

# Machine Learning Foundations

## **Introduction, Terminology and Setup**

Harish Guruprasad Ramaswamy  
IIT Madras

# Outline

1.What is Machine Learning?

2.The Wonders of Machine Learning

3.Data, Models and ML Tasks

4.Supervised Learning

    1. Regression

    2. Classification

5. Unsupervised Learning

    1. Dimensionality Reduction

    2. Density Estimation

# Machine Learning Definition

Machine learning (ML) is the study of computer algorithms that improve automatically through experience and by the use of data.

# ML Tasks You (Might Have) Performed Today

## Weather prediction



32

°C | °F

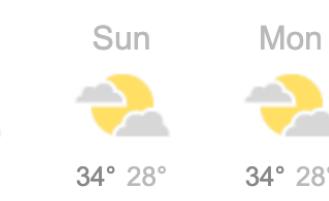
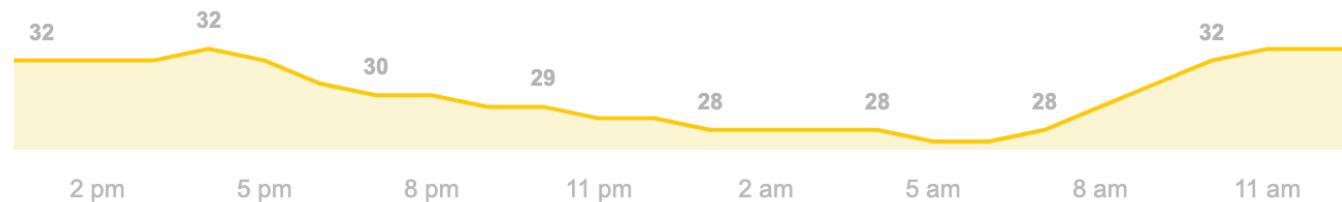
Precipitation: 0%  
Humidity: 66%  
Wind: 19 km/h

Andheri West, Mumbai,  
Maharashtra

Monday

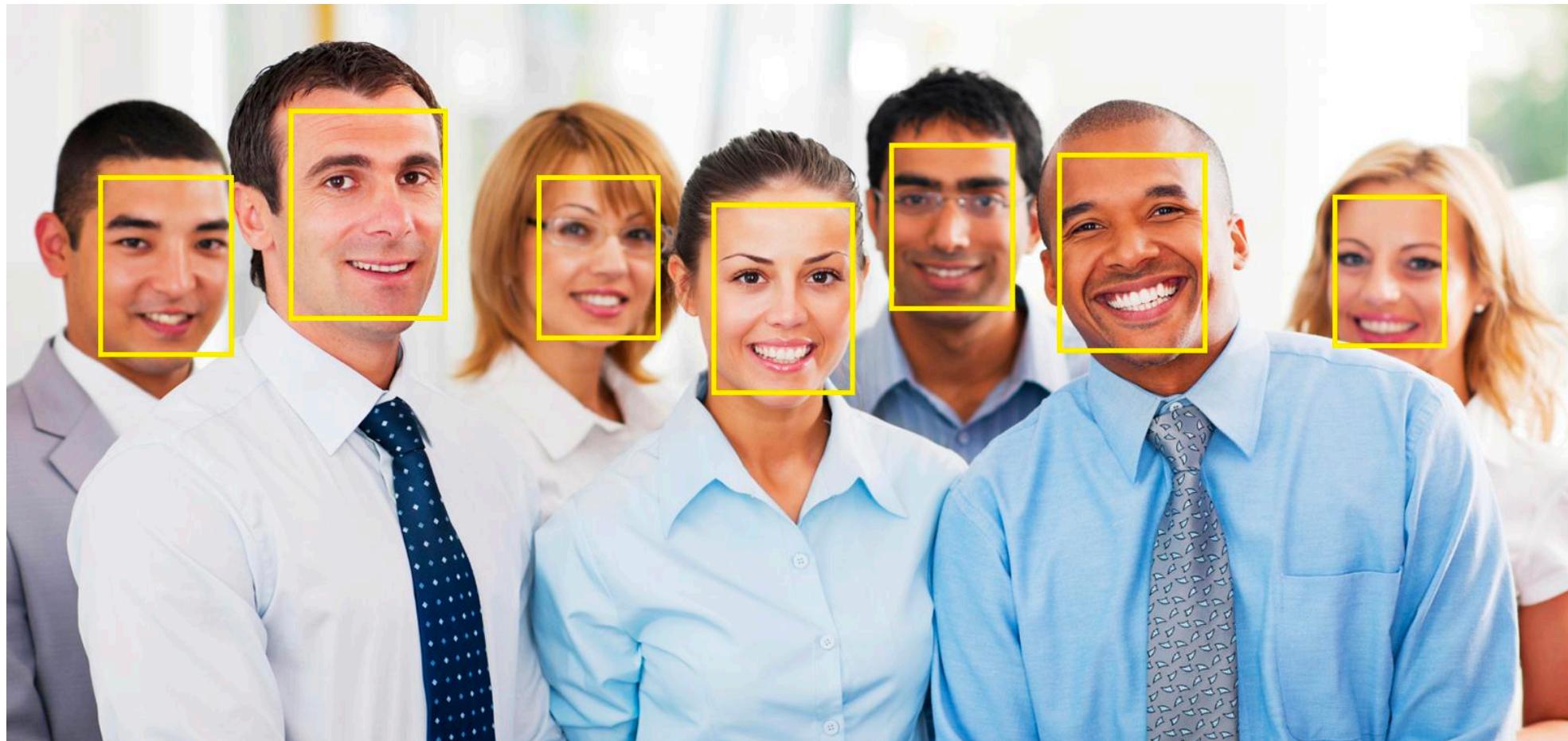
Mostly sunny

Temperature | Precipitation | Wind



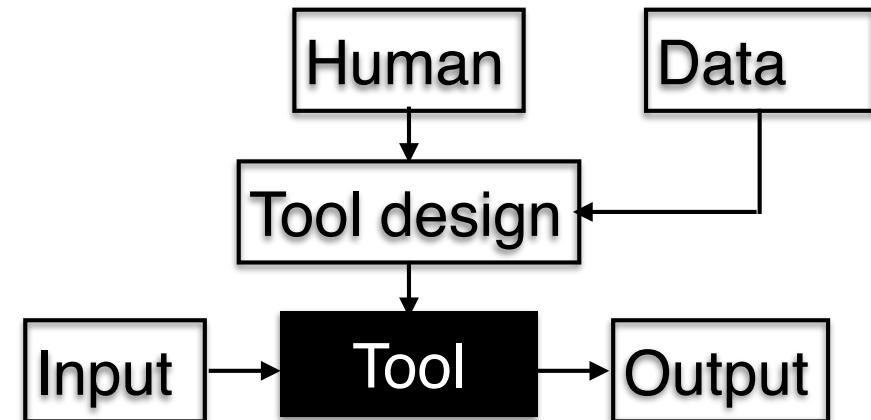
# ML Tasks You (Might Have) Performed Today

## Face Detection



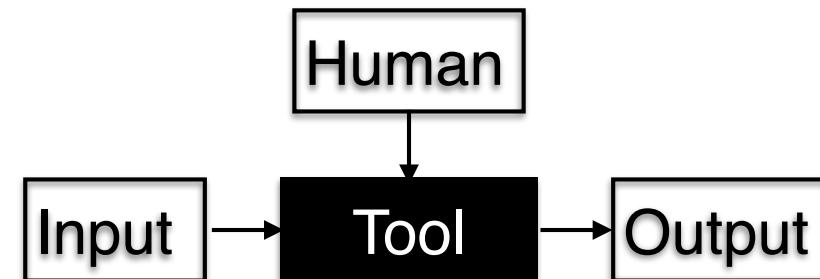
# Task Hierarchy

Machine Learning



Programming

Tool usage



Manual Labour



# Why and When Machine Learning?

## 1. Programming/Human Labour Fails

1. Scale/Speed/Cost of human labor
2. Inability to express rules using language.
3. Don't know the exact rules transforming input to output.

## 2. Machine Learning can succeed.

1. Have lots of example data
2. Have some structural idea on the rules

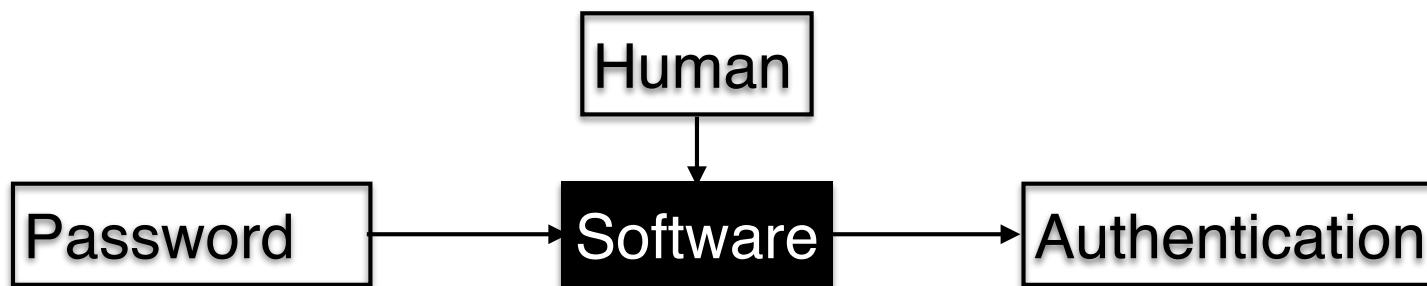
# Task Analysis : Password Verify

## Manual labour



---

## Programming



# Task Analysis : Password Verify

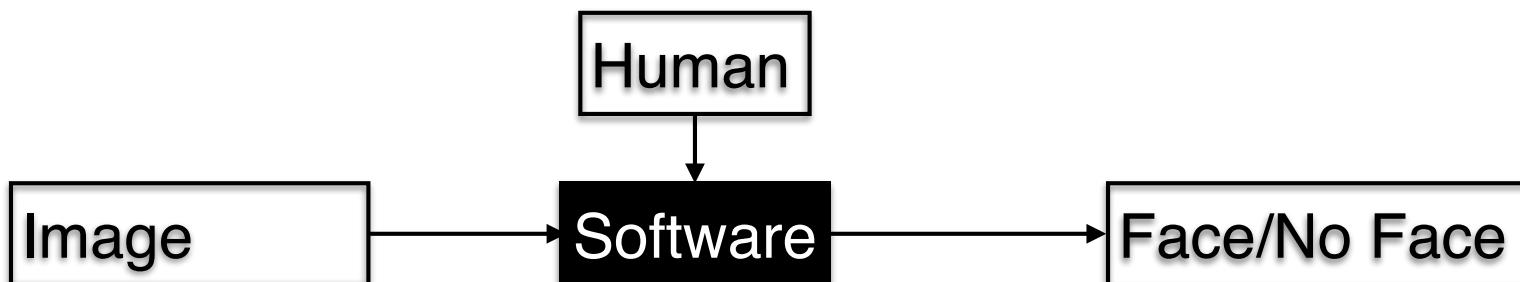
- Problems with Manual Labour
  - Scale: Having humans check login details of every login is impractical.
- Problems with programming
  - None.
- Machine Learning not required.

# Task Analysis : Face Detection

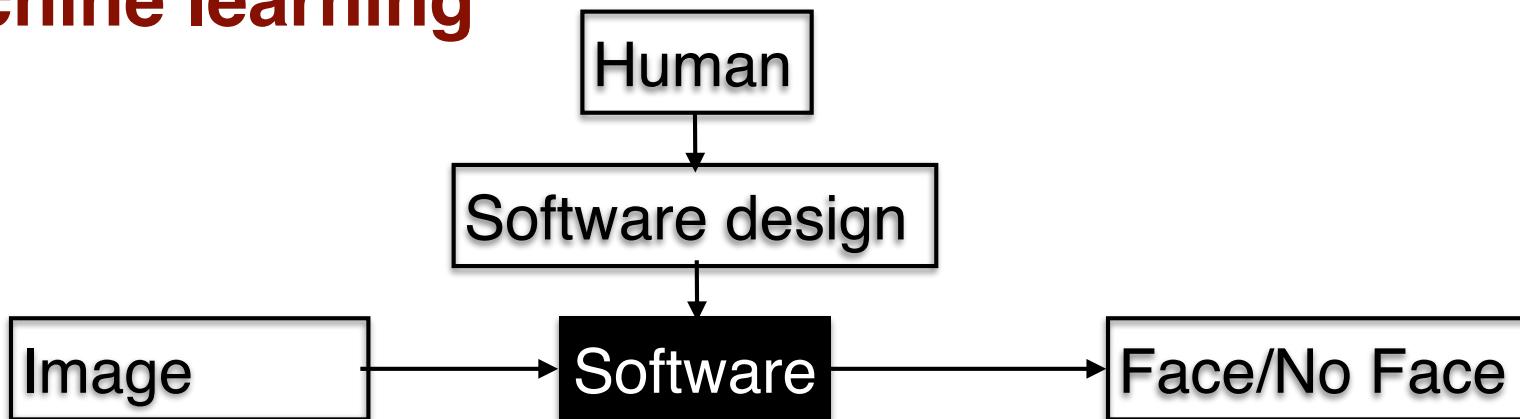
## Manual labour



## Programming



## Machine learning



# Task Analysis : Face Tagging

- Problems with Manual Labour
  - Scale: Having humans check all faces of every image is impractical.
- Problems with programming
  - Expressing face/not face in code is impossible.
- Case for Machine Learning
  - Lots of images available.

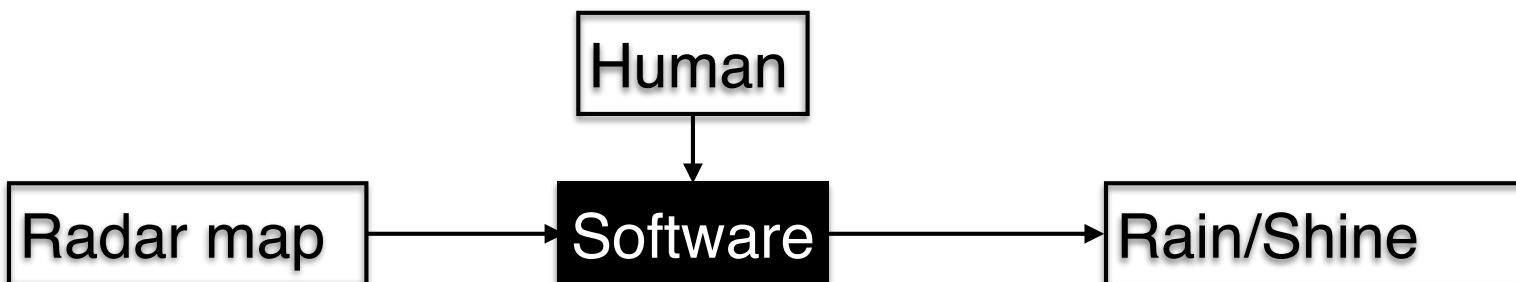
# Task Analysis : Weather Prediction

## Manual labour



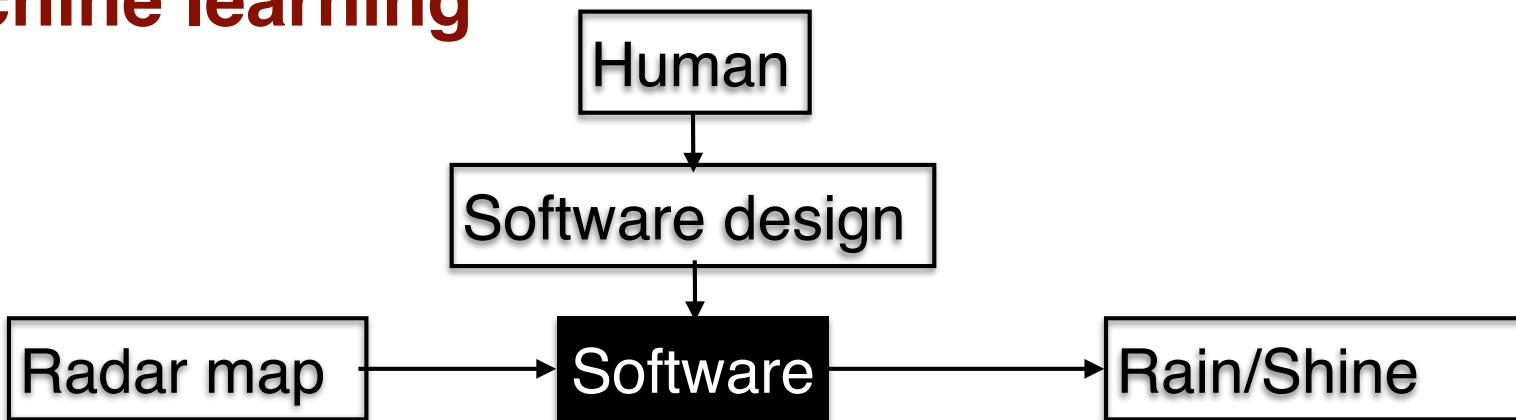
---

## Programming



---

## Machine learning



# Task Analysis : Weather Prediction

- Problems with Manual Labour
  - Humans just do not know the full rules and can't process that much information.
- Problems with programming
  - Cannot code unknown rules.
- Case for Machine Learning
  - Lots of weather data available.

# Outline

1.What is Machine Learning??

## **2.The Wonders of Machine Learning**

3.Data, Models and ML Tasks

4.Supervised Learning

    1. Regression

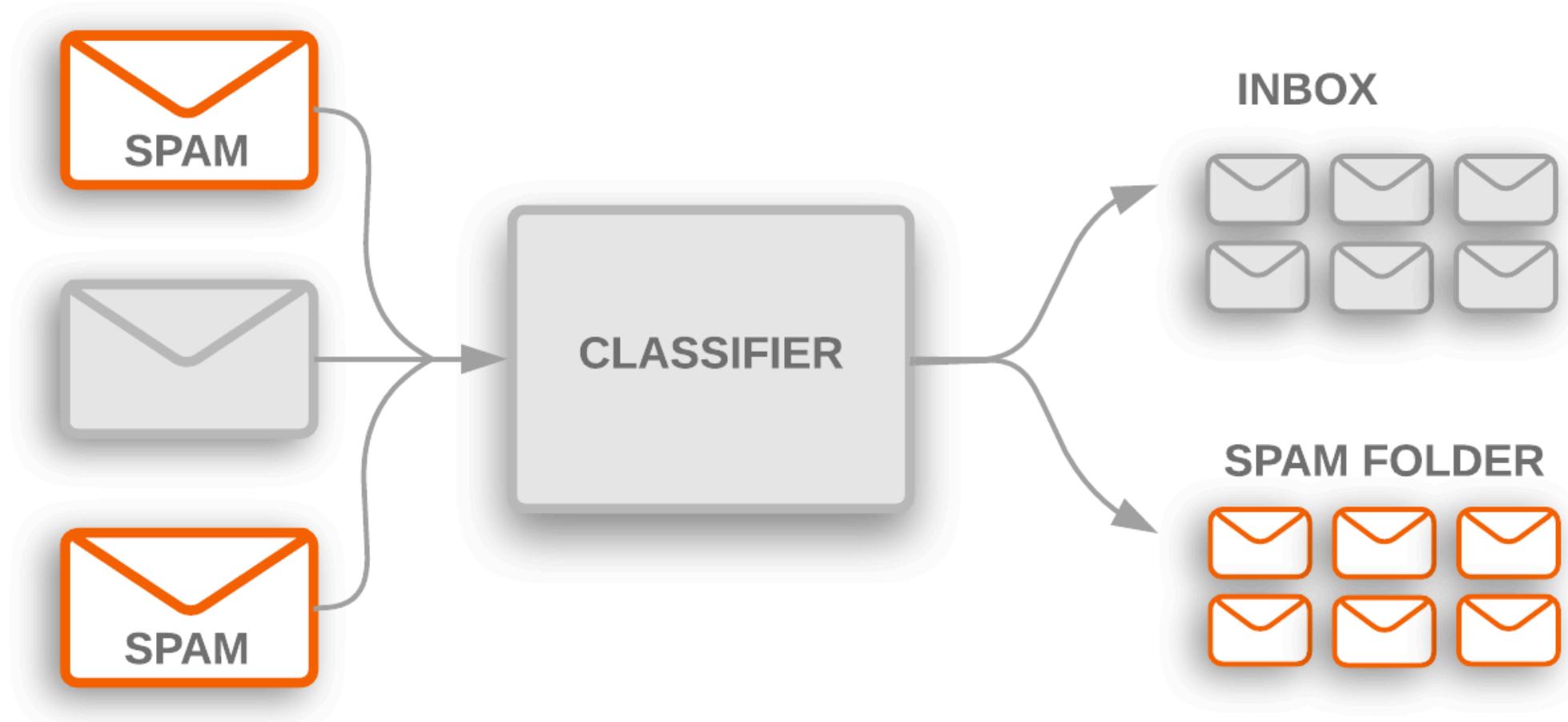
    2. Classification

5. Unsupervised Learning

    1. Dimensionality Reduction

    2. Density Estimation

# Machine Learning in your Inbox



# Machine Learning in your Shopping Cart

## Frequently Bought Together



Price For All Three: \$258.02

[Add all three to Cart](#)

- [This item: The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition \(Springer Series in Statistics\) by Trevor Hastie](#)
- [Pattern Recognition and Machine Learning \(Information Science and Statistics\) by Christopher M. Bishop](#)
- [Pattern Classification \(2nd Edition\) by Richard O. Duda](#)

## Customers Who Bought This Item Also Bought



[All of Statistics: A Concise Course in Statistical Inference](#) by Larry Wasserman  
 (8) \$60.00



[Pattern Classification \(2nd Edition\)](#) by Richard O. Duda  
 (27) \$117.25



[Data Mining: Practical Machine Learning Tools and Techniques](#) by Ian H. Witten and Eibe Frank  
 (29) \$41.55

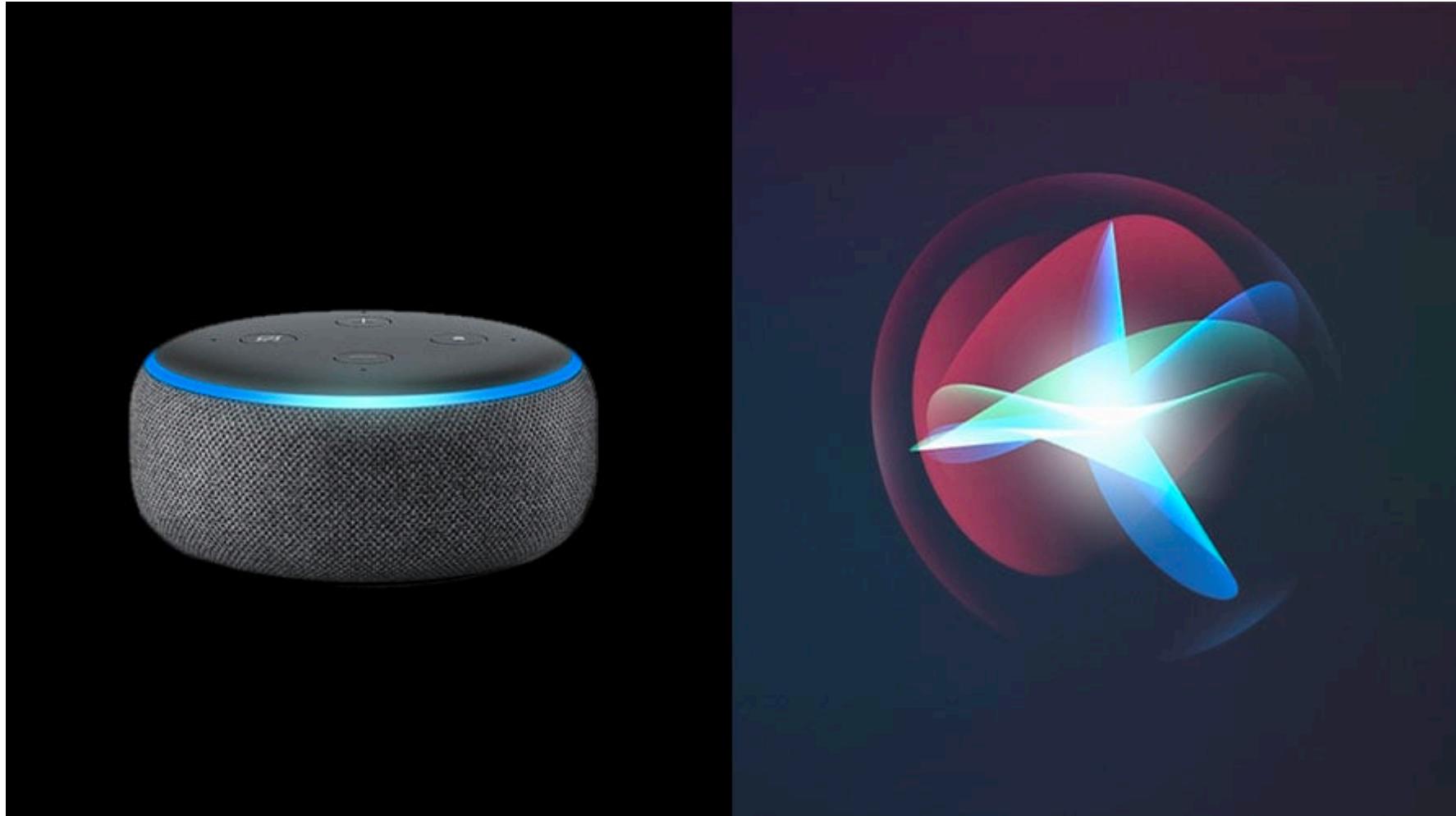


[Bayesian Data Analysis, Second Edition \(Texts in Statistical Science\)](#) by Andrew Gelman, John Carlin, Hal Stern, and Donald Rubin  
 (10) \$56.20

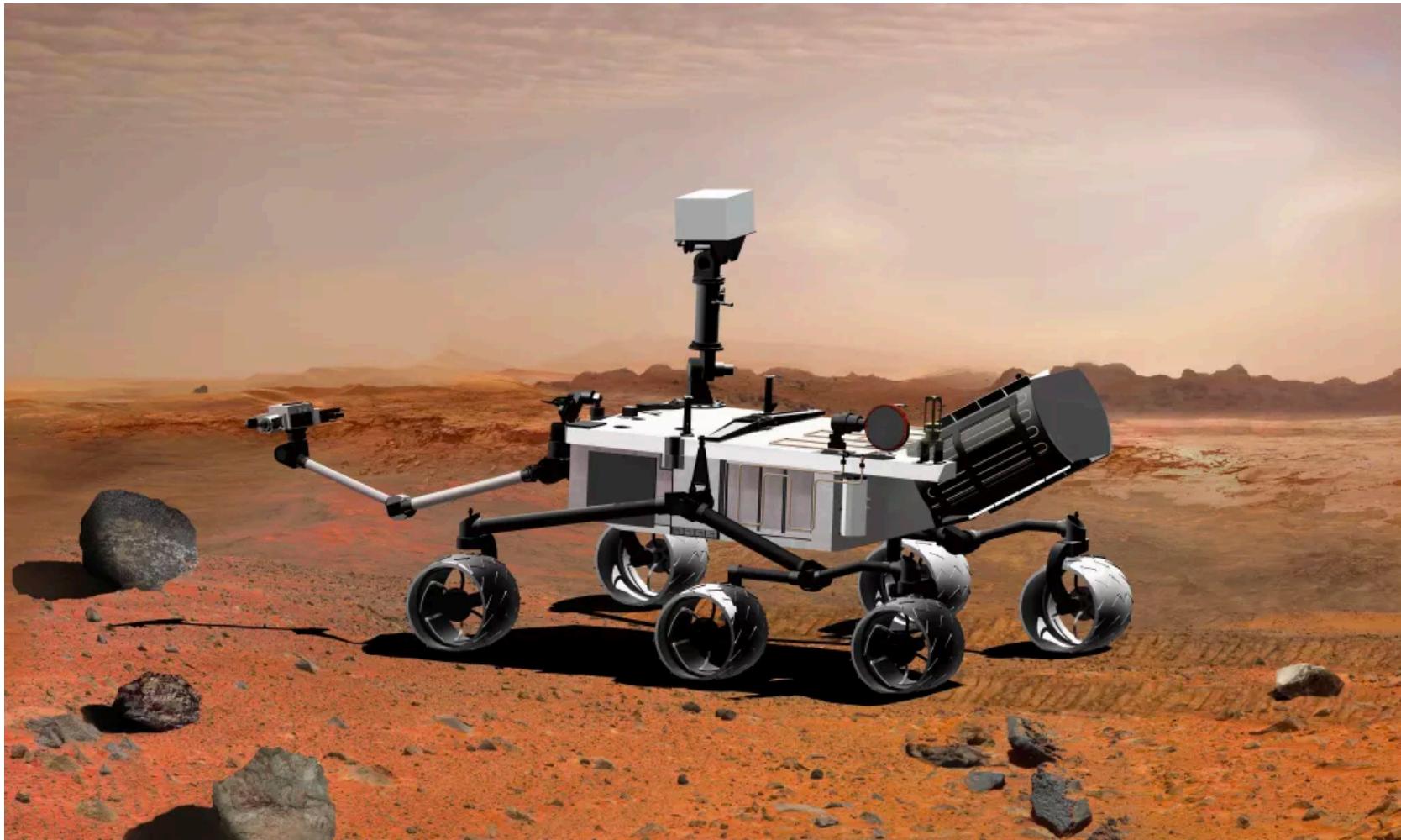


[Data Analysis Using Regression and Multilevel / Mixed Models](#) by Andrew Gelman and Jennifer Hill  
 (13) \$39.59

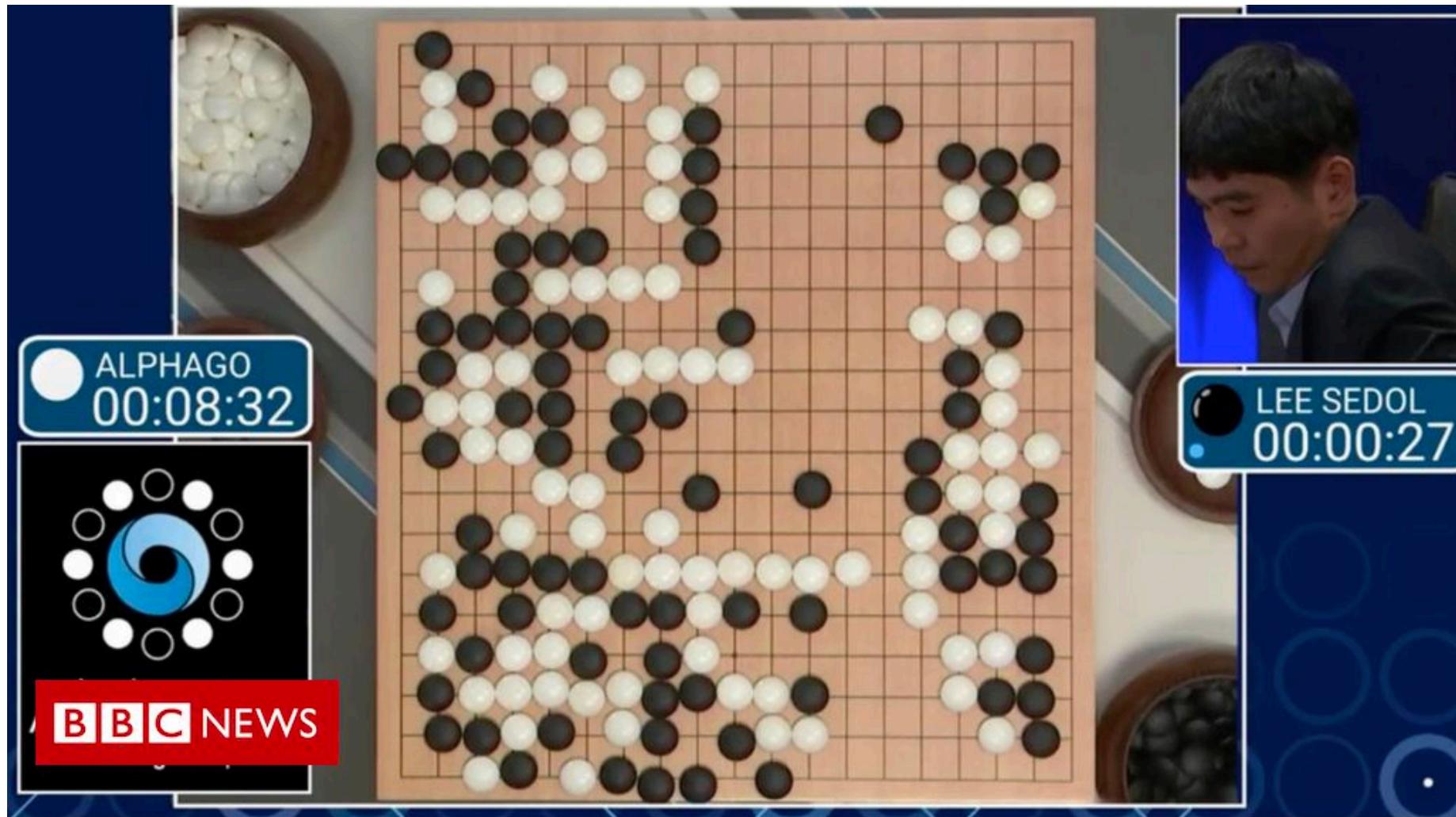
# Machine Learning in your Smart Assistant



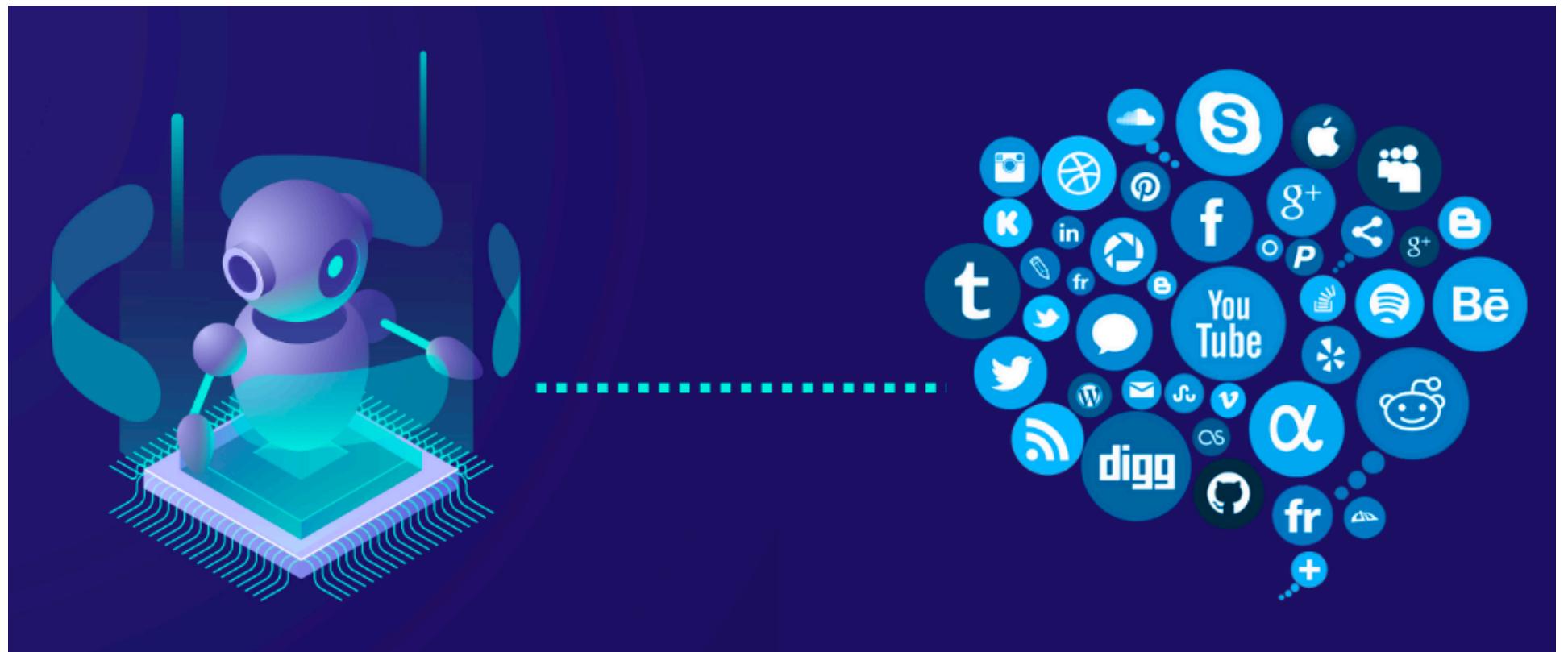
# Machine Learning in Robot Als



# Machine Learning in Games



# Machine Learning in Marketing



# Outline

1.What is Machine Learning??

2.The Wonders of Machine Learning

**3.Data, Models and ML Tasks**

4.Supervised Learning

    1. Regression

    2. Classification

5. Unsupervised Learning

    1. Dimensionality Reduction

    2. Density Estimation

# What is Data?

Data is a collection of vectors.

E.g.



3	9	1.9	5.0	House 1
2	7	2.1	3.2	House 2
4	12	2.8	6.6	House 3
5	16	0.9	9.8	House 4
5	15	3.1	8.5	House 5
4	11	1.6	6.9	House 6

Metadata is information on the data.

E.g. : (# rooms, Area in 100 sq.ft, Distance to metro in km, Price in 10 lakhs)

# What is a Model?

A model is a mathematical simplification of reality.

Some examples:

The Ideal Gas model

Inverse square law for gravitational attraction

Moore's Law for semiconductors

Cobb–Douglas model in Economics

*"All models are wrong, but some are useful"*

George Box

# Types of Models in ML

- Predictive Model
  - Regression Model
  - Classification Model
  - ....
- Probabilistic Model
  -

# Predictive Models

## Regression Model

Model the price of a house based on its area and distance to metro.

Example good model:

$$\text{Price} = 0.5 * \text{Area} - \text{Distance}$$

# Predictive Models

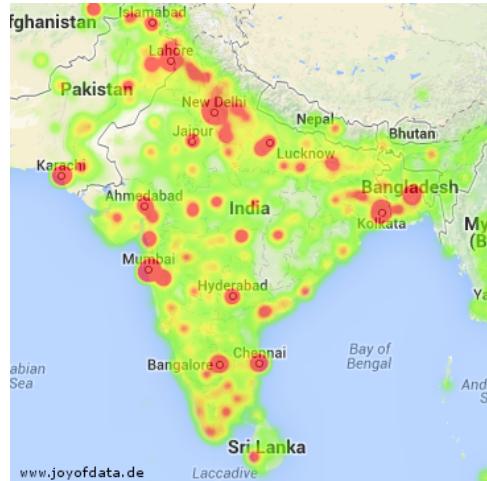
## Classification Model

Model whether a house is closer than 2kms to a metro based on price and area

Example good model:

Answer = Close if  $2 * \text{ROOMS} - \text{PRICE} < 1$   
Far otherwise

# Probabilistic Models



What is the probability that a randomly chosen person is in lat-long : (25N,30E) ?



"A formless void transforms total mysteries"

[RECEIVE MORE WISDOM...](#)

[Tweet the wisdom](#)

What is the probability that a given tweet was generated by Mr. Chopra?

# Learning Algorithms

Learning Algorithms: Data → Models

Choose from a collection of models, with same structure but different **parameters**.

E.g.

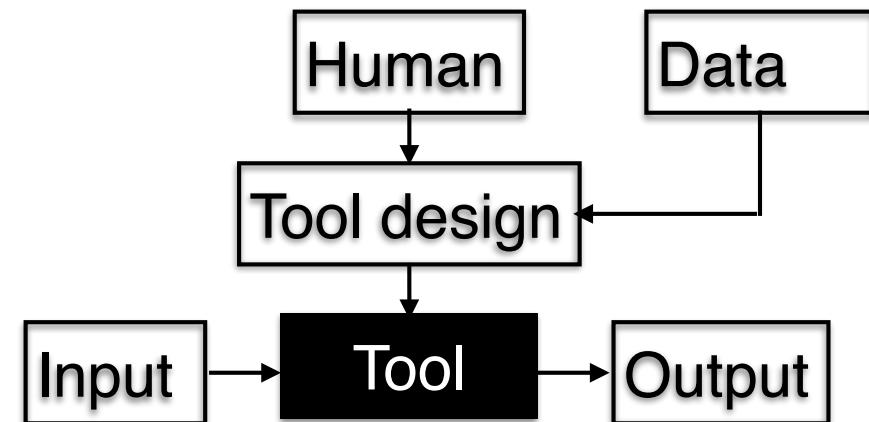
Price =  $a^*(\text{area}) + b^*(\#\text{ rooms}) + c^*(\text{distance to metro})$

Parameters: a,b,c

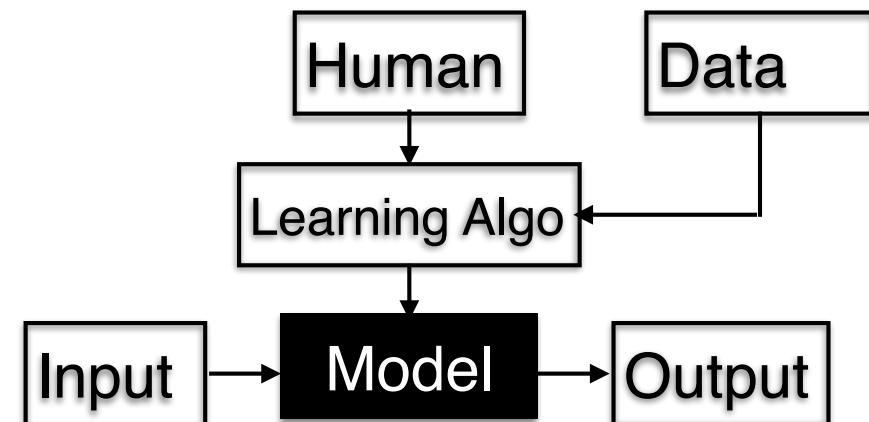
Use data to get the “**best**” parameters

# Machine Learning Tasks Revisited

Machine Learning



Machine Learning



# Outline

1.What is Machine Learning??

2.The Wonders of Machine Learning

3.Data, Models and ML Tasks

## **4. Supervised Learning**

1. Regression

2. Classification

5. Unsupervised Learning

1. Dimensionality Reduction

2. Density Estimation

# Notation

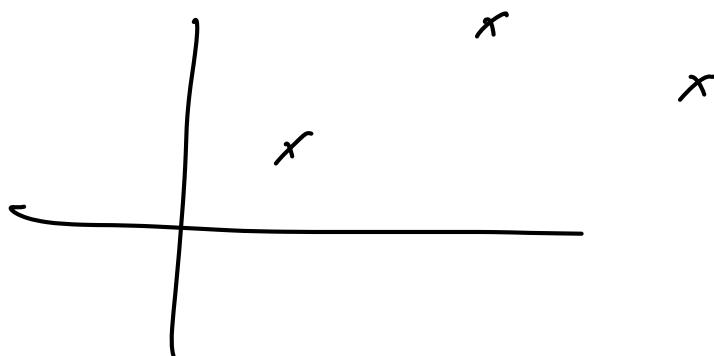
$$\begin{aligned}x^1 &= [1, 2, 3] \\x^2 &= [7, 8, 9] \quad x_2 = 8\end{aligned}$$

$$\begin{pmatrix} 1.3 \\ -7.6 \\ 5.9 \end{pmatrix} \in \mathbb{R}^3$$

- $\mathbb{R}$ : real numbers,  $\mathbb{R}_+$ : Positive reals,  $\mathbb{R}^d$ : d-dimensional vector of reals.
- $\mathbf{x}$  : vector.  $x_j$ :  $j^{\text{th}}$  co-ordinate.  $\|\mathbf{x}\|$ : Length of vector  $\mathbf{x}$ .  $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 7 \end{bmatrix}$
- $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n$ : Collection of  $n$  vectors.
- $x_j^i$  :  $j^{\text{th}}$  co-ordinate of  $i^{\text{th}}$  vector.
- $(x_1)^2$  : Square of the first co-ordinate of the vector  $\mathbf{x}$
- $1(2 \text{ is even}) = 1, 1(2 \text{ is odd}) = 0$ .

$$\|\mathbf{x}\|^2 = x_1^2 + x_2^2 + \dots + x_d^2$$

# Supervised Learning



- Supervised learning is curve-fitting.
- Given  $\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$
- Find a model  $f$  such that  $f(\mathbf{x}^i)$  is ‘close’ to  $y^i$

$f$  is a function

# Regression

- E.g. Predict house price from room, area, distance.
- Training data:  $\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$
- $\mathbf{x}^i \in \mathbb{R}^d, y^i \in \mathbb{R}$
- Algorithm outputs a model  $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- Loss  $\frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}^i) - y^i)^2 = \text{squared loss}$
- $f(\mathbf{x}) = \underbrace{\mathbf{w}^\top \mathbf{x} + b}_{\text{Linear Parameterisation}} = \sum_{j=1}^d w_j x_j + b = w_1 (\# \text{ rooms}) + w_2 (\text{area}) + w_3 (\text{distance}) + b$

# Regression Illustration 1

$d = 1$

$x$	$y$	$f$	$g$
[1]	2.1		
[2]	3.9		
[3]	6.2		
[6]	11.5		
[7]	13.9		

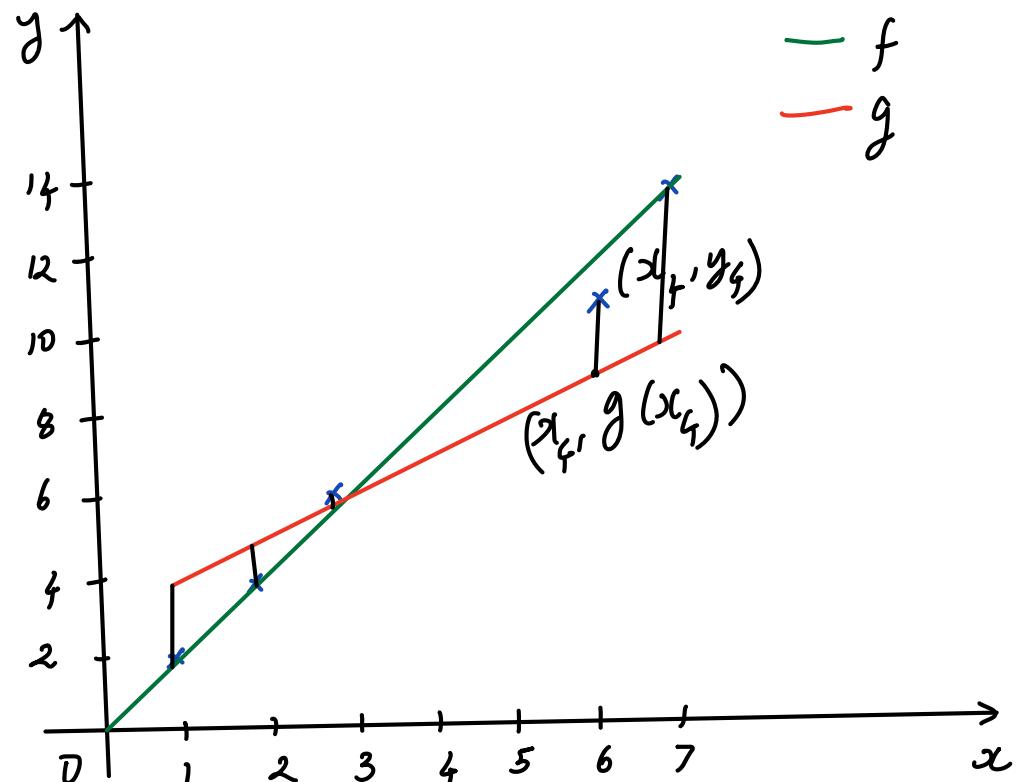
$$f(x) = 2x,$$

$$g(x) = x + 3$$

$$\text{Loss } [f] =$$

=

$$\text{Loss } [g] =$$



# Regression Illustration 2

Rooms	Area	Distance	Price
3	9	1.9	<b>5.0</b>
2	7	2.1	<b>3.2</b>
4	12	2.8	<b>6.6</b>
5	16	0.9	<b>9.8</b>
5	15	3.1	<b>8.5</b>
4	11	1.6	<b>6.9</b>

$$f = 2 * \text{Rooms} - 0.5 * \text{dist}$$

$$g = \text{Rooms} + 2 * \text{dist}$$

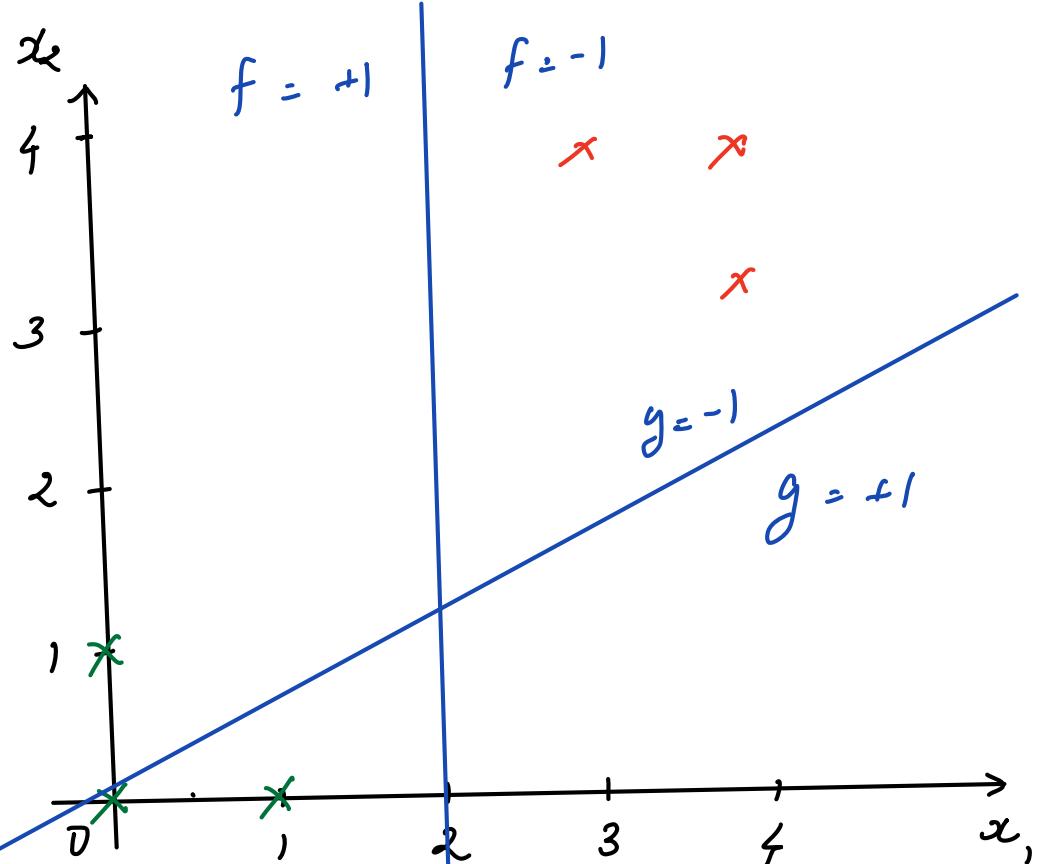
# Classification

- E.g. Predict if rooms>3 from area and price.
- Training data:  $\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$
- $\mathbf{x}^i \in \mathbb{R}^d, y^i \in \{+1, -1\}$
- Algorithm outputs a model  $f : \mathbb{R}^d \rightarrow \{+1, -1\}$
- Loss  $\stackrel{[f]}{=} \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(\mathbf{x}^i) \neq y^i)$  = Fraction of training data classified wrongly by  $f$
- $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$

Linear separator

# Classification Illustration 1

$x$	$y$	$f$	$g$
$[0, 0]$	+1		
$[0, 1]$	+1		
$[1, 0]$	+1		
$[4, 4]$	-1		
$[3, 4]$	-1		
$[4, 3]$	-1		



$$f(x) = \text{Sign}(x_1 - x_2)$$

$$g(x) = \text{Sign}(x_1 - 2x_2)$$

$$\text{Loss}[f] = \frac{1}{6}(0) = 0$$

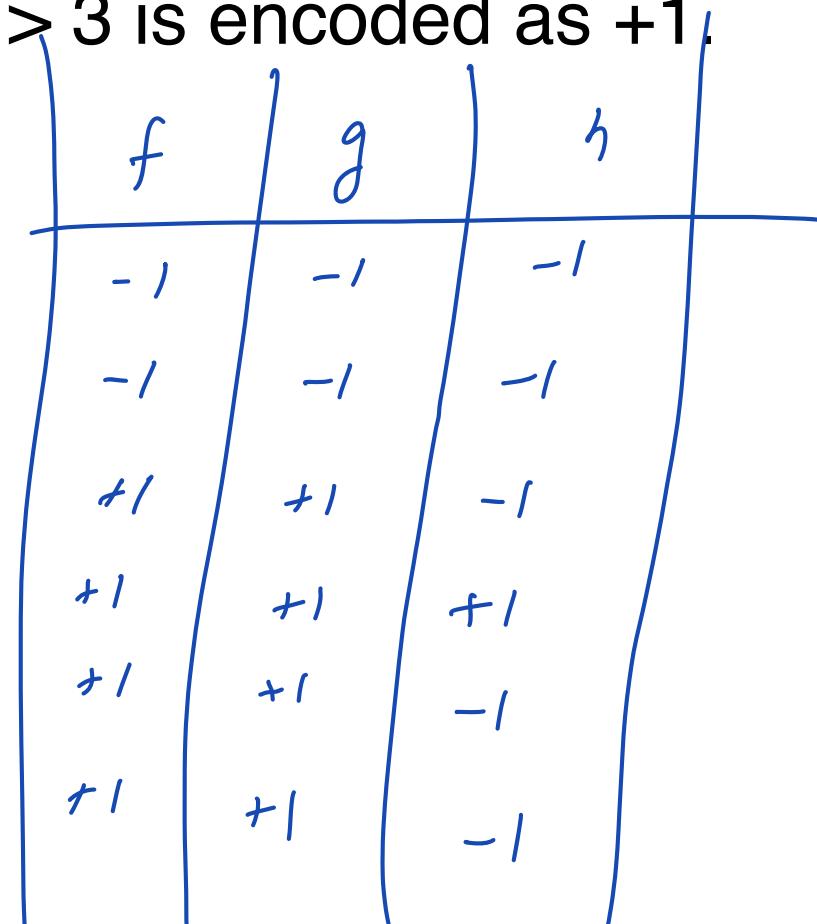
$$\text{Loss}[g] = \frac{1}{6}(1) = \frac{1}{6}$$

# Classification Illustration 2

Area	Price	Rooms
------	-------	-------

9	5.0	-1
7	3.1	-1
12	6.9	+1
16	9.7	+1
15	8.5	+1
11	7.1	+1

Rooms=1 or 2 or 3 is encoded as -1.  
Rooms > 3 is encoded as +1.



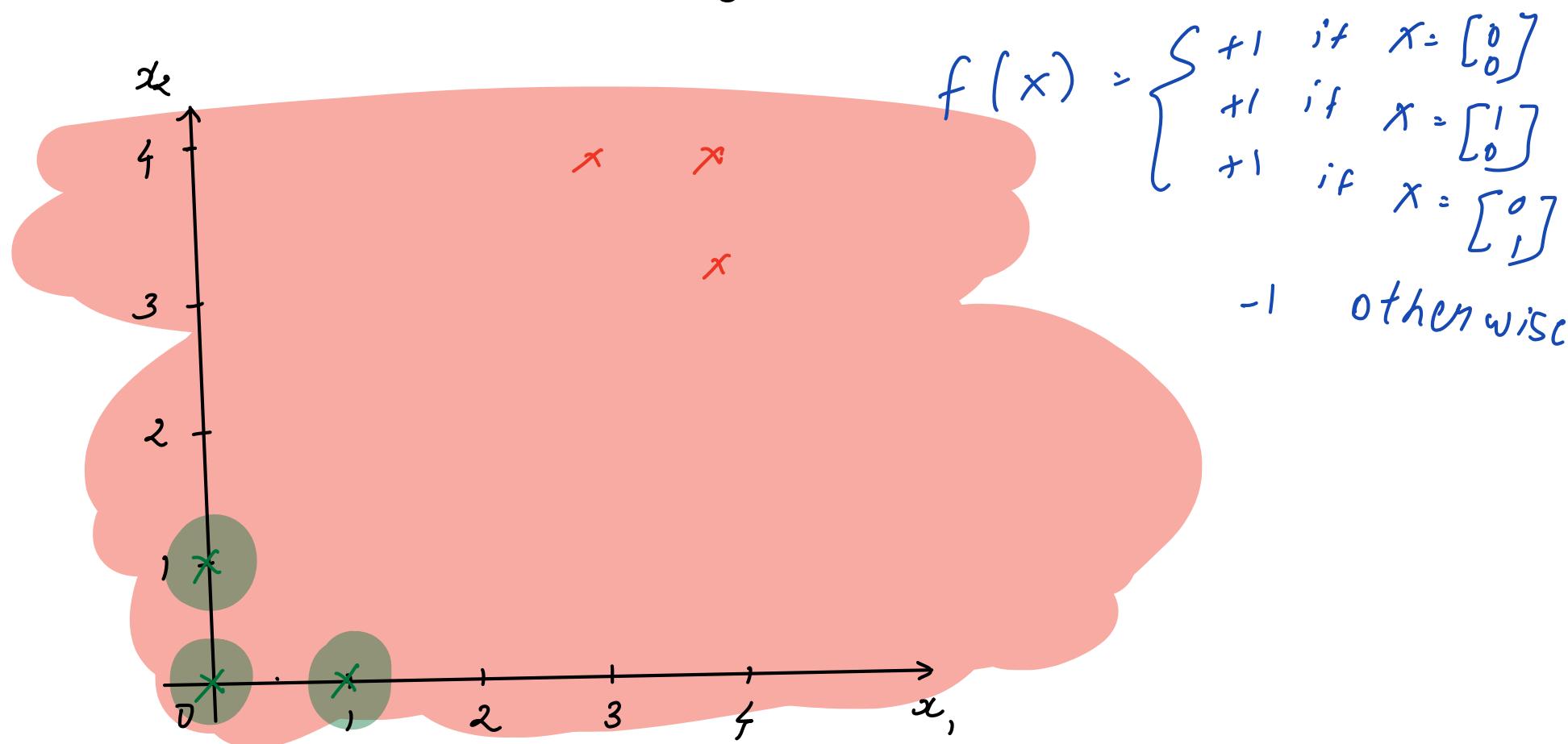
$$f(x) = \text{sign}(\text{area} - 10)$$

$$g(x) = \text{sign}(\text{price} - 6)$$

$$h(x) = \text{sign}(\text{price} - 9)$$

# Evaluating Learned Models : Test Data

- Learning algorithm uses training data  $(x^1, y^1), \dots, (x^n, y^n)$  to get model  $f$ .
- But evaluating the learned model must **not** be done on the training data itself.
- Use test data that is **not** in the training data for model evaluation.



# Model Selection : Validation Data

- Learning algorithms just find the “best” model in the collection of models given by the human.
- How to find the right collection of models?
- This is called model selection, and it is done by using another subset of data called **validation data** that is distinct from train and test data.

$$\text{Price} = w_1 * (\# \text{rooms}) + w_2 (\text{area}) + w_3 (\text{distance}) + b$$

# Outline

- 1.What is Machine Learning??
- 2.The Wonders of Machine Learning
- 3.Data, Models and ML Tasks
- 4.Supervised Learning
  1. Regression
  2. Classification
- 5. Unsupervised Learning**
  1. Dimensionality Reduction
  2. Density Estimation

# Unsupervised Learning

- Unsupervised learning is ‘understanding data’
- Data:  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$
- $\mathbf{x}^i \in \mathbb{R}^d$
- Build models that compress, explain and group data.

# Unsupervised Learning Application

Tweet 1



⋮

Tweet 999999



Group the million tweets into 10 manageable groups

# Dimensionality Reduction

$$10^4 \times 10^6 \rightarrow 10^6 \times 100$$

E.g.: Represent a million gene expression levels of a million people, using just 100 numbers per person.

Dimensionality reduction: compression and simplification.

# Dimensionality Reduction

- Data:  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$
- $\mathbf{x}^i \in \mathbb{R}^d$
- Encoder  $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$
- Decoder  $g : \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$
- Goal :  $g(f(\mathbf{x}^i)) \approx \mathbf{x}^i$
- Loss =  $\frac{1}{n} \sum_{i=1}^n \|g(f(\mathbf{x}^i)) - \mathbf{x}^i\|^2$

$$f : \mathbb{R}^d \xrightarrow{Wx + b} \mathbb{R}^{d'}$$
$$g : \mathbb{R}^{d'} \xrightarrow{Vu + c} \mathbb{R}^d$$

Handwritten notes:

- $d' \ll d$
- $d' \times d'$
- $d \times d$

# Dimensionality Reduction Illustration

$$d=2, d'=1, n=4$$

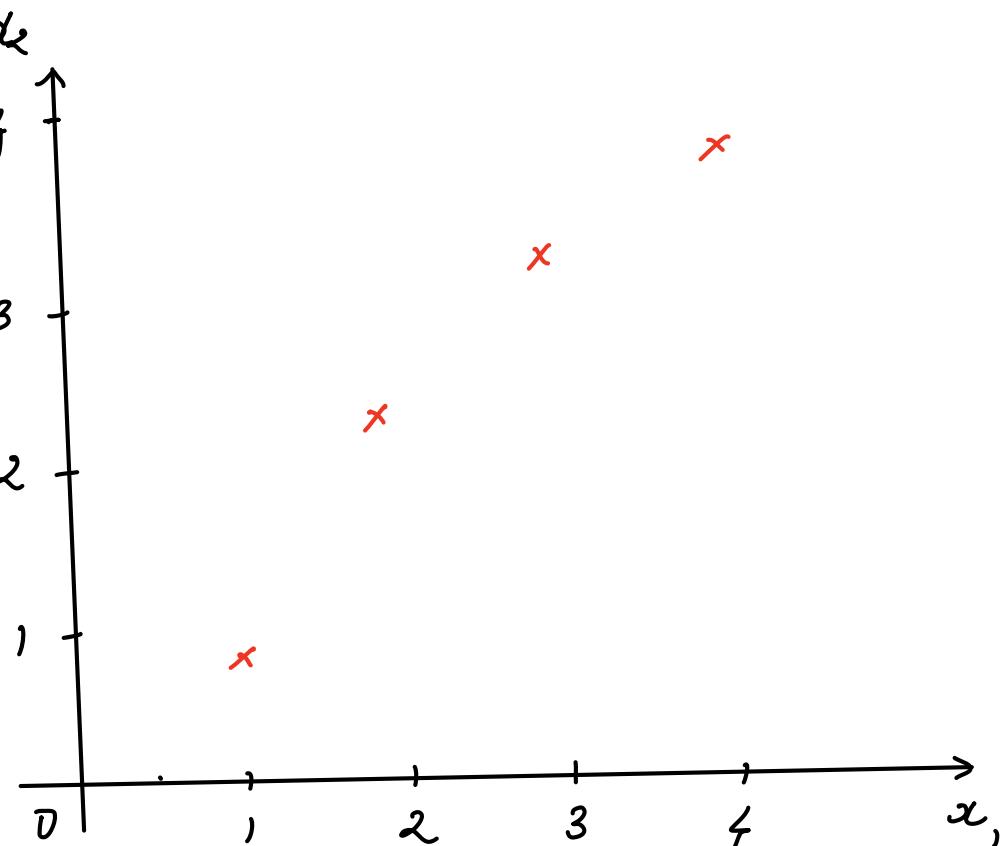
$[1, 0.8]$	$f$ 0.2	$g$ $[0.2, 0.2]$
$[2, 2.2]$	-0.2	$[-0.2, -0.2]$
$[3, 3.2]$	-0.2	$[-0.2, -0.2]$
$[4, 3.8]$	0.2	$[0.2, 0.2]$

$$f(x) = x_1 - x_2$$

$$g(u) = [u, u]$$

$$\tilde{f}(x) = \frac{x_1 + x_2}{2}$$

$$\tilde{g}(u) = [u, u]$$



# Dimensionality Reduction Illustration

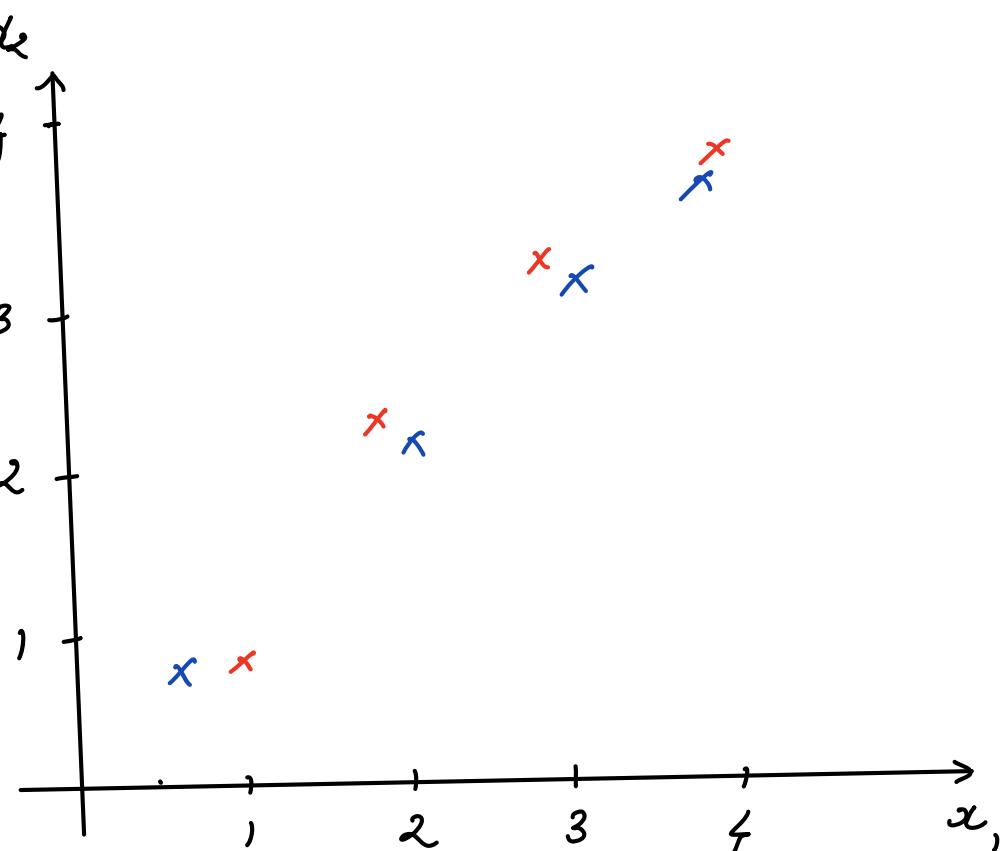
$d=2$	$d'=1$	$n=4$
$[1, 0.8]$	$\tilde{f}$	$[0.9, 0.9]$
$[2, 2.2]$	$0.9$	$[2.1, 2.1]$
$[3, 3.2]$	$2.1$	$[3.1, 3.1]$
$[4, 3.8]$	$3.1$	$[3.9, 3.9]$
	$3.9$	

$$f(x) = x_1 - x_2$$

$$g(u) = [u, u]$$

$$\tilde{f}(x) = \frac{x_1 + x_2}{2}$$

$$\tilde{g}(u) = [u, u]$$



$10^6$   
 $\tilde{f}, \tilde{g}$   
 $10 \times 100$

# Density Estimation

E.g.: Assuming tweets from an account are independently generated randomly. Create a robot account that generates more such tweets.

$$f(\text{Tweet}) = \frac{\text{Score of the tweet}}{\text{Chopra's Twitter stream}}$$

wisdomofchopra.com

It has been said by some that the thoughts and tweets of Mr. Chopra are indistinguishable from a set of profound sounding words put together in a random order, particularly the tweets tagged with "#cosmisconsciousness". This site aims to test that claim! Each "quote" is generated from a list of words that can be found in Chopra's Twitter stream randomly stuck together in a sentence.

"A formless void transforms total mysteries"

[RECEIVE MORE WISDOM...](#)

 [Tweet the wisdom](#)

Disclaimer: This is intended for entertainment purposes only. It in no way reflects the thoughts of any real person.

# Density Estimation

"A formless void transforms total mysteries"

[RECEIVE MORE WISDOM...](#)

 [Tweet the wisdom](#)

To generate such sentences randomly, we need to be able to assign a probability score to every possible 128 character sentence, giving high scores to those that are likely to be from the original source.

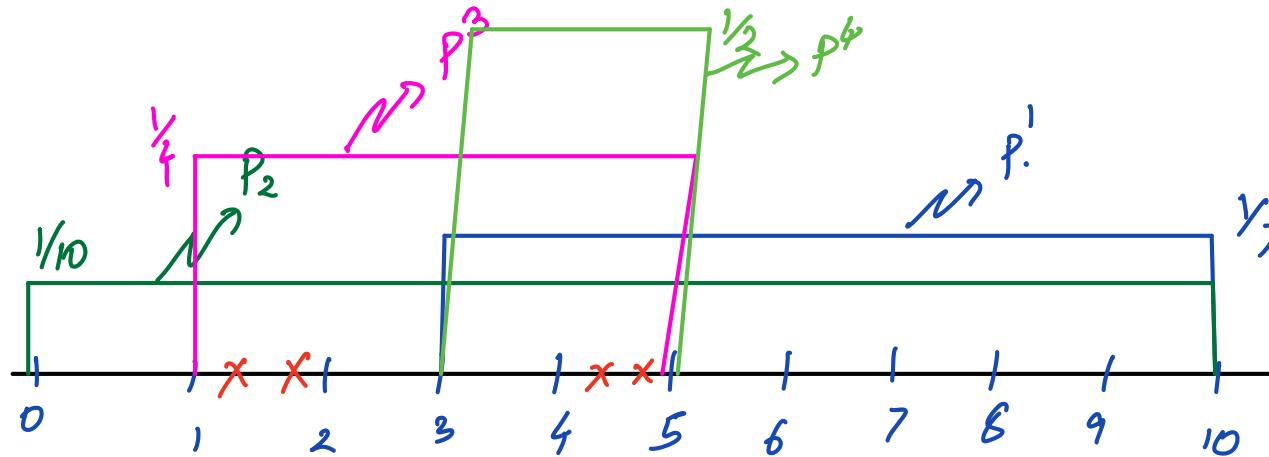
A density estimation model takes in several samples from a random source, and outputs a model that assigns a probability score to every possible instance.

# Density Estimation

- Data:  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n\}$
- $\mathbf{x}^i \in \mathbb{R}^d$
- Probability mapping  $P : \mathbb{R}^d \rightarrow \mathbb{R}_+$  that ‘sums’ to one.
- Goal :  $P(\mathbf{x})$  is large if  $\mathbf{x} \in \text{Data}$ , and low otherwise.
- Loss =  $\frac{1}{n} \sum_{i=1}^n -\log(P(\mathbf{x}^i))$   $P(\mathbf{x}^i)$  is large.

$$P(\text{anything}) = 10^{10}$$

# Density Estimation Illustration 1



$$\begin{aligned}x^1 &= [1.2] \\x^2 &= [1.9] \\x^3 &= [4.3] \\x^4 &= [4.8]\end{aligned}$$

$P^1 = \text{Uniform in } [3, 10]$   $0, 0, \frac{1}{7}, \frac{1}{7}$

$P^2 = \text{Uniform in } [0, 10]$   $\frac{1}{10}, \frac{1}{10}, \frac{1}{10}, \frac{1}{10}$

$P^3 = \text{Uniform in } [1, 5]$   $\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}$

$P^4 = \text{Uniform in } [3, 5]$   $0, 0, \frac{1}{2}, \frac{1}{2}$

$\text{loss}[P^4] = \text{loss}[P^1] = \infty > \text{loss}[P^2] > \text{loss}[P^3]$

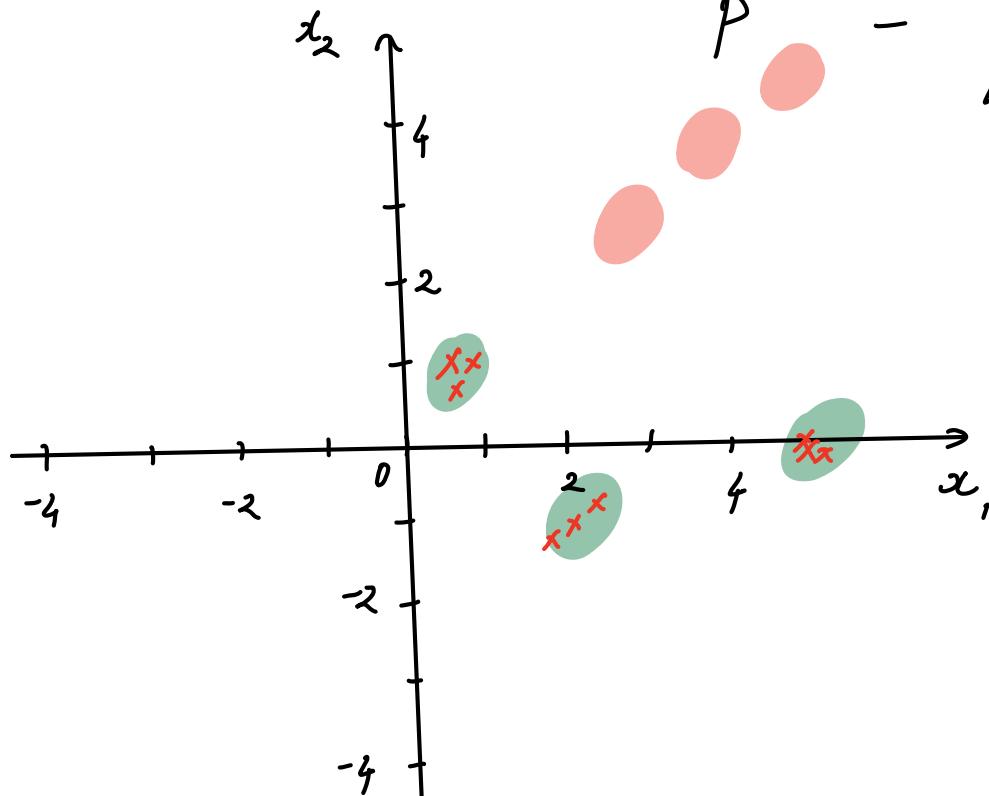
# Density Estimation Illustration 2

$$d=2$$

- $(1.1, 1.3)$
- $(0.9, 0.7)$
- $(2.1, -1)$
- $(5.1, 0.1)$
- $(2.2, -0.9)$
- $(5.1, 0.0)$
- $(0.9, 1.2)$
- $(1.9, -1.1)$
- $(4.8, -0.1)$

Gaussian Mixture Model

$$\begin{aligned} P^1 &= \begin{bmatrix} 3 \\ 3 \end{bmatrix}, \begin{bmatrix} 5 \\ 5 \end{bmatrix}, \begin{bmatrix} 4 \\ 4 \end{bmatrix} \\ P^2 &= \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \begin{bmatrix} 5 \\ 0 \end{bmatrix} \end{aligned}$$



Clustering

# Density Estimation Illustration 2