Regression Models belong to <u>Supervised Learning Class</u>
of Algorithms

*What is regr.*

numerical input → model $\beta$ → numerical

\* Assume the relationship
input ~ output is linear

$(X^{(i)}, y^{(i)})$ $\quad$ $X^{(i)} \in \mathbb{R}^n$ $\quad$ $y^{(i)} \in \mathbb{R}$

$y^{(i)} \approx \beta \cdot X^{(i)}$ $\qquad$ model parameter

$y^{(i)} \approx \beta \cdot X^{(i)} + \alpha$
$\qquad \uparrow$ intercept

What do we learn from the Data?
We determine $\beta$!

*How $\beta$ is determined*

$$TSE_\beta = \sum_{i=1}^{N} | y^{(i)} - \beta \cdot X^{(i)} |^2$$

we choose $\beta$ s.t. $\qquad$ $\boxed{TSE_\beta \text{ is minimized}}$

$$\beta_0 = \arg\min_\beta TSE_\beta$$ $\qquad$ $TSE : \mathbb{R}^n \to \mathbb{R}$
$\qquad \qquad \qquad \qquad \qquad \qquad \downarrow$
$\qquad \qquad \qquad \qquad \qquad \qquad \beta$

$TSE$ is a function in $\beta$
if $D = \{ (X^{(i)}, y^{(i)}) \in \mathbb{R}^n \times \mathbb{R} \}$ $\beta$ fixed

There are two main Reasons
why we do ML

① Make predictions!

$(i) \to \boxed{M} \to (P)$

② To understand the data
particularly to understand
the functional relationship
between input and output!

Data $\longrightarrow$ we decide on the model type $\longrightarrow$ we determine the best model param.

structured Num. Data $\rightsquigarrow$ Regression $\rightsquigarrow$ the model param. $\boxed{\beta}$

Today the $\beta$ (components of $\beta$) will tell us the functional linear relationship between the features and the output
predictors        regressor

Require an analysis of the model and the relationship between predictors & regressor

$$\beta = \langle 1, 2, 3 \rangle$$

$f_1 \quad f_2 \quad f_3$

$$f_1 < f_2 < f_3$$

| $f_1$ | $f_2$ | $f_3$ | $y$ |
|---|---|---|---|
|  |  |  |  |

Doesn't mean $f_3$ more important.

One has to look at the distribution of each of these columns.

For example if $\text{Data} \xrightarrow{\pi_3} \mathbb{R}$    $\leftarrow$ proj on the third feature

$[0, 0.1]_{\pi_2}$

as opposed to $\text{Data} \xrightarrow{\pi_2} \mathbb{R}$

$[0, 10000]$

## Categorical Variable Encodings

Categorical $\rightarrow$ A   B $\leftarrow$ numerical

| A | B |
|---|---|
| $a_1$ | $b_1$ |
| $a_2$ | $b_2$ |
| $\vdots$ | $\vdots$ |
| $a_n$ | $b_n$ |

Mean

I want to encode A using values in B

① Split the data set B using the labels in a.

② Bag$_a$   $\forall a \in A$   then

represent each $a \in A$ using mean($Bag_a$)

(Freq) ① is the same
② instead of mean use the # of elements

| A | B |
|---|---|
| a | 10 |
| b | 11 |
| a | 12 |
| b | 13 |

(mean)

$A = \{a, b\}$

$Bag_a = \{10, 12\} \longrightarrow a \longleftrightarrow 11$

$Bag_b = \{11, 13\} \longrightarrow b \longleftrightarrow 12$

(Freq) No need for B

A | a | b | a | b |

A | a | a | a | b |

$\begin{matrix} a \longleftrightarrow 2 \\ b \longleftrightarrow 2 \end{matrix}$ ← represented by the same number.

$\begin{matrix} a \longleftrightarrow 3 \\ b \longleftrightarrow 1 \end{matrix}$

---

Regression with Categorical Variables

Regression must use numerical vectors as input

(numerical input) → | model $_\beta$ | → (numerical output)

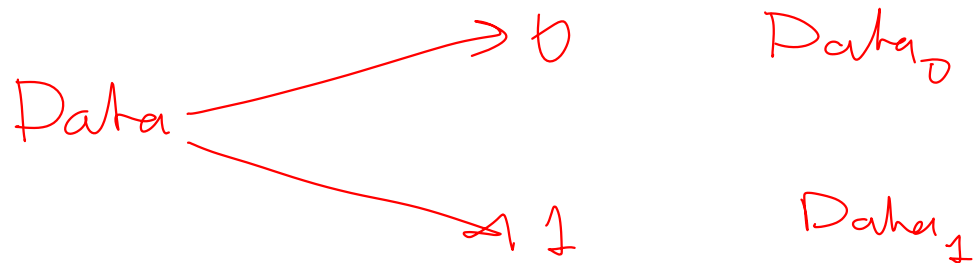| | $X_1$ | $X_2$ | | $X_n$ |
|---|---|---|---|---|
| $d_1$ | | | | |
| $d_2$ | | ∿ | | |
| $d_N$ | | | | |

If any of $X_i$ is categorical we must convert/encode it as a numerical variable

# Logistic Regression

Input $\in \mathbb{R}^n$ x | Categorical |

Output $\in \{0, 1\}$

Binary.

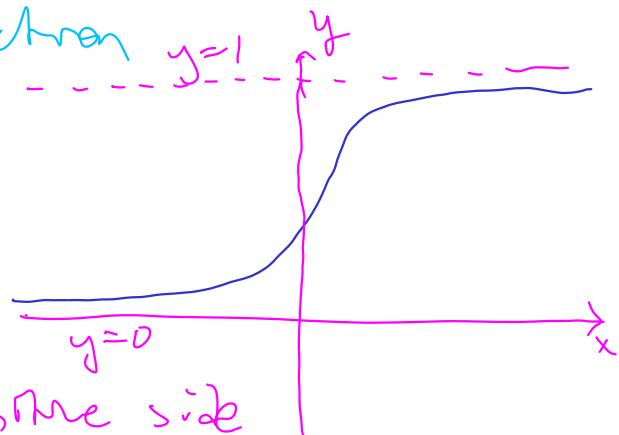In other words I want to split my data set into __two__ disjoint subsets

Data $\longrightarrow$ 0    Data$_0$

Data $\longrightarrow$ 1    Data$_1$

| Idea | Create a linear regression model

so that
$$\text{model}_\beta(x) > 0 \longrightarrow \text{Data}_1 \quad \boxed{\text{model}_\beta}$$
$$\text{model}_\beta(x) < 0 \longrightarrow \text{Data}_0$$

Decision on whether $x \in \text{Data}_1$ or $x \in \text{Data}_0$
is done via the sign of model$_\beta$ output.

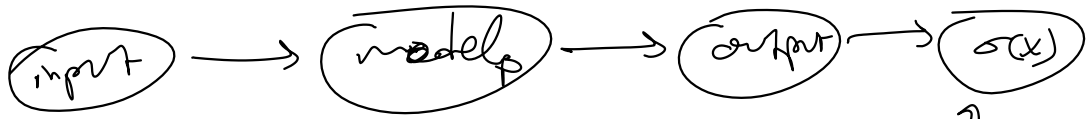sigmoid, logistic function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

why this function
① it rapidly increases to 1
as soon as x passes to the positive side
② it rapidly decreases to 0
as soon as x passes to negative region

model $\qquad \beta \cdot X + \alpha$ $\qquad$ numerical values in $\mathbb{R}$

$\rightsquigarrow \qquad \sigma(\beta \cdot X + \alpha)$

$$\boxed{\text{input}} \longrightarrow \boxed{\text{model}_\beta} \longrightarrow \boxed{\text{output}} \longrightarrow \boxed{\sigma(x)}$$

this number here is a number btw 0 and 1

## The interpretation

Recall ① if the output is close to 0 the input belongs to $Data_0$

② if the output is close to 1 the input belongs to $Data_1$

read $\sigma(\beta \cdot X^{(i)} + \alpha)$ as the $\boxed{\text{probability}}$

that $X^{(i)}$ belongs to $Data_1$!