

DEFINITION [ARTHUR SAMUEL 1959 - PIONEER OF AI
COINED THE TERM MACHINE LEARNING]

MACHINE LEARNING GIVES THE COMPUTERS THE ABILITY
TO LEARN WITHOUT EXPLICITLY BE PROGRAMMED.

DEFINITION (TOM MITCHELL 1998 - WRITTEN ONE OF THE FIRST BOOKS OF ML)

TOM: "THE PREVIOUS DEFINITION IS NOT WELL POSED!"



A MACHINE LEARNING ALGORITHM IS SAID TO LEARN
FROM EXPERIENCE T WITH TO RESPECT A TASK H

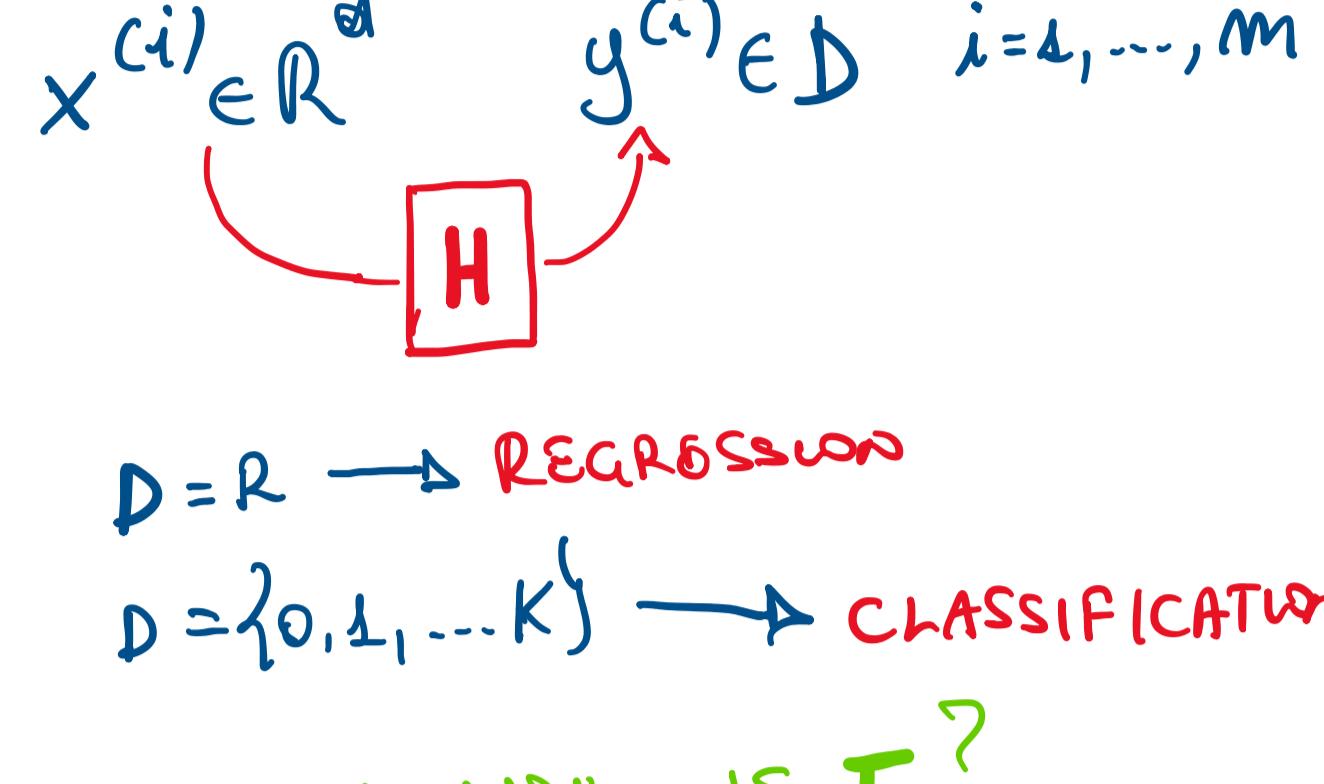
AND SOME MEASURES OF PERFORMANCE P IF ITS
PERFORMANCES ON H, AS MEASURED BY P,
IMPROVE WITH EXPERIENCE T.

WHAT "EXPERIENCE" MEANS?
WHAT "TASK" MEANS?
WHAT "PERFORMANCE" MEANS?

HOW THESE THREE
ARE RELATED EACH OTHER?

TYPE OF LEARNING ALGORITHMS

SUPERVISED LEARNING



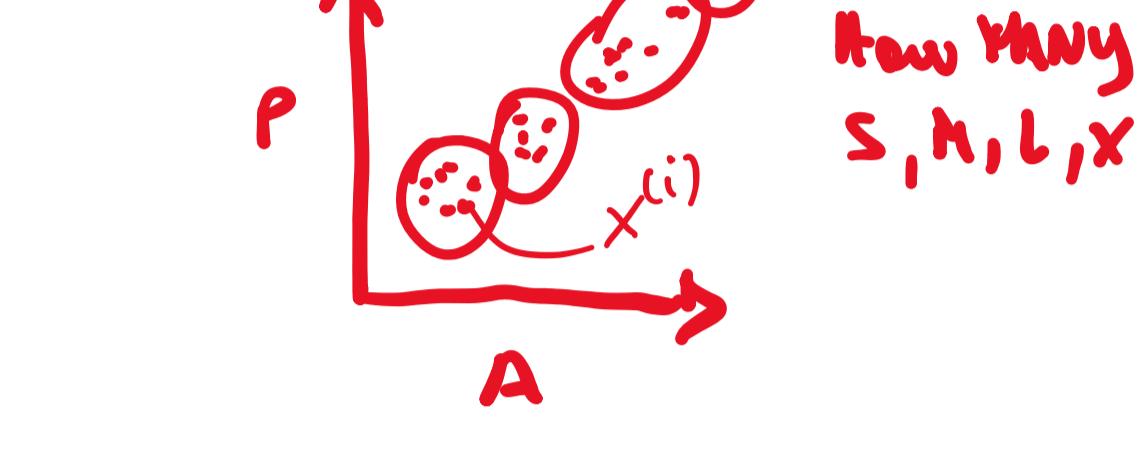
QUESTIONS: WHO IS T?
WHO IS F?
WHO IS P?
WHO IS X?
WHO IS Y?

UNSUPERVISED LEARNING

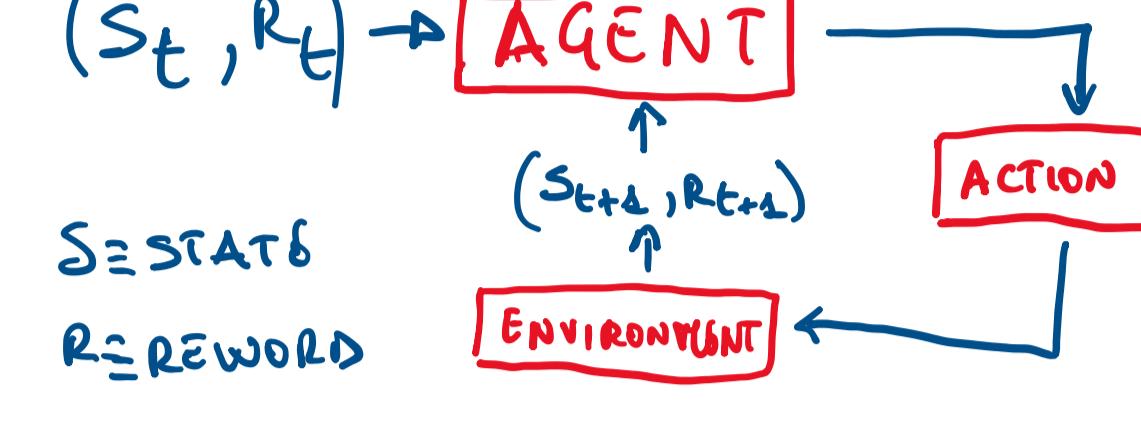
FIND H WHICH IS
ABLE TO MODEL DATA

$$x^{(i)} \in \mathbb{R}^d \quad i=1, \dots, m$$

• CLUSTERING



REINFORCEMENT LEARNING



• POLICY LEARNING

SELF-SUPERVISED

A MIX OF SUPERVISED AND UNSUPERVISED
LEARNING

MOTIVATION: PRODUCING LABELED DATA IS EXPENSIVE
(FOR SOME DOMAIN VERY HARD, EG. MEDICAL)

LEAR HOW INFANTS: LEARNING IN AN UNSUPERVISED WAY
EXPLOITING SUPERVISION GIVEN BY DATA
EG. PUZZLE \rightarrow RELATIVE POSITION

KEY INGREDIENTS OF ML ALGORITHM

• MODEL = PARAMETRIZED FUNCTION (HYPOTHESIS) H

• COST FUNCTION J

• PARAMETERS θ

• TRAINING ALGORITHM / LEARNING PROCEDURE A

• TRAINING SET X_{TRAIN}

$X = X_{\text{TRAIN}} \cup X_{\text{VAL}} \cup X_{\text{TEST}}$ (EXPERIMENTAL DATASET)

• VALIDATION SET X_{VAL}

$X_{\text{TRAIN}} \cap X_{\text{VAL}} = \emptyset \quad X_{\text{VAL}} \cap X_{\text{TEST}} = \emptyset$

• TEST SET X_{TEST}

$X_{\text{TRAIN}} \cap X_{\text{TEST}} = \emptyset$

• EVALUATION MEASURES P

NOTATION

INPUT $x^{(i)} = [x_1^{(i)}, \dots, x_d^{(i)}]^T; \quad X^{(i)} = \begin{bmatrix} x_{1,1}^{(i)}, x_{1,2}^{(i)}, \dots, x_{1,d}^{(i)} \\ \vdots \\ x_{n,1}^{(i)}, x_{n,2}^{(i)}, \dots, x_{n,d}^{(i)} \end{bmatrix}; \quad X^{(i)} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$

EXAMPLES: A PATIENT MEDICAL RECORD EXAMPLES: A GRAY SCALE IMAGE

$\theta = [\theta_0, \theta_1, \dots, \theta_d]^T$ PARAMETERS OF THE MODEL

$x_j^{(i)} \quad i=1, \dots, m \quad j=1, \dots, d$
SIZE OF THE TRAINING SET
INPUT COMPONENT OF A ONE-DIMENSIONAL VECTOR

TENSOR

EXAMPLE: A SERIES OF RGB IMAGES OVER TIME

FOR MATH CONVENIENCE, WE WILL USUALLY EXTEND $x^{(i)}$ ADDING $x_0^{(i)} = 1$

$$x^{(i)} = [x_0^{(i)}, x_1^{(i)}, \dots, x_d^{(i)}]^T$$

ASSOCIATED TO θ_0

HOW CAN WE REPRESENT THE EXPERIENCE T?

$$\begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_d^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_d^{(m)} \end{bmatrix} \quad \text{TRAINING SET
(IN THIS CASE OUR PATTERNS ARE VECTORS)}$$

DATA NORMALIZATION

$$\hat{x}_j = \frac{x_j - x_j^{\min}}{x_j^{\max} - x_j^{\min}} \in [0, 1]$$

HAVING
FEATURES
TRAINING SET
PARTICLES

$$x_j^{\min} = \min \{x_j^{(i)} \mid i=1, \dots, m\}$$

$$x_j^{\max} = \max \{x_j^{(i)} \mid i=1, \dots, m\}$$

MOST USED: ZERO MEAN NORMALIZATION

$$\tilde{x}_j = \frac{x_j - \mu_j}{\sigma_j}$$

DATA WILL HAVE ZERO MEAN
AND STANDARD DEVIATION 1

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

MEAN

$$\sigma_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2}$$

HOW "SHAPE" OF MY DATA DISTRIBUTION CHANGE?

WHY THIS IS USEFUL?

