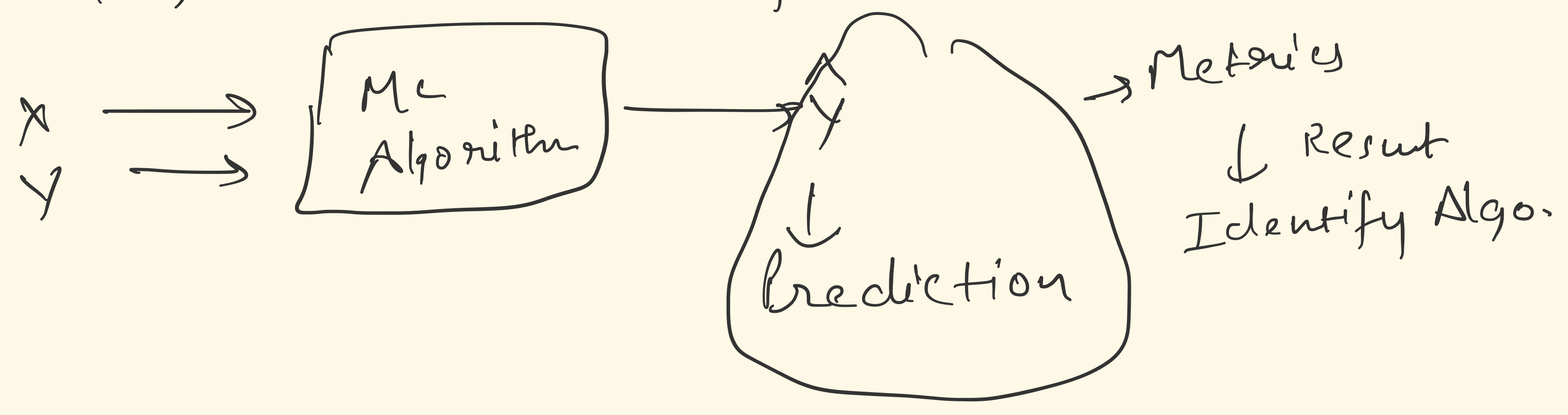
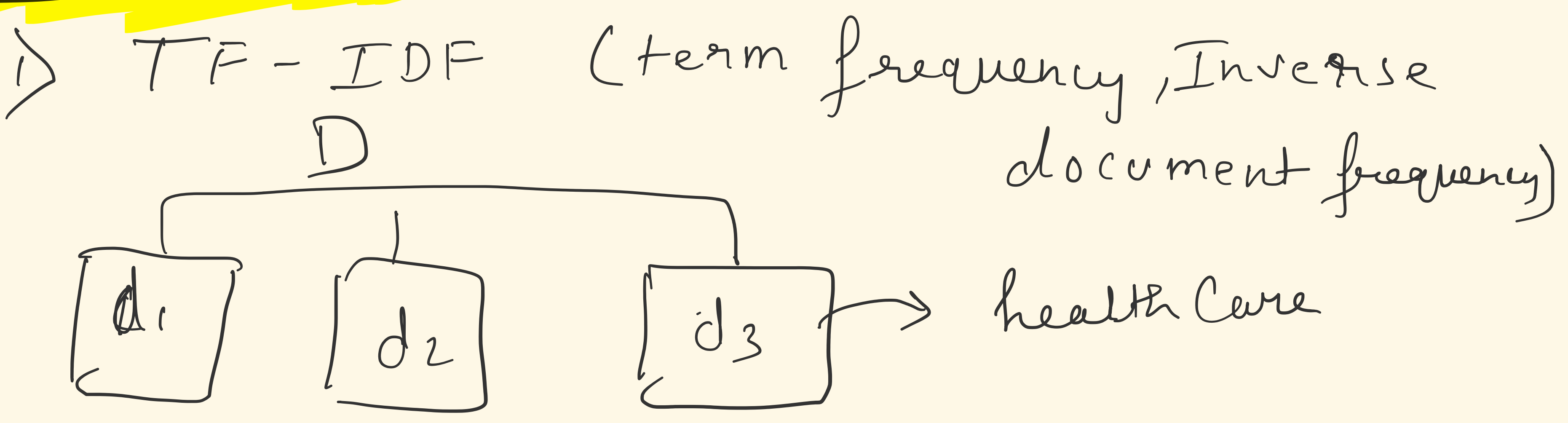


### M.L training

- \* Everything should be numeric
- \* Data is always non-linear in real world
- \* find features  $\rightarrow \mathbf{X} \in \mathbb{R}^D$   
 $\hookrightarrow [10 \times 5]$
- \* find labels  $y \in \mathbb{R} \setminus \{0, 1\}$
- \* Problem  $\in [\text{Regression, Classification}]$
- \* Apply multiple Algo and find the best



### \* Text-preprocessing



### Term frequency

$tf(t, d_3) \Rightarrow \frac{5}{25} = \left(\frac{1}{5}\right) = .2$

↓  
term  
↓  
word.

$\Rightarrow \frac{20}{25}$

$\left\{ \frac{\text{total time word occur in document } d}{\text{total number of word in document } d} \right\}$

$\left[ \begin{matrix} is, the, a \\ \uparrow \quad \uparrow \quad \uparrow \end{matrix} \right]$

### Inverse document frequency

$idf(t, D) \Rightarrow \left( \frac{N}{\text{number of time term occur in docs}} \right) \log$

$\Rightarrow \frac{3}{3} \rightarrow \log\left(\frac{3}{3}\right) \rightarrow 0$