# LOGISTIC REGRESSION

→ concept matters ; tool/package doesn't matter

### LINEAR DATA



$y = \beta_0 + \beta_1 x$

### NON-LINEAR DATA



$y = \beta_0 + \beta_1 x$

Age vs Buy

Linear Regression fails here

---

Sales — Buying vs Not Buying

Marketing — Response vs No Response

Credit card & Loans — Default vs Non Default
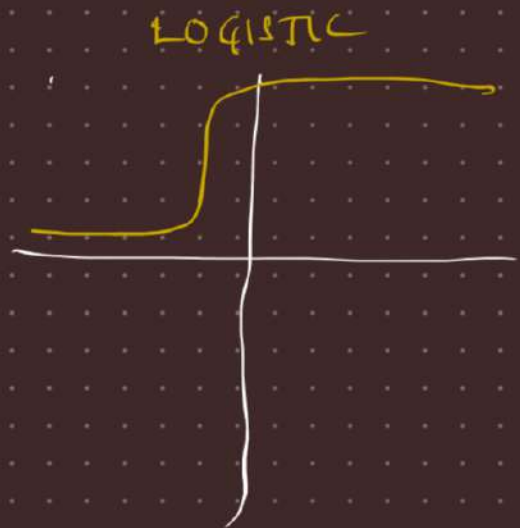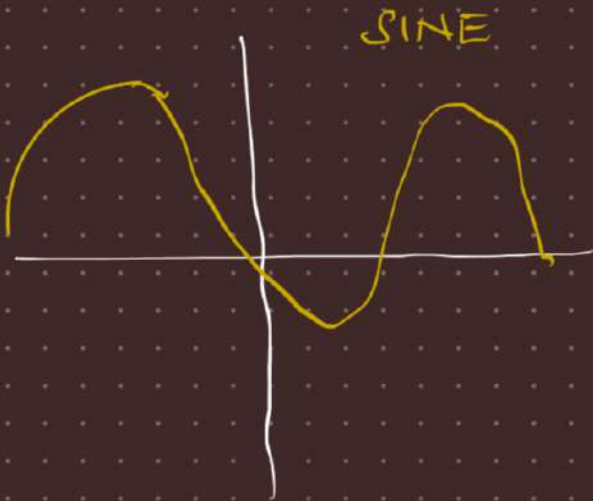
Website — Click vs No click
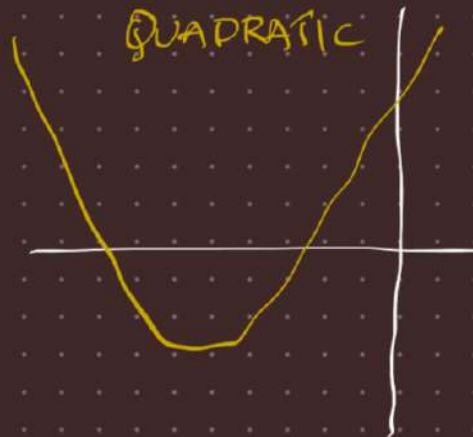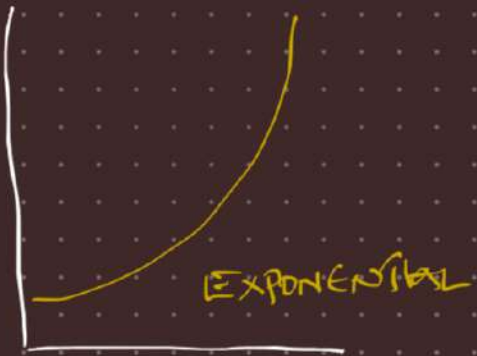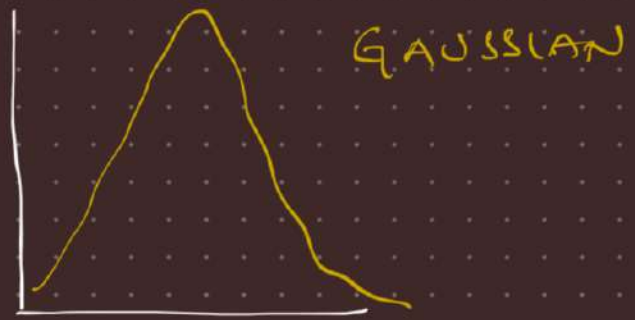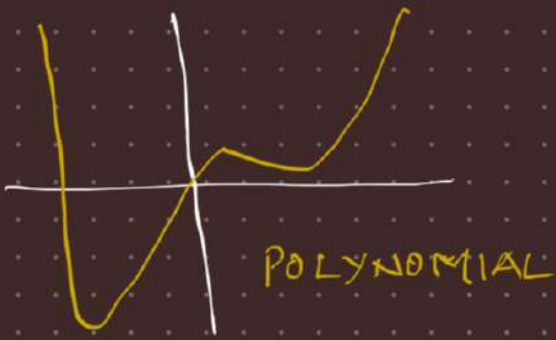
→ CLASSIFICATION PROBLEM

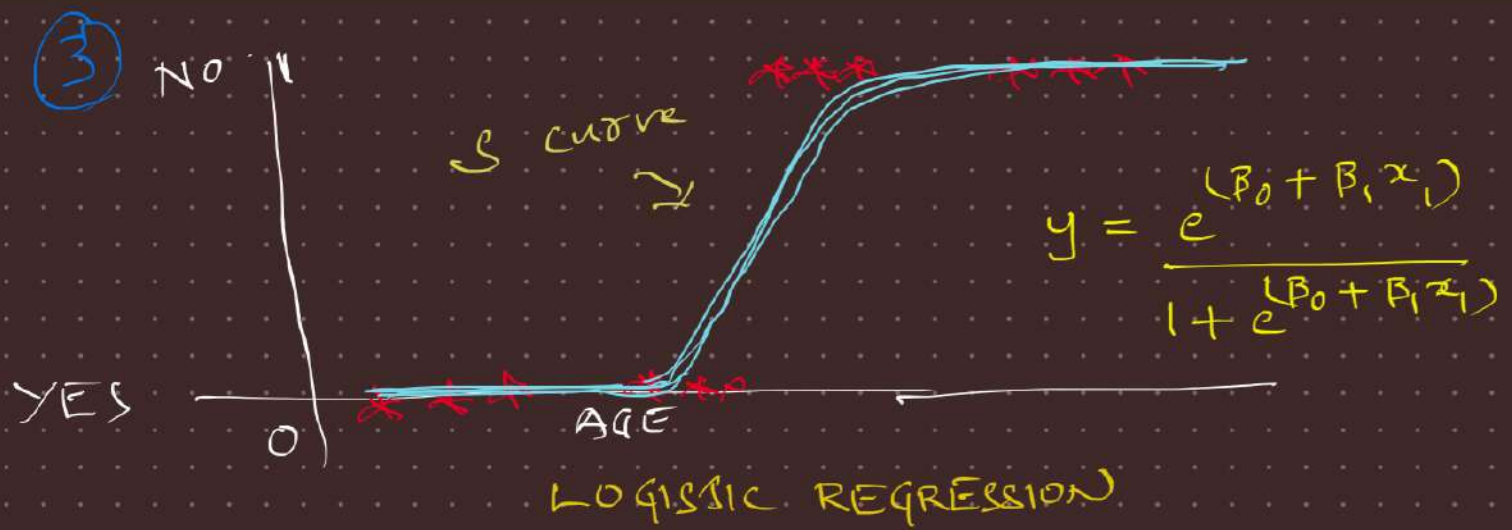LIN. REG — a good fit if you are predicting a continuous value (target variable)

If you are predicting class 0 or class 1 ?

↳ Non linear function

(2)

# Some Non linear functions

POLYNOMIAL

GAUSSIAN

EXPONENTIAL

QUADRATIC

SINE

LOGISTIC

③ NO

S curve

$$y = \dfrac{e^{(\beta_0 + \beta_1 x_1)}}{1 + e^{(\beta_0 + \beta_1 x_1)}}$$

YES

0   AGE

LOGISTIC REGRESSION

Note :

Looking at problem statement
at target variable   } should be able to decide

Regression   Classification
problem      problem

Exercise : Problem statements

① Predicting Loss %          →   Regression

② Predicting Buying Vs Not buying → classification

③ Predicting no. of customers → regression

④ Predicting response vs no response → classification

⑤ Predicting Revenue → regression

⑥ Predicting the product price → regression

⑦ Predicting Attrition vs. Retention → classification

⑧ Predicting click Vs No click → classification

⑨ Predicting Fraud vs Non Fraud → classification

⑩ Predicting the amount of fraud → Regression

# MULTIPLE LOGISTIC REGRESSION

real world $\rightarrow$ target variable depends on multiple features

$$y = \frac{e^{\beta_0 + B_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \ldots - \beta_k x c_k}}{1 + e(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + - \ldots \beta_k x_k)}$$

---

Note :

customer
churned out $\}$ $\left. \begin{array}{l} \text{customer} \\ \text{leaving} \end{array} \right.$

Model building must be always followed by Model validation.

Earlier, for Linear Regression, model validation $\rightarrow R^2$

For Logistic Regression (classification problem)

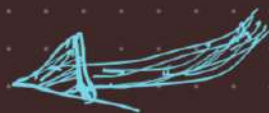Model validation :  $\longrightarrow$ next page

# LOGISTIC REGRESSION: MODEL VALIDATION

→ model.predict( )

| $x_1, x_2, x_3$ | $\cdots\cdots$ | $x_K$ | $y_{actual}$ | $y_{pred}$ |
|---|---|---|---|---|
| — — | | — | 1 | 1 |
| — — | | — | 0 | 1 |
| — — | | — | 1 | 1 |
| — — | | — | 0 | 0 |
| — — | | — | 1 | 0 |
| — — | | — | 0 | 0 |



Confusion matrix

|  | | $y_{pred}$ | |
|---|---|---|---|
| | | 0 | 1 |
| $y_{actual}$ | 0 | a | b |
| | 1 | c | d |

$$\text{Accuracy} = \frac{a+d}{(a+b+c+d)}$$

What is the guarentee that this "Accuracy" remains the same for new data.

To solve this problem, we need to do "Cross Validation"

Lets say, dataset = 10,000 records

80,000                    20,000
(training data)           (testing data)

Train data Accuracy = 92%  } this model
Test data Accuracy = 60%   } is suffering
                             from
                             **OVERFITTING**

Train Accuracy = 60%  } this model is
                      } suffering from
                        **UNDERFITTING**

Overfitting : Train accuracy is very high
                but test accuracy is
                significantly low

Thats the reason, validation alone is
not sufficient, you need to do
cross validation.

How do you do that ?

Divide the data into train data & test data

In sklearn, → model_selection.train_test_split