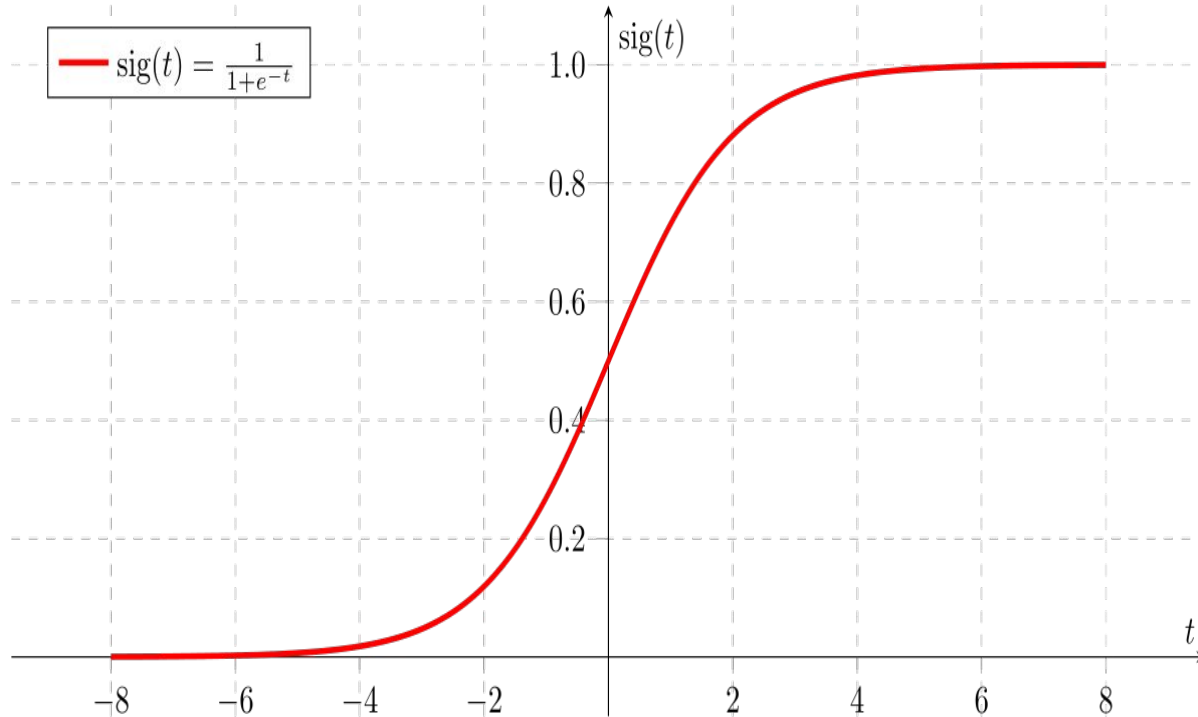# Logical Rhythm

Class #2: 27th August 2019

# Logistic Regression (Classification Algorithm)



$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

**sig(z) = 1/(1+e⁻ᶻ)**

**Y = sig (∑wᵢxᵢ)**

"*Probability of output to be categorically 1 for given value of x.*"

**Reasons for not using Linear Regression** -

Value not in finite range.

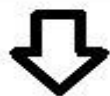Works with luck for best fit.

**Logistic regression cost function**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[ \sum_{i=1}^{m} \boxed{y^{(i)} \log h_\theta(x^{(i)})} + \boxed{(1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))} \right]$$
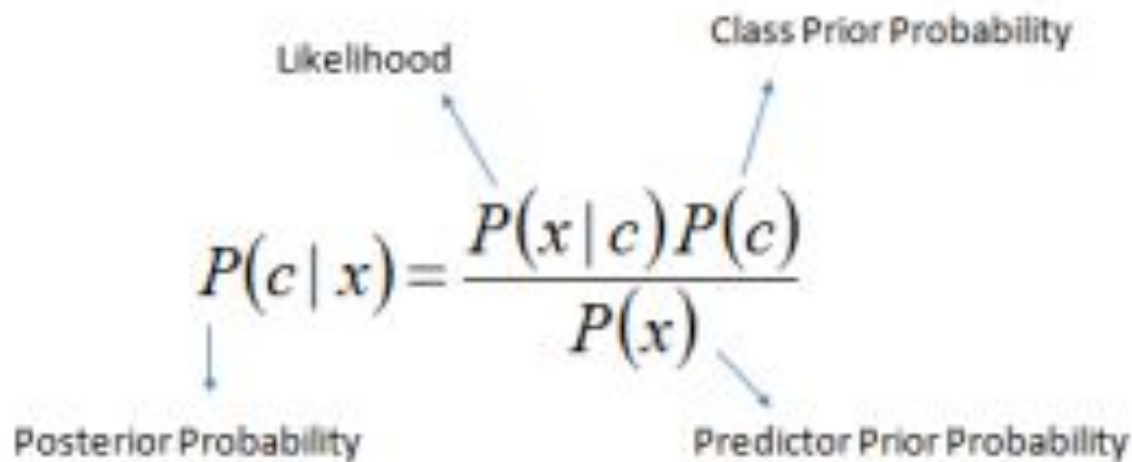
0 ⬆ If actual y=1

⬇ 0   If actual y=0

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

Note: $y = 0$ or $1$ always

# Bayes' Theorem

$$P(c \mid x) = \frac{P(x \mid c)P(c)}{P(x)}$$

Likelihood

Class Prior Probability

Posterior Probability

Predictor Prior Probability

# Derivation

- P(A and B) = P(A)*P(B|A)
- P(A and B) = P(B)*P(A|B)
- P(A)*P(B|A) = P(B)*P(A|B)
- **P(B|A) = P(B)*P(A|B)/P(A)**

# Multivariate Naive Bayes

$$P(y|X) = \frac{P(X|y)P(y)}{P(X)}$$

$$X = (x_1, x_2, x_3, ....., x_n)$$

$$P(y|x_1, ..., x_n) = \frac{P(x_1|y)P(x_2|y)...P(x_n|y)P(y)}{P(x_1)P(x_2)...P(x_n)}$$

$$P(y|x_1, ..., x_n) \propto P(y) \prod_{i=1}^{n} P(x_i|y)$$

# Additional Information

- **Pros and Cons**
  - **Pros**
    - Fast and accurate when independent features are involved
    - Works better for categorical valued features, as numeric features involved normal distribution assumption
  - **Cons**
    - Features are rarely independent in real-life problems
    - **Zero frequency** problem

- **Applications:** Recommender System, Text Classification, Sentiment Analysis