

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 EXPERIENCES T WITH RESPECT TO A TASK H

AND SOME MEASURES OF PERFORMANCE P IF ITS

PERFORMANCES ON H, AS MEASURED BY P,

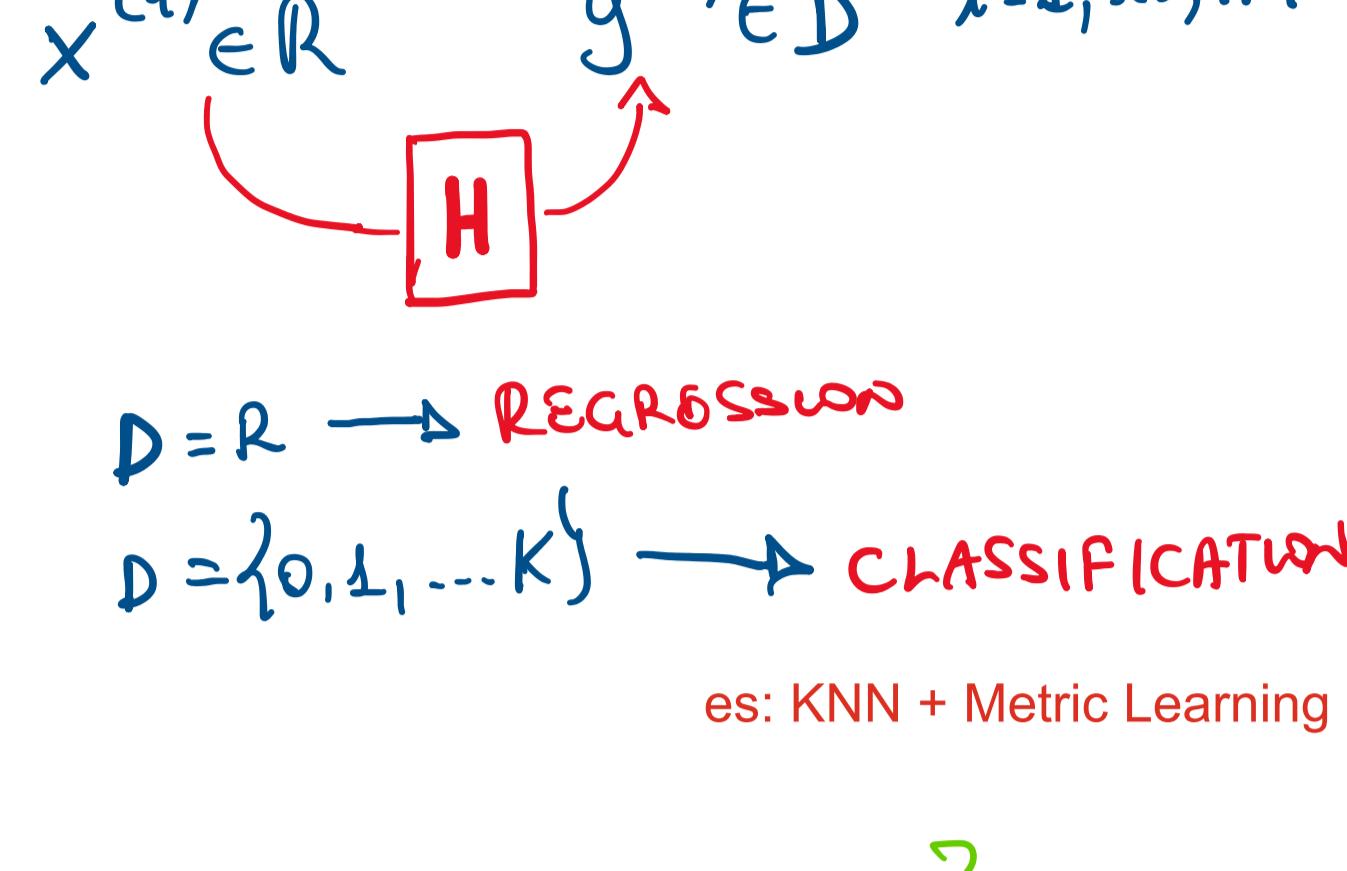
IMPROVE WITH EXPERIENCES.

WHAT "EXPERIENCES" MEAN?
WHAT "TASK" MEANS?
WHAT "PERFORMANCE" MEANS?

HOW THESE THREE
ARE RELATED EACH OTHER? See Notes in file Basic_Concepts_2.pdf

TYPE OF LEARNING ALGORITHMS

SUPERVISED LEARNING

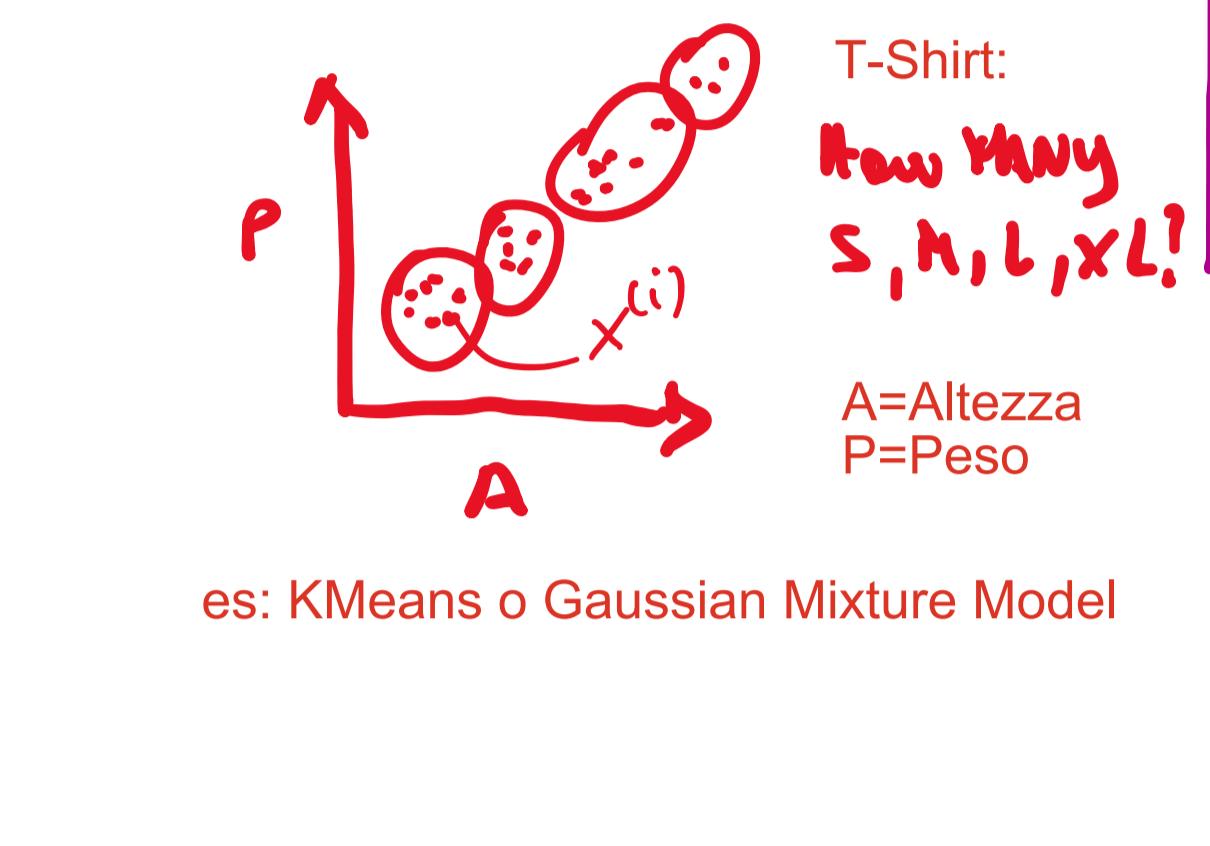


QUESTIONS: WHAT IS T?
WHAT IS H?
WHAT IS P?
WHAT IS X?
WHAT IS Y?

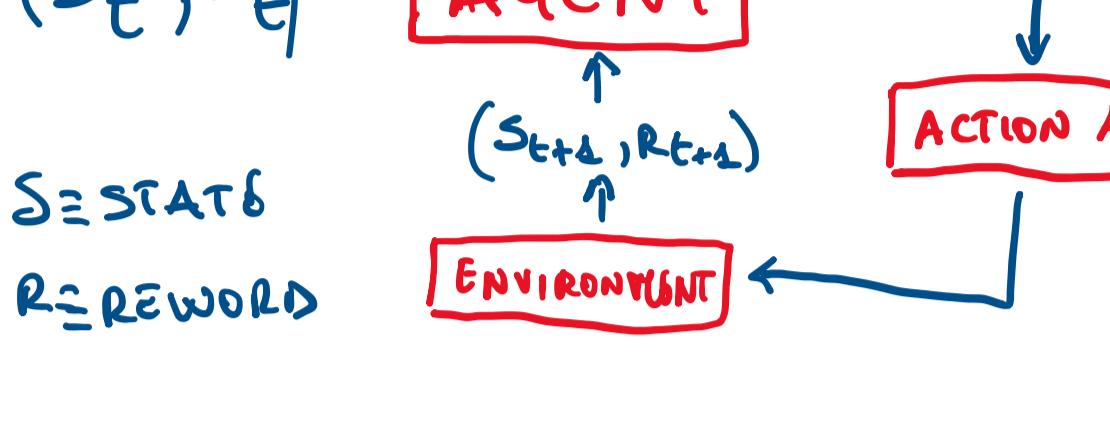
UNSUPERVISED LEARNING

FIND H WHICH IS
ABLE TO MODEL DATA

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



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}

• VALIDATION SET X_{VAL}

• TEST SET X_{TEST}

• EVALUATION MEASURES P

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

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

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

NOTATION

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

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

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 EXPERIENCES T?

$$\begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_d^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_d^{(m)} \end{bmatrix}$$

TRAINING SET
(IN THIS CASE OUR PATTERNS ARE VECTORS)

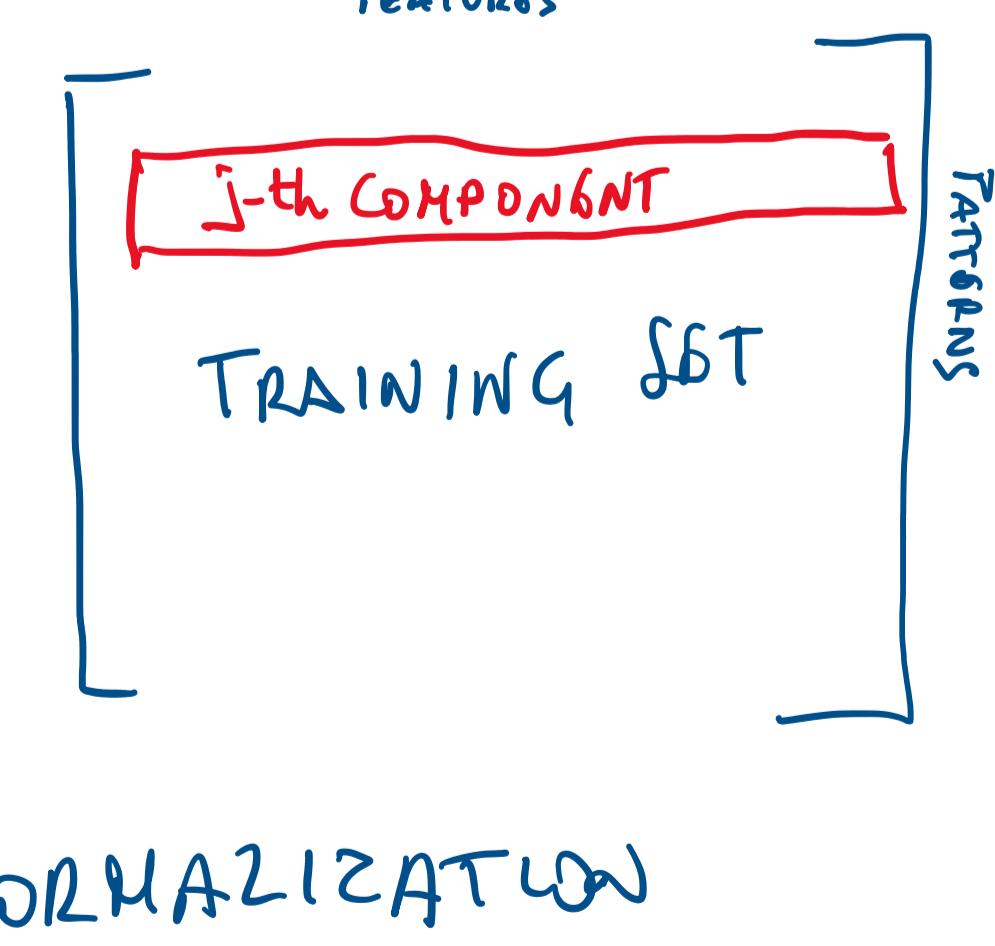
DATA NORMALIZATION

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

HAPPING

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

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



MOST USED: ZERO MEAN NORMALIZATION

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

DATA WILL HAVE ZERO MEAN
AND STANDARD DEVIATION 1

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

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

HOW "SHAPE" OF MY DATA DISTRIBUTION CHANGE?
WHY THIS IS USEFUL?

