$, $|x|$ doesn't matter
projection will be 0.

$x: \theta^T x / \|\theta\| = a$ will be the dotted line.

let $x : \Theta^T x / \|\Theta\|$ be $f(x, \Theta)$

$f(x, \Theta) > b - b$

$: \Theta^T x / \|\Theta\| = -b = x : \Theta^T x + b\|\Theta\| = 0$

$x : \Theta^T x + \Theta_0 = 0$

$f(x, \Theta) < -b$

$\Theta_0 = b\|\Theta\|$

- Linear classifier

$$h(x) = \text{sign}(\Theta^T x + \Theta_0)$$

$$= \begin{cases} +1 & \text{if } \Theta^T x + \Theta_0 \geq 0 \\ -1 & \text{if } \Theta^T x + \Theta_0 \leq 0 \end{cases}$$

$$h(x; \underset{\uparrow \nearrow}{\Theta, \Theta_0})$$

parameters (not input)

$H = $ set of all $h$.

The double bars on $\|\theta\|$ represent a norm (vector)
It calculates the magnitude of that vector in multi
dimensional space

✻ **How good is a classifier**

- It should predict well on future data
- How good is a classifier at a single point? Loss $L(g, a)$
  g: guess ; a: actual.

- Our guess should be closer to the actual value.

  Example: 0-1 Loss.

  $$L(g, a) = \begin{cases} 0 & \text{if } g = a \\ 1 & \text{else.} \end{cases}$$

  : Asymmetric loss.

  $$L(g, a) = \begin{cases} 1 & \text{if } g = 1, \; a = -1 \quad (\text{Newborn didn't have seizure, diagnosed else}) \\ 100 & \text{if } g = -1, \; a = 1 \quad (\text{Newborn had seizure, diagnosed else}) \\ 0 & \text{else.} \quad \quad \quad \quad \quad \quad (\text{correct diagnosis}) \end{cases}$$

- Test error ($n'$ new points) $\mathcal{E}(h) = \dfrac{1}{n'} \sum\limits_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$

- Training error $\mathcal{E}_n(h) = \dfrac{1}{n} \sum\limits_{n=1}^{n} L(h(x^{(i)}), y^{(i)})$

- Prefer $h$ to $\tilde{h}$ if $\mathcal{E}_n(h) < \mathcal{E}_n(\tilde{h})$

## learning a classifier.

Recall: $x \longrightarrow \boxed{h} \longrightarrow y$

New:

$$D_n \longrightarrow \boxed{\begin{array}{c}\text{learning}\\\text{algorithm}\end{array}} \longrightarrow h$$

Ex

for $j = 1, \ldots$ 1 million

Randomly sample $\{\theta^{(j)}, \theta_0^{(j)}\}$

Set $h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$

Ex — learning_alg $(D_n; \underset{\nwarrow \text{hyperparameter}}{k \leq k \text{ trillion}})$

Set $j^* = \underset{j \in \{1, \ldots k\}}{\text{argmin}} \, E_n(h^{(j)})$

Return $h^{(j^*)}$