

MACHINE LEARNING

BY,
SANDESH G

CONTENTS

Vs OF BIG DATA

CONCEPT

PROBLEM SOLVING APPROACHES

CLASSIFICATION

ALGORITHMS

MODEL EVALUATION

DATA SCIENCE METHODOLOGIES

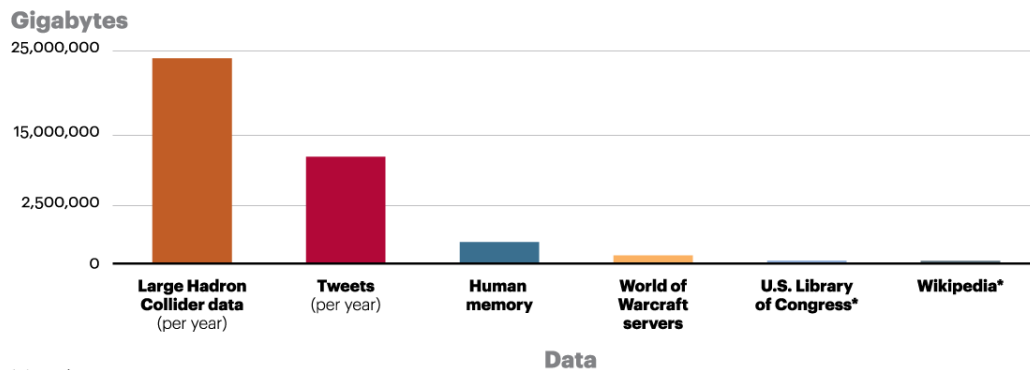
Vs of BIG DATA

VOLUME



Figure

The LHC collects about 25 million gigabytes of data per year



*Binary data

Note: All numbers are approximate.

Source: "Particle Physics Tames Big Data," Leah Hesla, *Symmetry*, 1 August 2012

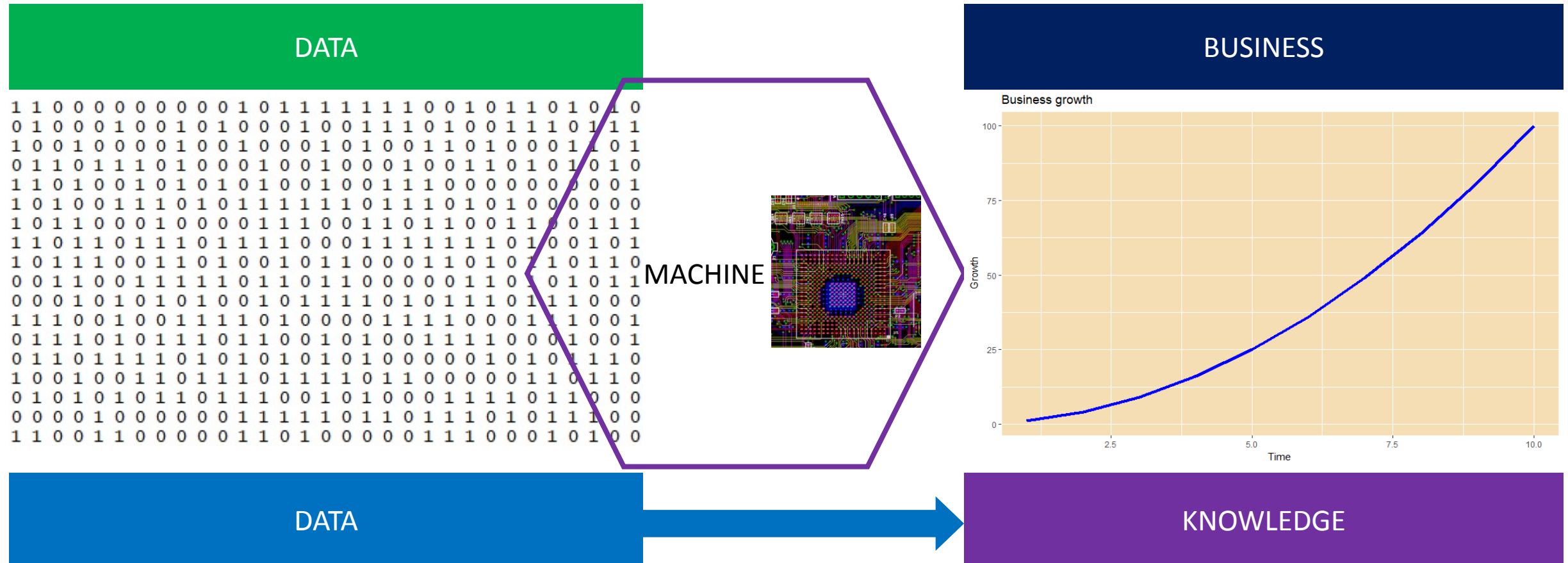
VARIETY



VELOCITY



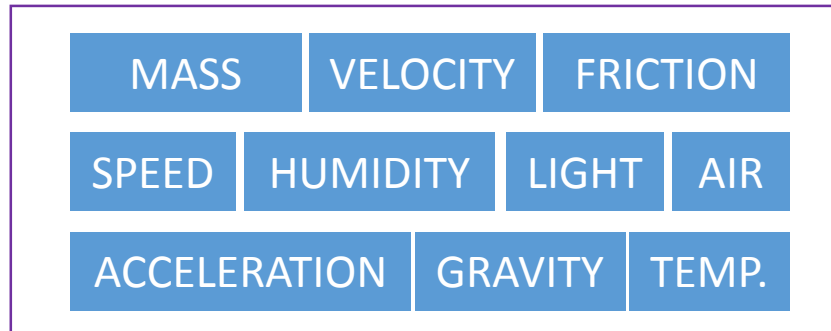
CONCEPT



WIKIPEDIA: Machine learning is a field of computer science that gives computer systems the ability to "learn" (i.e. progressively improve performance on a specific task) with data, without being explicitly programmed.

CONCEPT

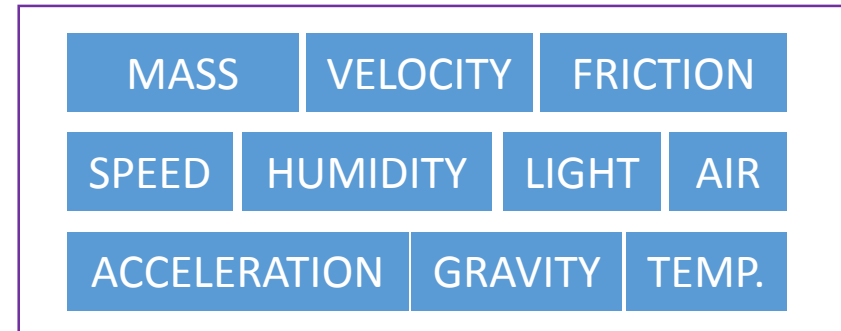
SIR ISSAC NEWTON - SCIENTIST



MATHEMATICS - CALCULUS

$\text{FORCE} = \text{MASS} \times \text{ACCELERATION}$

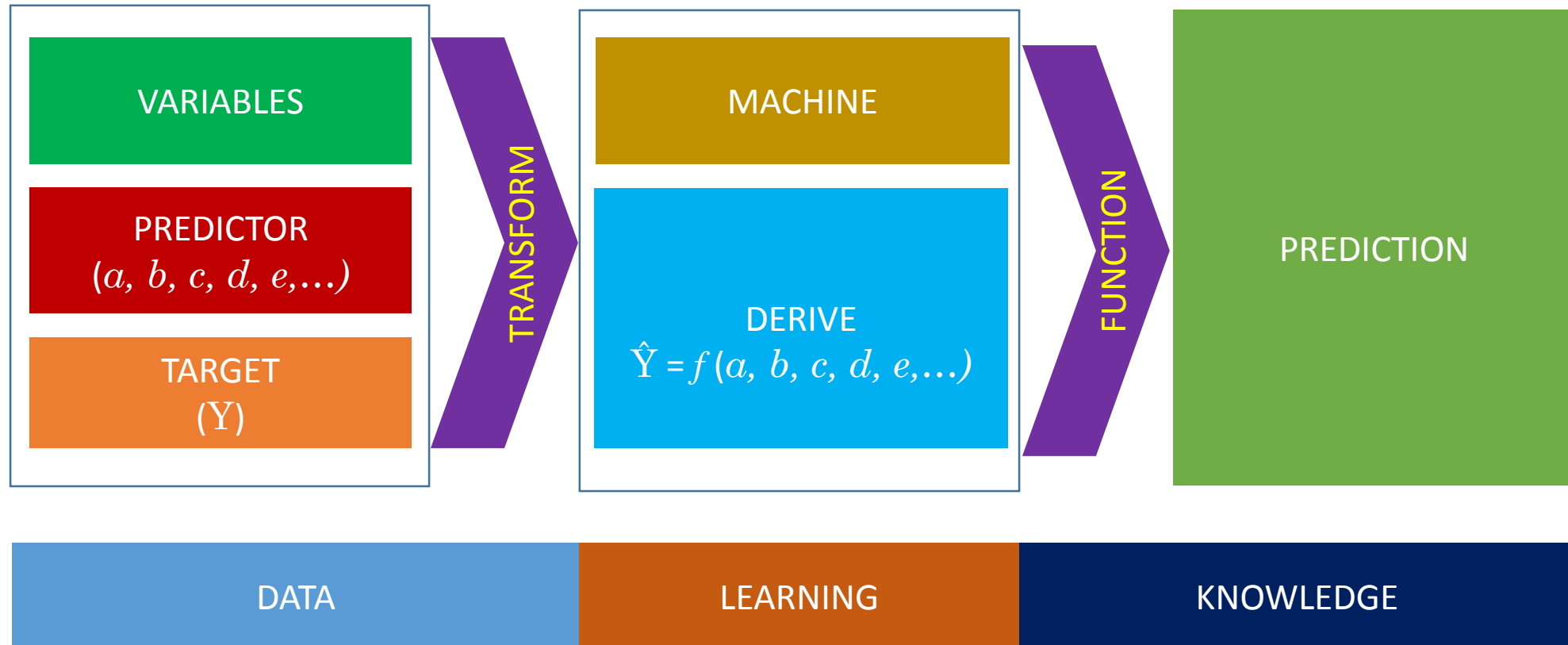
DATA SCIENTIST



MATHEMATICS - ALGORITHMS

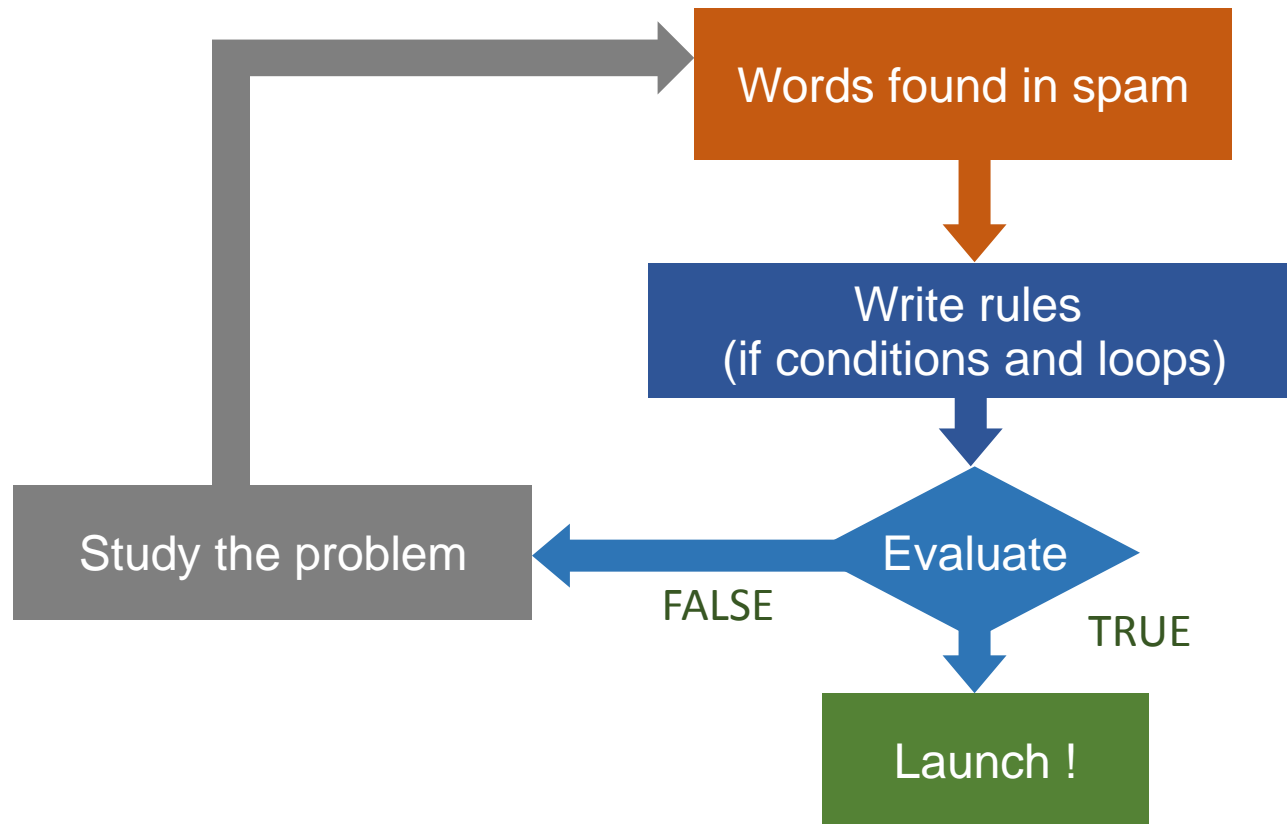
$\text{FORCE} = \text{MASS} \times \text{ACCELERATION}$

CONCEPT



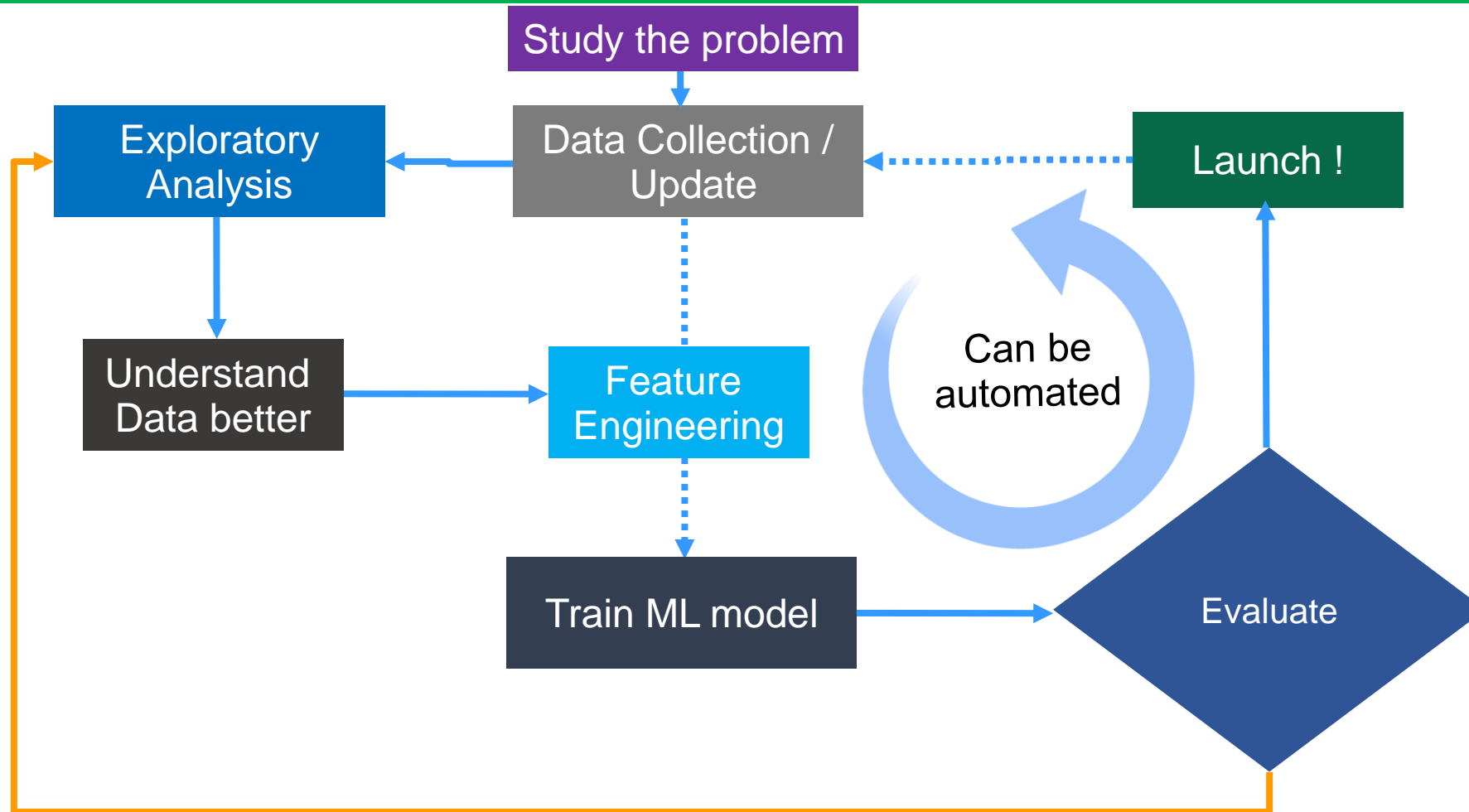
PROBLEM SOLVING APPROACHES

TRADITIONAL



PROBLEM SOLVING APPROACHES

MACHINE LEARNING



CLASSIFICATION

SUPERVISED

REGRESSION

CLASSIFICATION

UNSUPERVISED

CLUSTERING

ASSOCIATION RULE MINING

DIMENSIONALITY REDUCTION

REINFORCED

BRUTE FORCE

MONTE CARLO METHOD

ALGORITHMS

REGRESSION

LINEAR REGRESSION

NEURAL NETWORKS

DECISION TREE

RANDOM FOREST

CLASSIFICATION

K NEAREST NEIGHBOURS

NEURAL NETWORKS

LOGISTIC REGRESSION

SUPPORT VECTOR MACHINES

DECISION TREE

RANDOM FOREST

CLUSTERING

K MEANS CLUSTERING

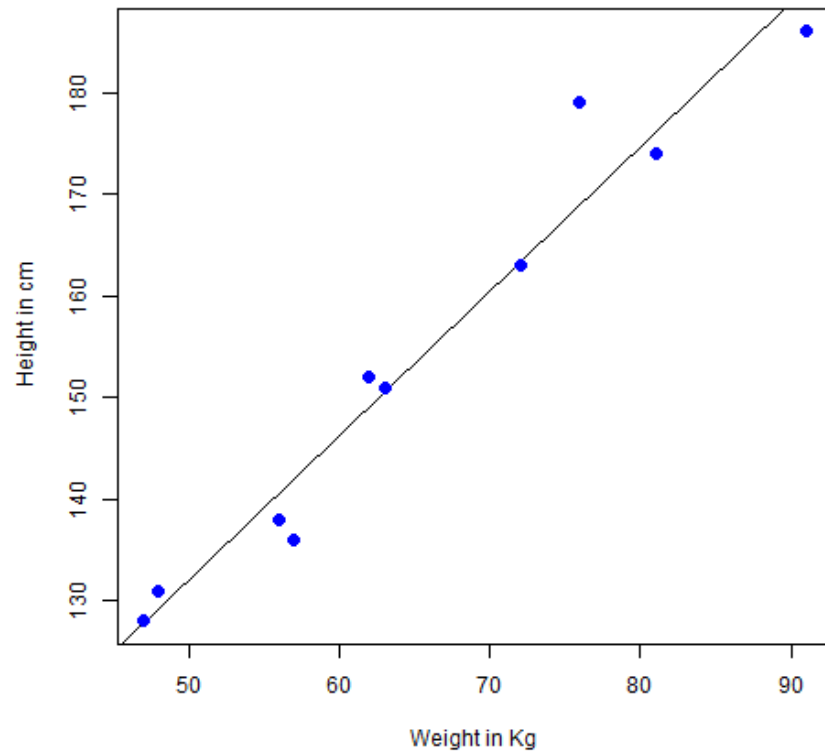
HIERARCHICAL CLUSTERING

ALGORITHMS

REGRESSION

LINEAR REGRESSION

Height & Weight Regression



$$Y = AX + B$$

MEAN

HEIGHT = NUMBER X WEIGHT + BIAS

MEDIAN

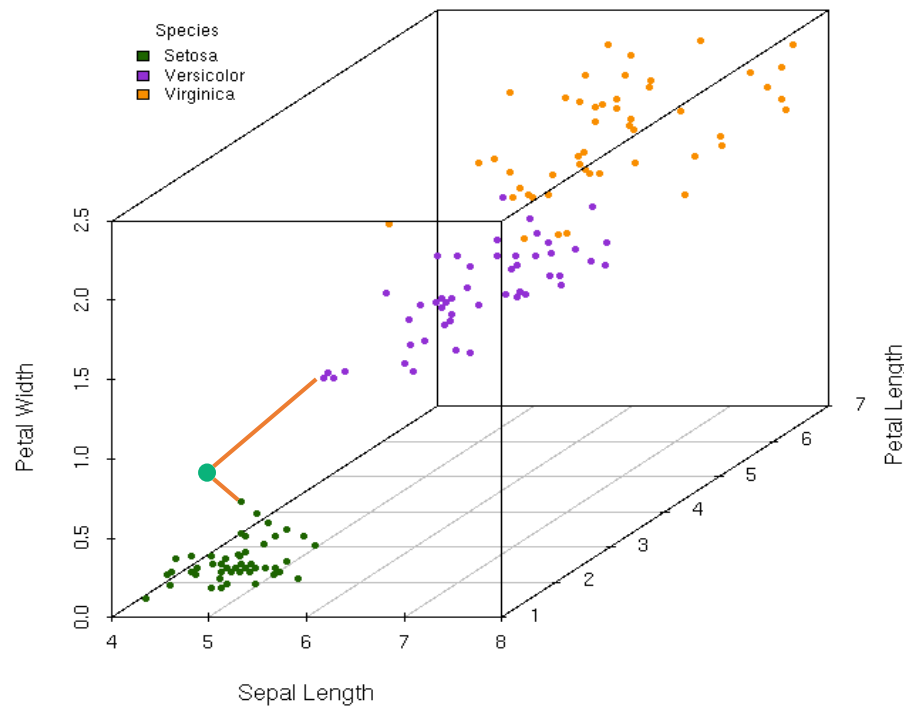
```
x <- seq(-5,5,0.1)
y <- 2 * x + 3
plot(x,y,type = "l", col = "blue")
grid()
```

ALGORITHMS

CLASSIFICATION

K - NEAREST NEIGHBOURS

3-D Scatterplot of Iris Data



DISTANCE MEASURES

EUCLIDEAN

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$

$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

MANHATTAN

$$d_1(\mathbf{p}, \mathbf{q}) = \|\mathbf{p} - \mathbf{q}\|_1 = \sum_{i=1}^n |p_i - q_i|$$

ALGORITHMS

CLUSTERING

K - MEANS

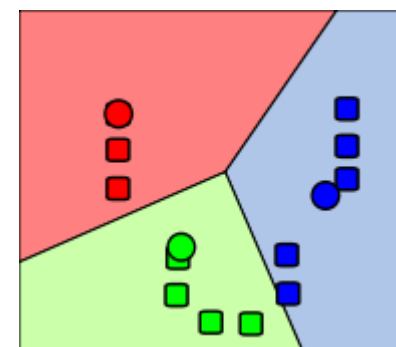
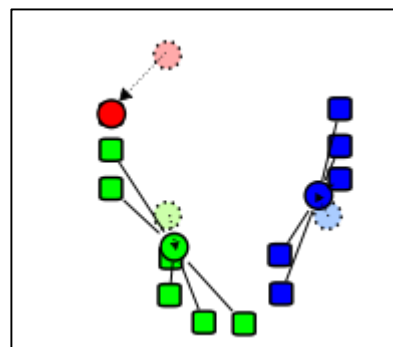
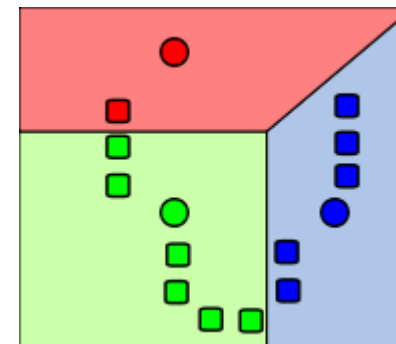
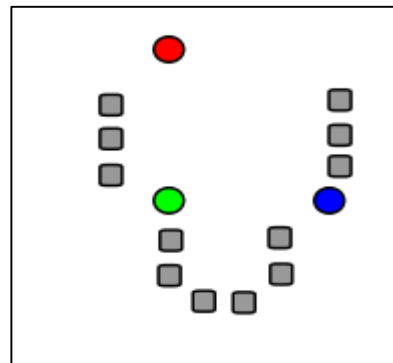
RANDOM INITIALIZATION OF
CENTROIDS

CALCULATE DISTANCES TO
NEAREST POINTS

CALCULATE MEAN OF DISTANCES

SHIFT TO AVERAGE

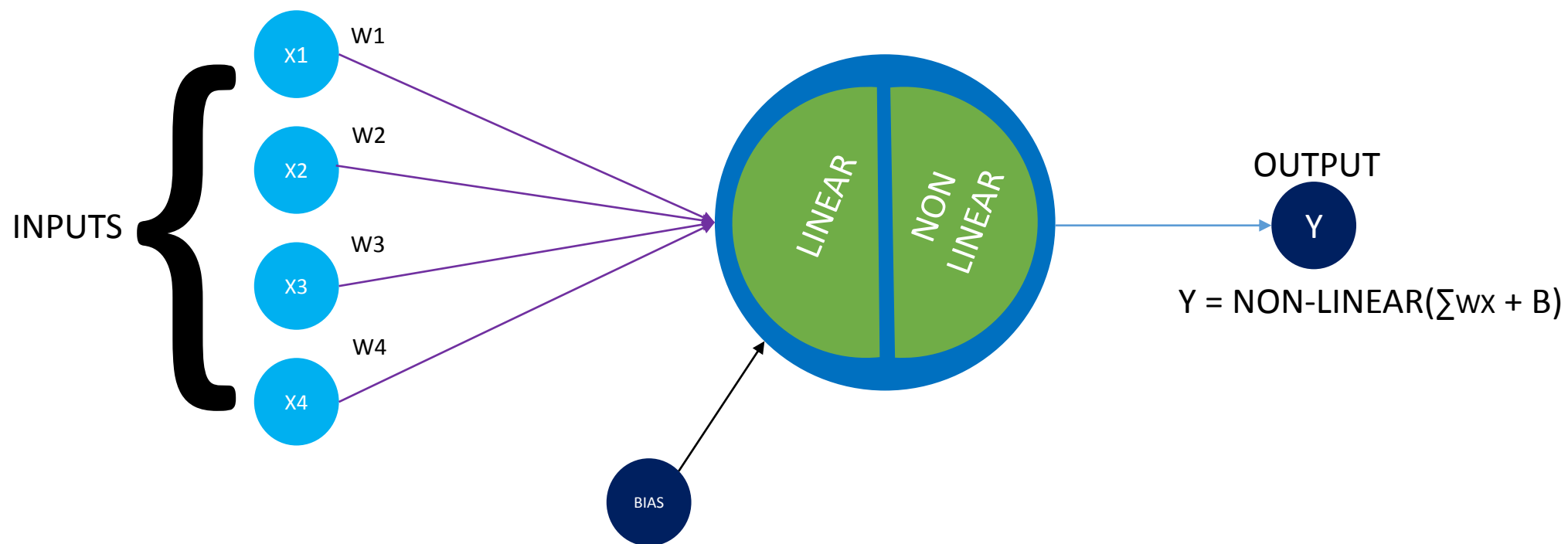
BACK TO STEP 2



ALGORITHMS

NEURAL NETWORKS

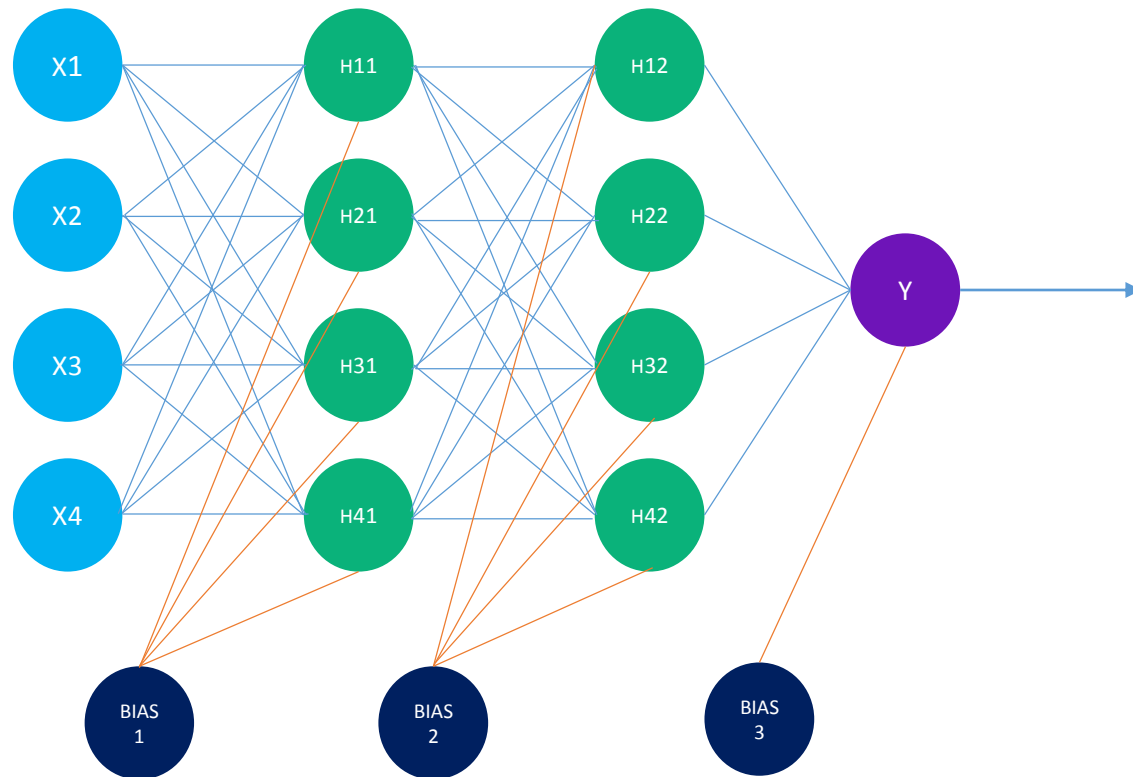
PERCEPTRON



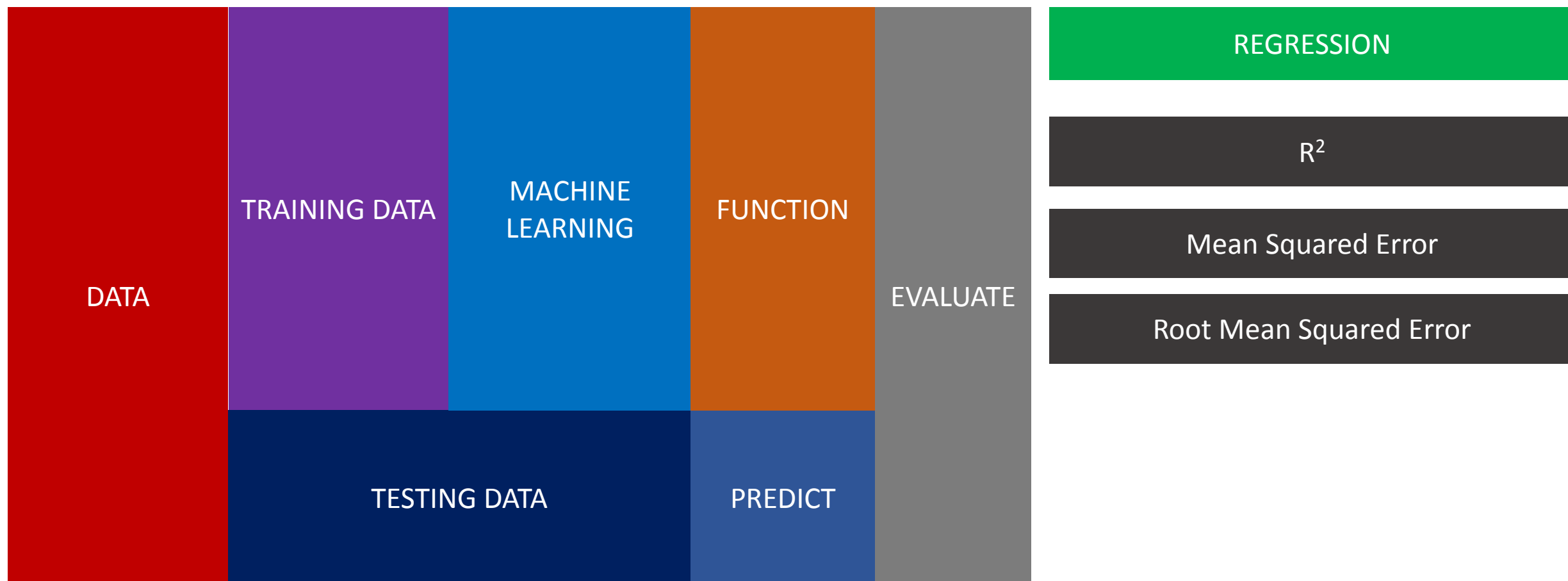
ALGORITHMS

NEURAL NETWORKS

MULTI LAYER PECEPTRON



MODEL EVALUATION



MODEL EVALUATION

CLASSIFICATION

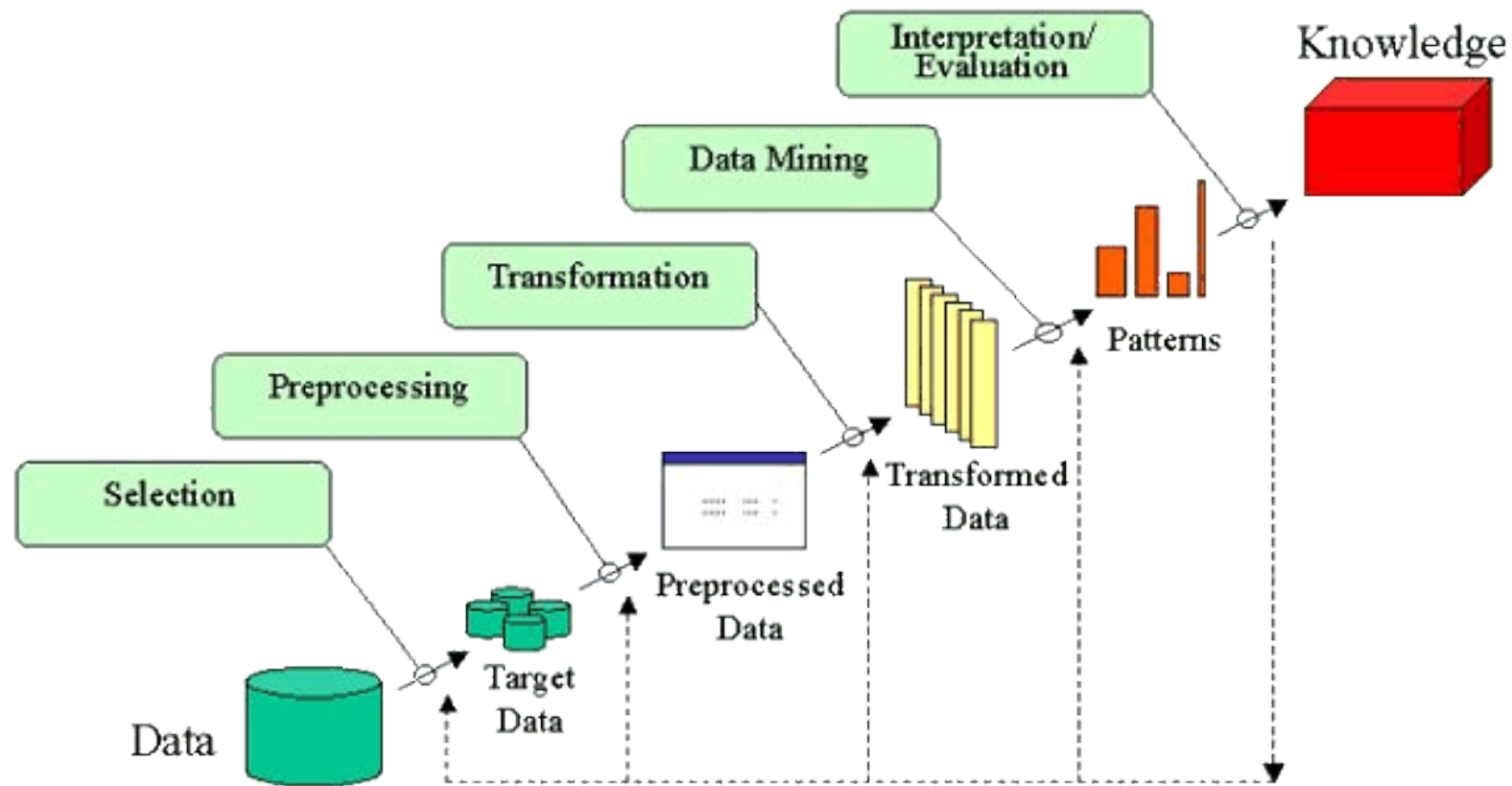
		ACTUAL		
		1	0	
PREDICTED	1	TRUE POSITIVE	FALSE POSITIVE Type I Error	Precision = TP/ PTP
	0	FALSE NEGATIVE Type II Error	TRUE NEGATIVE	
		Recall = TP/ TTP		

$$F1 \text{ MEASURE} = 2 \frac{\text{PRECISION} \times \text{RECALL}}{\text{PRECISION} + \text{RECALL}}$$

$$\text{ACCURACY} = \frac{\text{TP} + \text{TN}}{\text{ALL}}$$

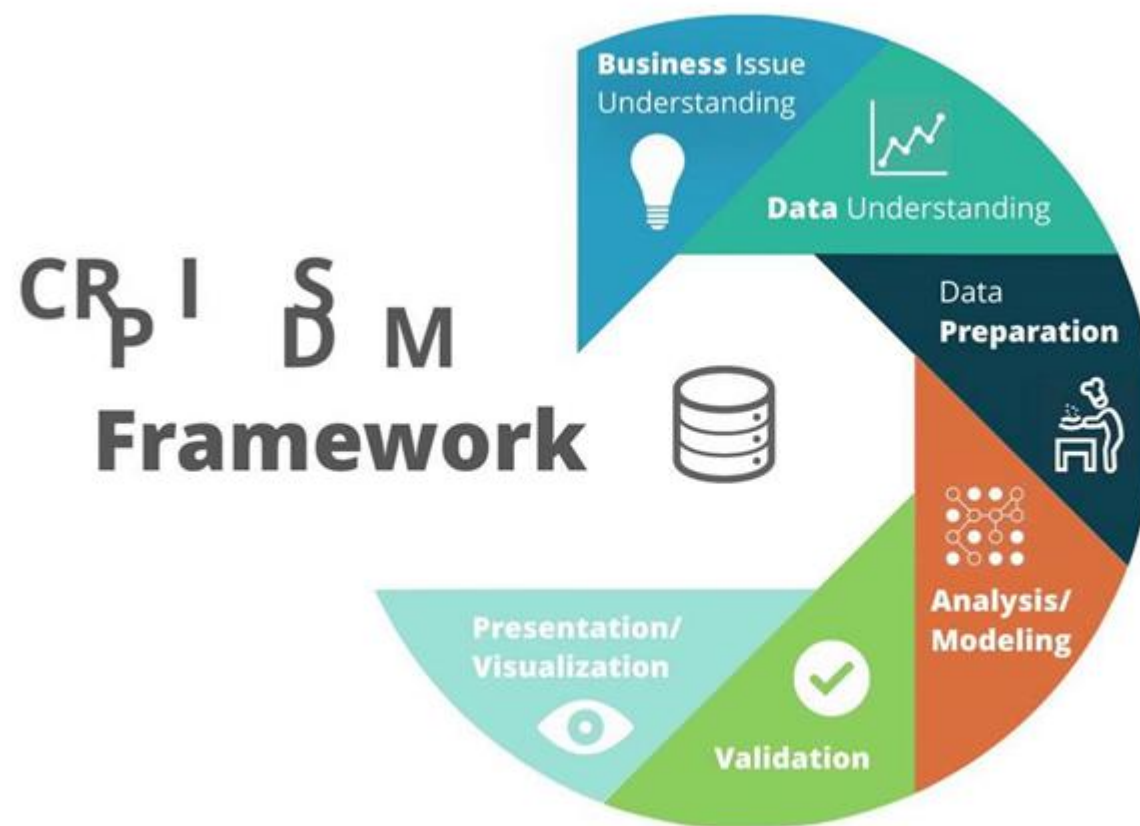
DATA SCIENCE METHODOLOGY

KNOWLEDGE DISCOVERY IN DATABASES



DATA SCIENCE METHODOLOGY

CROSS INDUSTRY STANDARD PROCESS FOR DATA MINING



Q & A

THANK YOU