

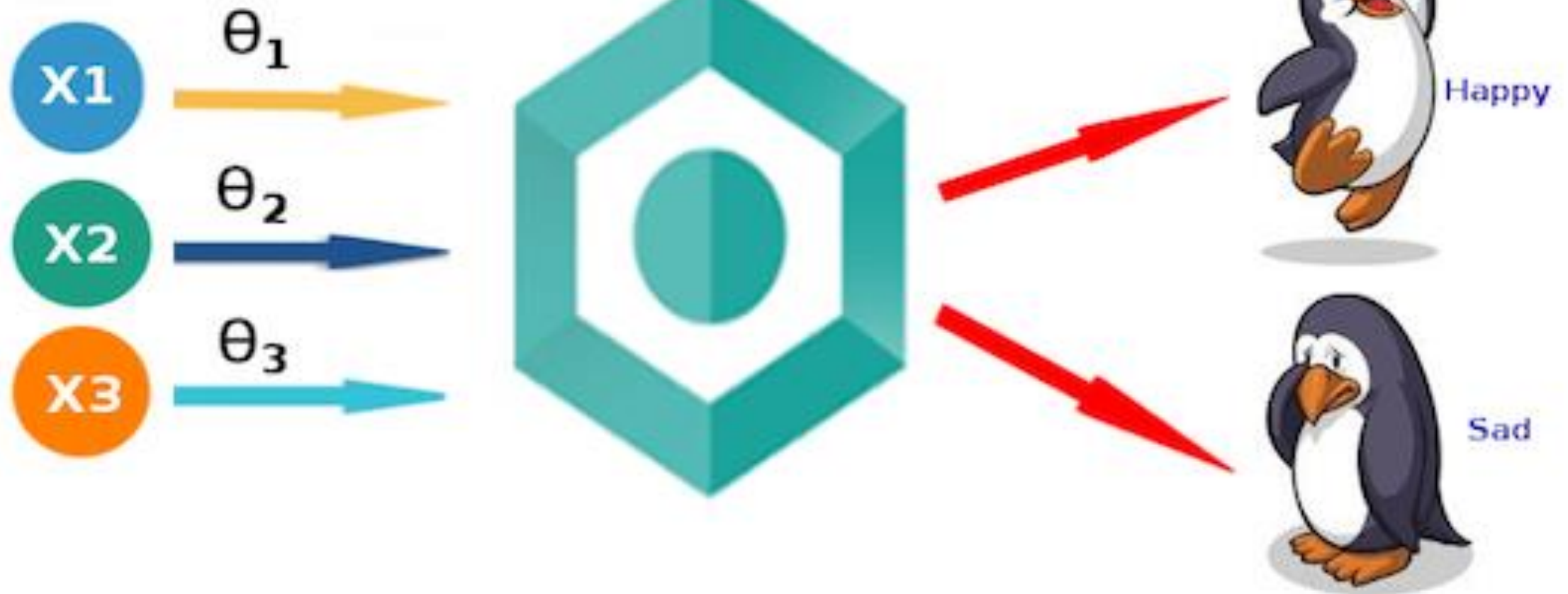
# Practical Machine Learning

## Day 8: Mar23 DBDA

Kiran Waghmare

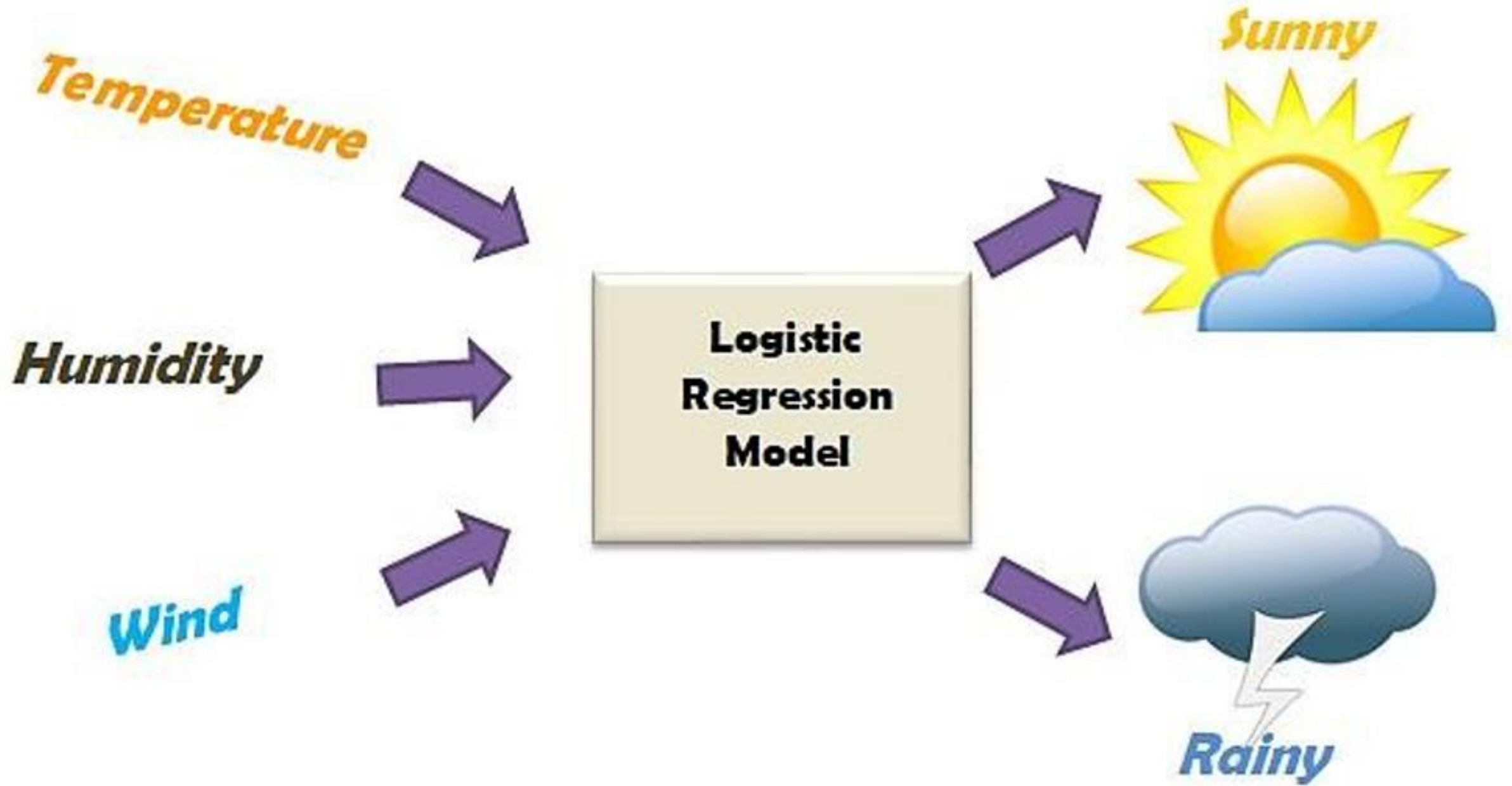
# Agenda

- Logistic Regression
- Classification
- Measures for classification
- KNN

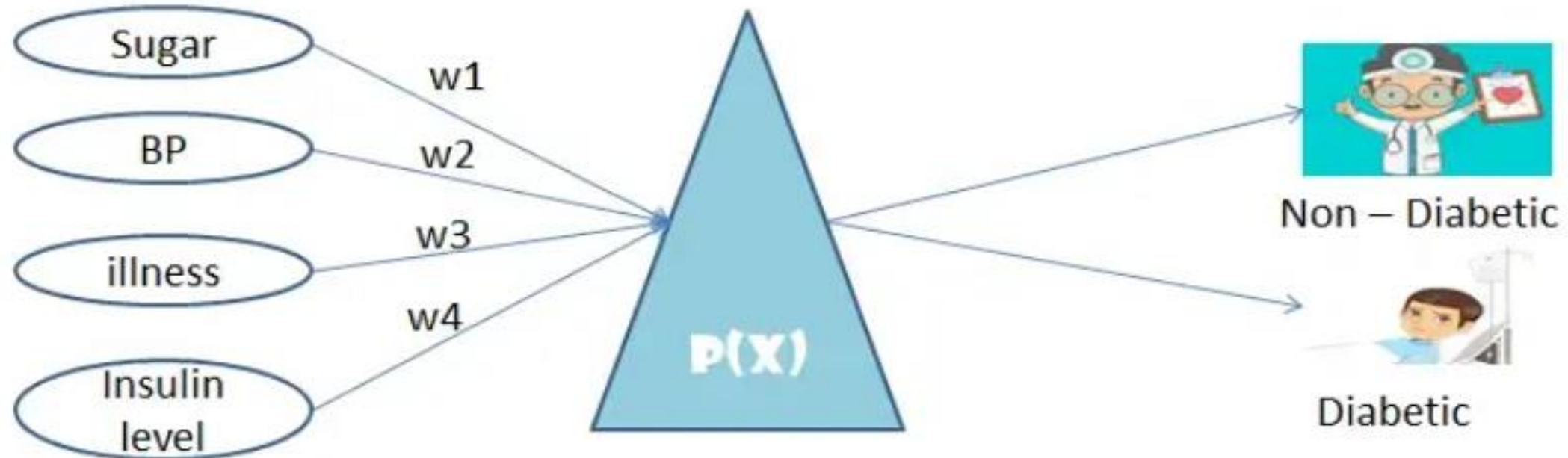


@dataaspirant.com

Inputs:  $X_1$ ,  $X_2$ ,  $X_3$  || Weights:  $\theta_1$ ,  $\theta_2$ ,  $\theta_3$  || Outputs: Happy or Sad



# LOGISTIC REGRESSION MODELLING

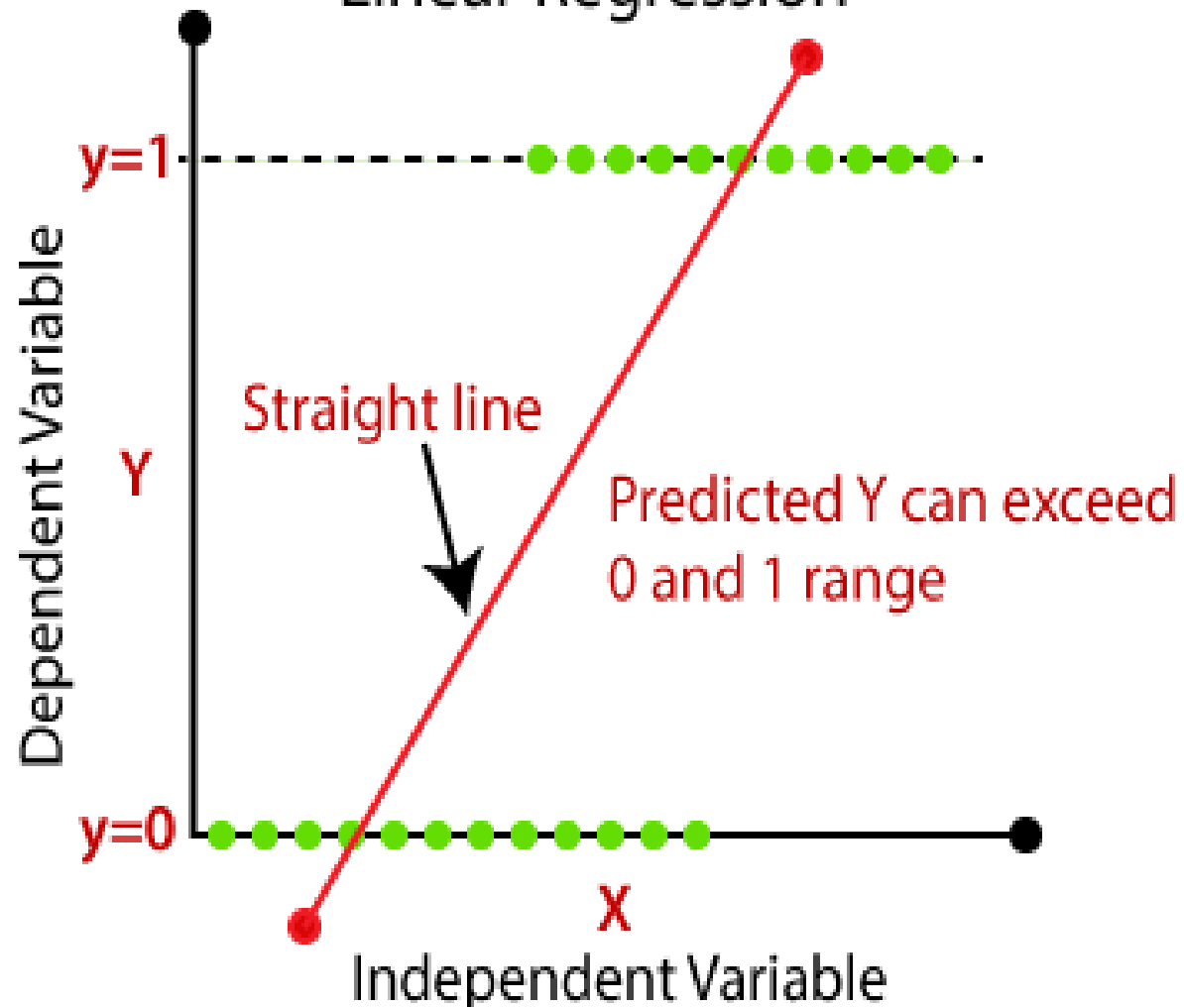


$w1, w2, w3, w4$  – Amount of each individual medical problem

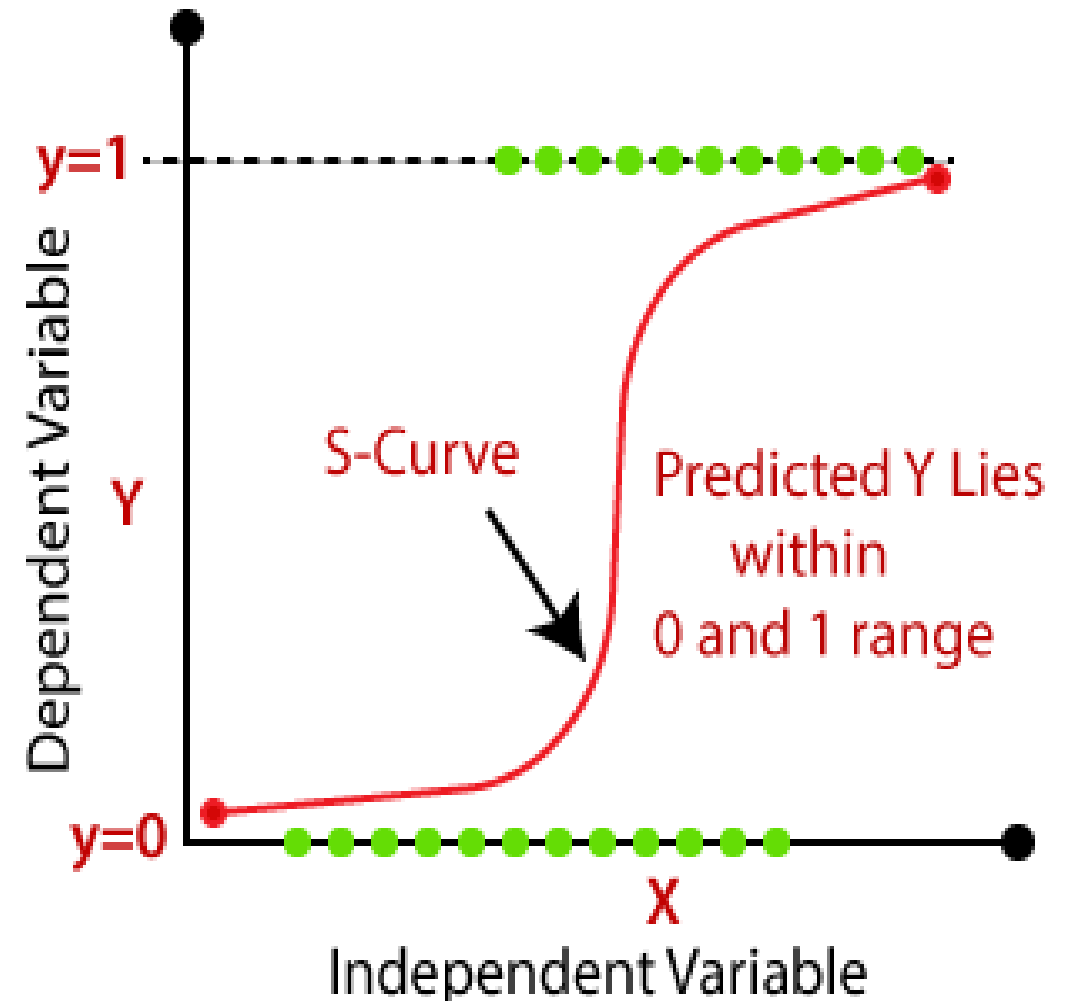
$P(x)$  – Probability Calculation

Logistic Regression

## Linear Regression

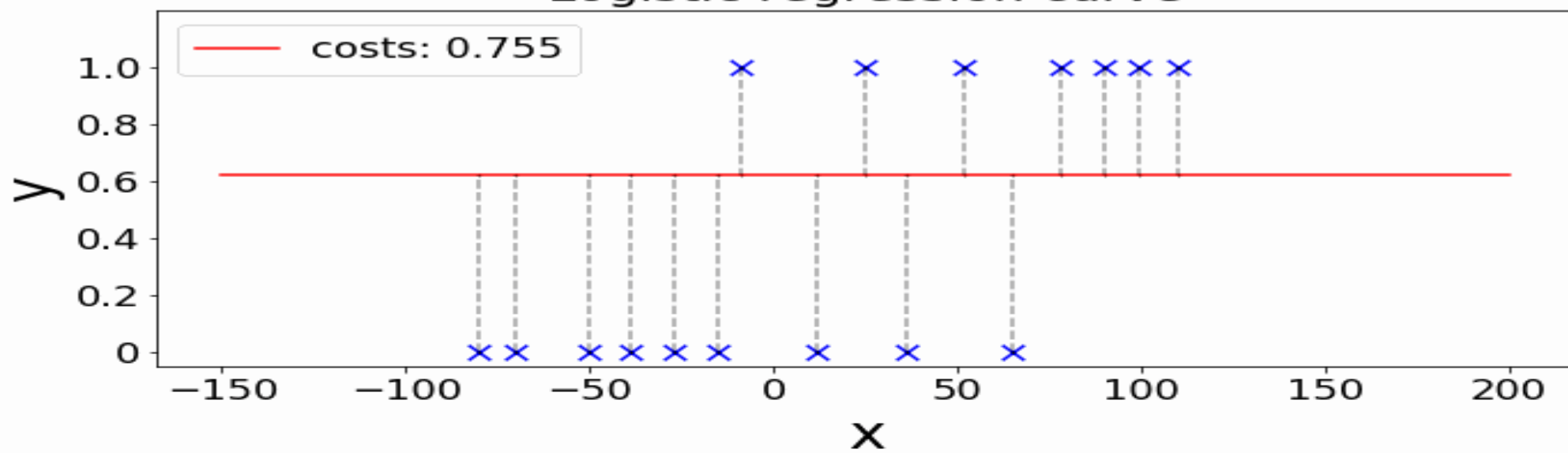


## Logistic Regression

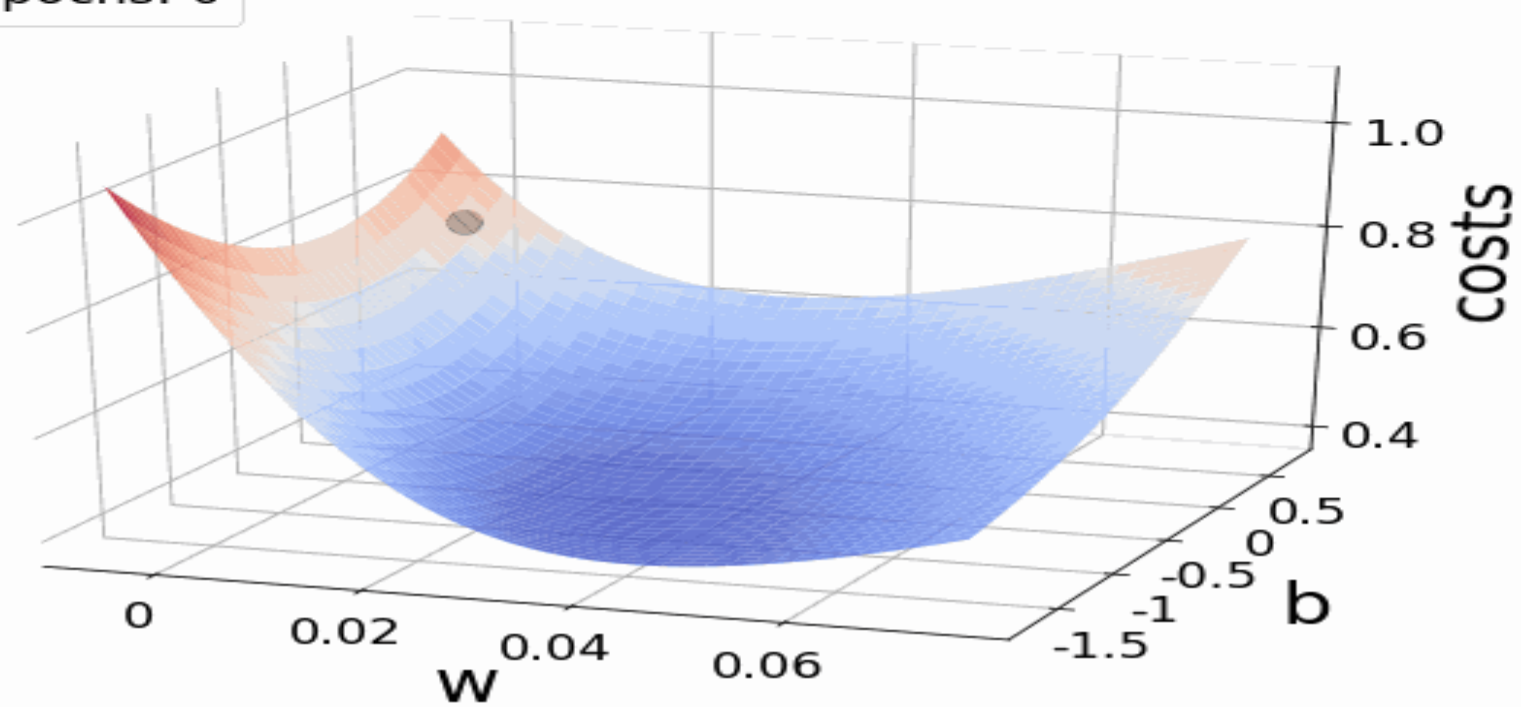


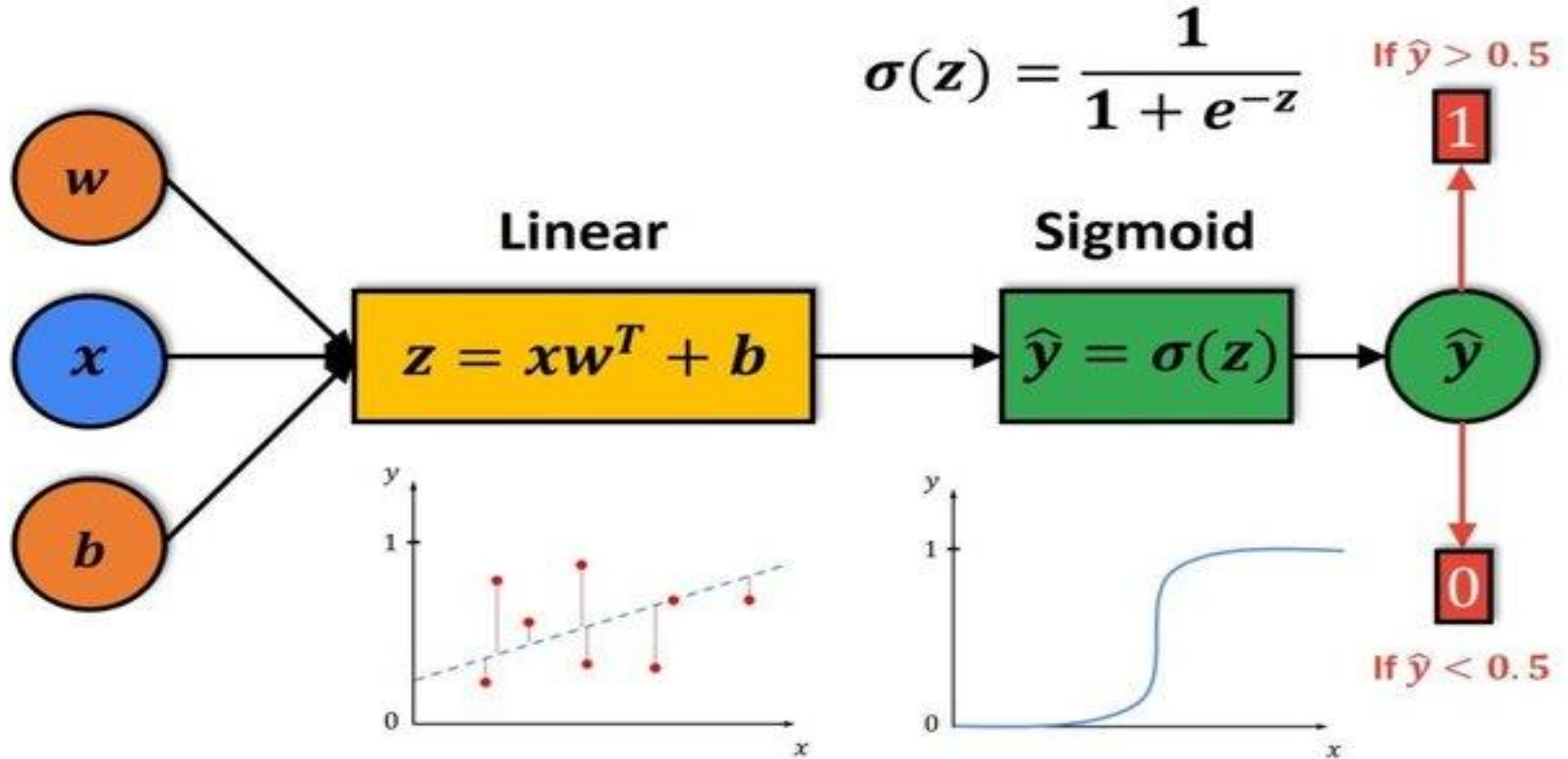


Logistic regression curve



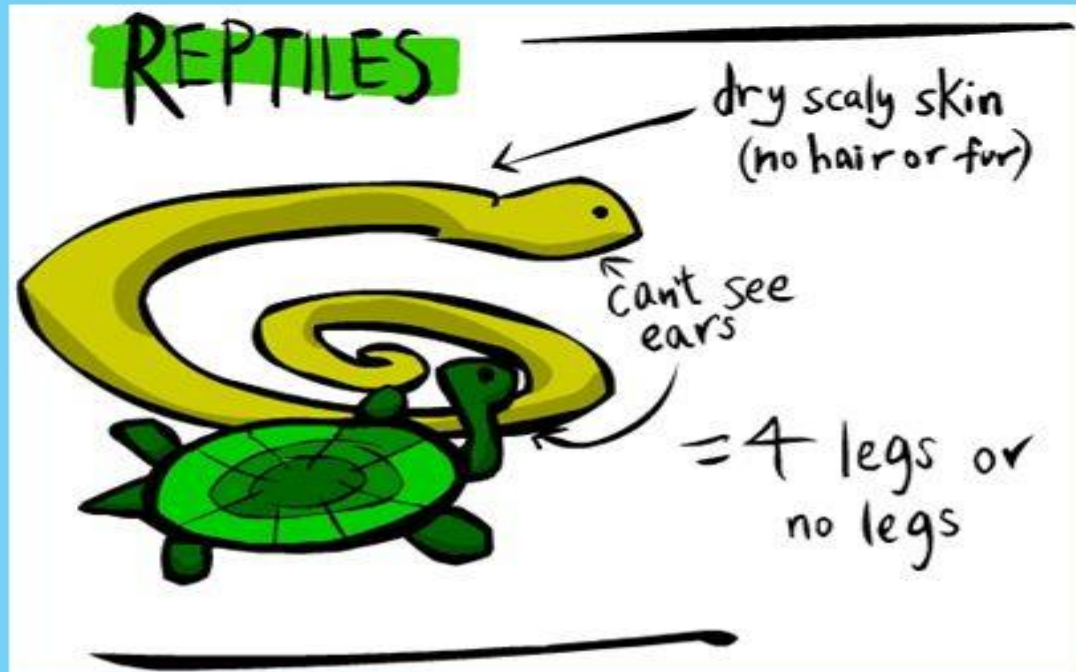
epochs: 0



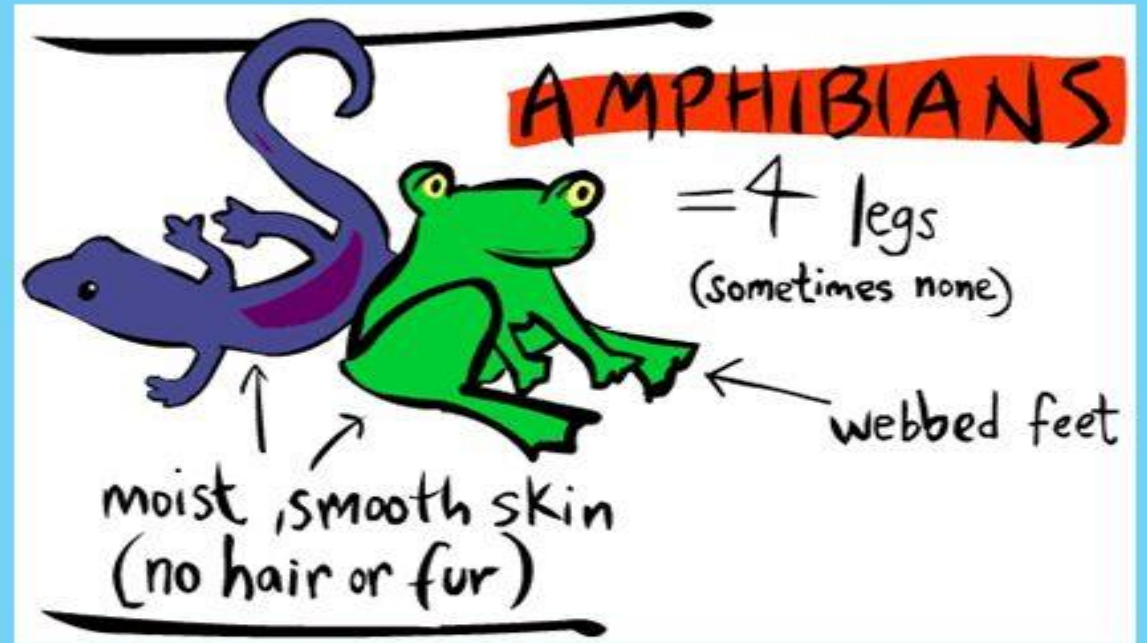




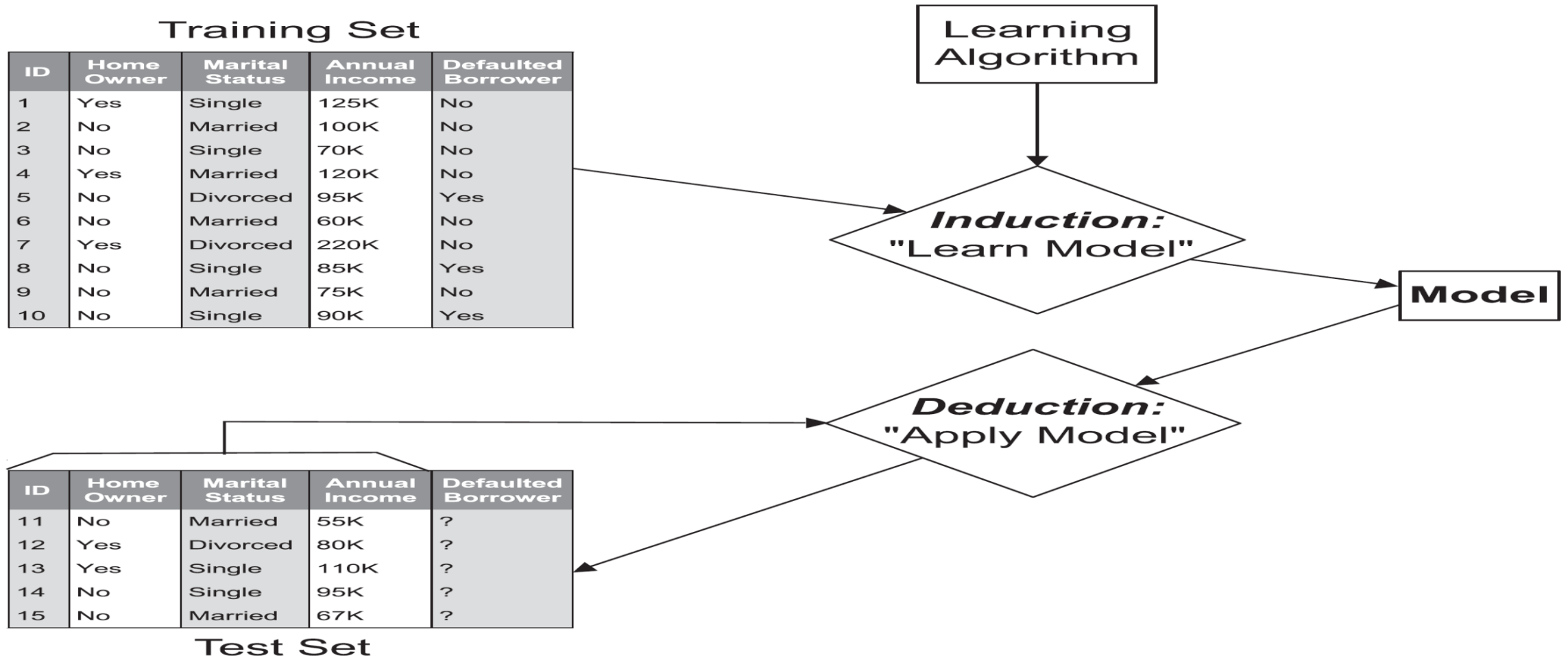
# Reptiles



# Amphibians



# General Approach for Building Classification Model



**Figure 3.3.** General framework for building a classification model.

# Face Recognition

Training examples of a person

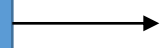
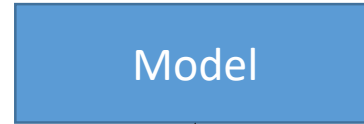
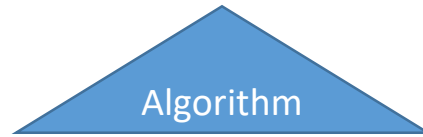


Test images



AT&T Laboratories, Cambridge UK  
<http://www.uk.research.att.com/facedatabase.html>

# Classification



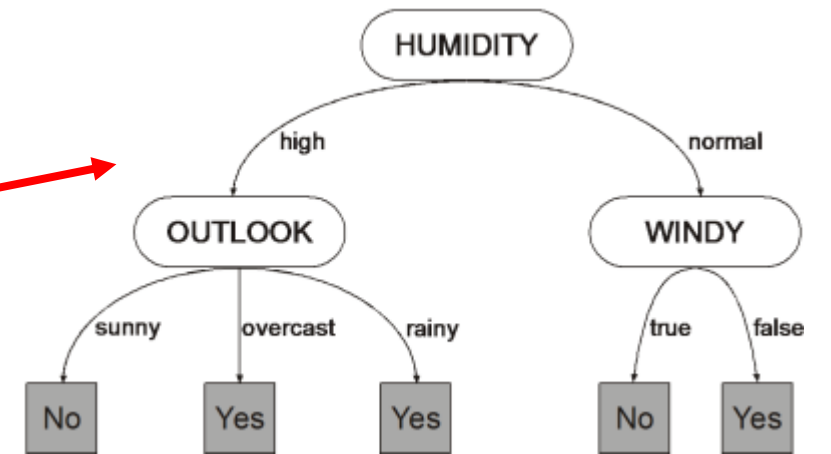
lion





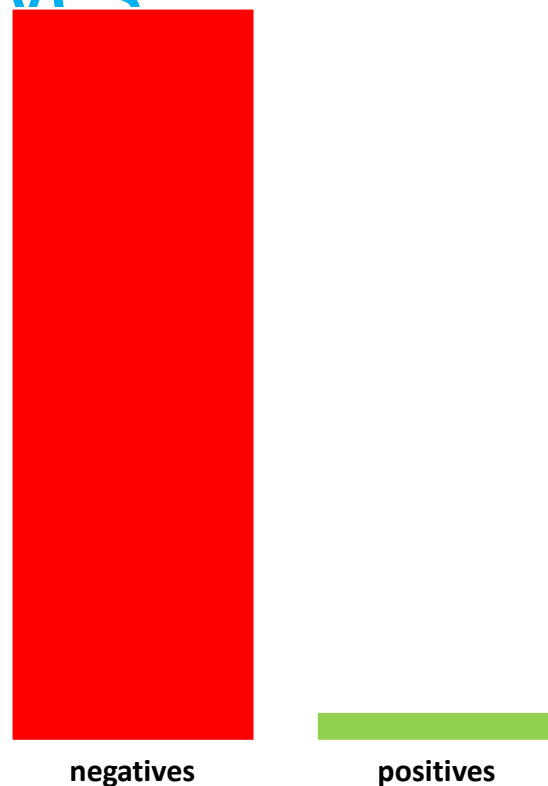
# The resulting *model* is also called the *hypothesis*.

Given a model space and an optimality criterion, a *model* satisfying this criterion is sought.



Optimal tree!

After plotting your class distribution you see that you have **thousands of negative examples but just a couple of positives**





# Performance metrics

- Most of the time accuracy will not be enough to assess performance.

- $accuracy = \frac{TP+TN}{P+N}$

Percentage of correctly classified instances.

- $sensitivity = \frac{TP}{P}$

The proportion of positives that are correctly identified as such.

- $precision = \frac{TP}{TP+FP}$

Equivalently, it is the fraction of relevant instances among the selected ones.

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Matthews correlation coefficient (takes into account imbalance)

# Confusion Matrix

- Confusion Matrix:

ACTUAL CLASS	PREDICTED CLASS	
	Class=Yes	Class=No
Class=Yes	a	b
	c	d

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

# Accuracy

ACTUAL CLASS	PREDICTED CLASS	
	Class=Yes	Class=No
Class=Yes	a (TP)	b (FN)
	c (FP)	d (TN)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

- Most widely-used metric:

$$\textit{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\textit{Precision} = \frac{TP}{TP + FP}$$

$$\textit{Recall} = \frac{TP}{TP + FN}$$

$$\textit{F1-score} = \frac{2 \times \textit{Precision} \times \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$

# Problem Statement

- **Titanic dataset**
- **Explore:** How does each feature relate to whether a person survives/alive?
- Do the EDA in more detail than usual and explain the results!
  - Splitting: 80-20, stratify: y, random\_state = 0
- **Preprocessing:**
  - \* Drop decks
  - \* Fill in the missing value using a simple imputer
  - \* One hot encoding: sex, alone
  - \* Ordinal encoding: class
  - \* Binary encoding: embark town
- **Model selection:**
  - \* Evaluation metrics used: F1\_score
  - Logistic Regression