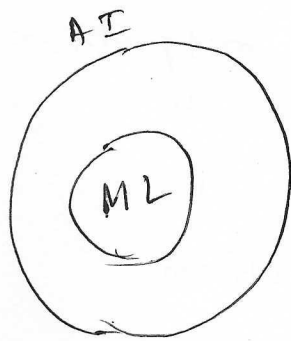


①

Introduction

AI \Rightarrow The broad goal

Building machines that can perform tasks that normally require human intelligence.

Ex: Reasoning & problem solving, understanding language, acting in an environment (robots, agents).

ML \Rightarrow A subset of AI.

Techniques that allow machines to learn patterns from data instead of being explicitly programmed.

Ex: Learning to classify images, translating languages, recommender systems.

[All ML is AI, but not all AI is ML]

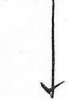
Before ML dominated:

- AI systems were rule based.

- Ex: medical expert systems

If symptoms A AND B \rightarrow disease X

- Problems: hard to scale, fragile, required human experts to encode knowledge.



ML shift: Instead of rules, learn from examples.

Ex: Given 1 million labelled cases, learn disease prediction.

What is learning?

- The acquisition of knowledge or skills through experience, study, or being taught.

What is ML?

- Field of study that gives computers the ability to learn without being explicitly programmed.

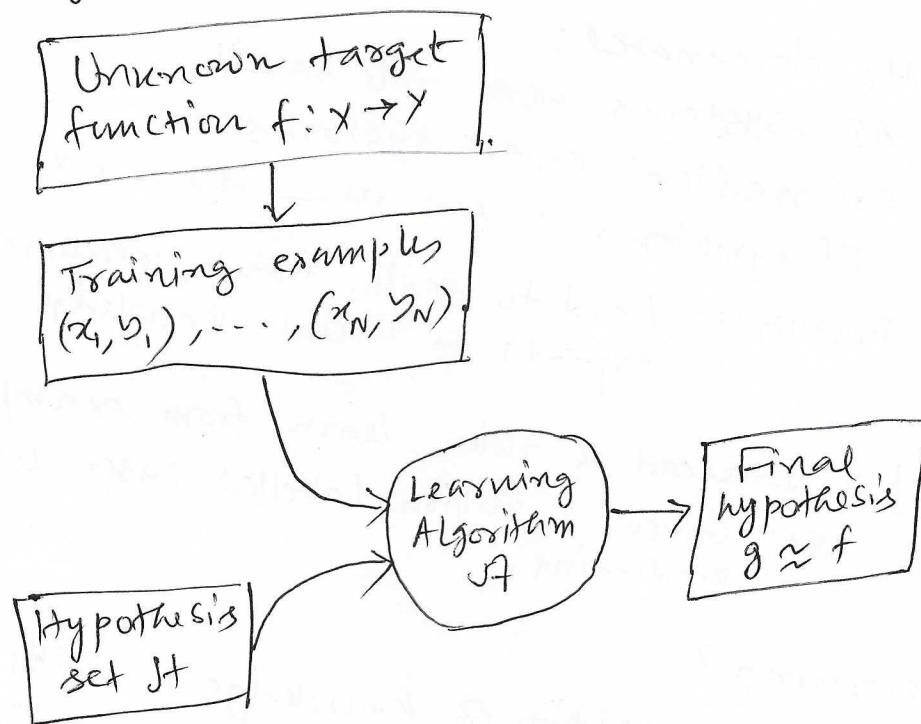
Types of Learning

- Supervised: training data includes desired outputs.
- unsupervised: " " does not include " "
- weakly/semi supervised: " " includes a few " "
- Reinforcement: Rewards from sequence of actions.

Supervised Learning

Given examples of a function $(x, f(x))$, predict the function $f(x)$ for new examples. The function can be:

- Discrete \rightarrow classification (e.g. spam/not spam)
- Continuous \rightarrow regression (housing price, weather)
- probability \rightarrow probability estimation (e.g. patient symptom)



Training vs Testing

What do we want?

- good performance on training data?
- No, good performance on unseen test data.

Training data \rightarrow given to us for learning the function f

Testing data \rightarrow used to see if you have learnt anything

Usually a dataset is split into train & test
 \downarrow \rightarrow
annotate Check performance.

②

Probability Review.

Consider a non-deterministic event A (boolean variable)

What does $P(A)$ mean?

Frequentist view: Limiting frequency of a repeating non-deterministic event.

$$\lim_{N \rightarrow \infty} \frac{\# A \text{ is true}}{N}$$

Bayesian view: $P(A)$ is your "belief" about A .

Axioms of probability.

- $0 \leq P(A) \leq 1$
- $P(\text{empty set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Sample space \rightarrow space of events

Random Variables \rightarrow Mapping from events to numbers (set of possible values from a random experiment)

discrete \swarrow continuous

Probability distribution \rightarrow List all possible outcomes of a random variable along with their corresponding values. ex: $\{H, T\}$, $\{1:1/6, 2:1/6, \dots, 6:1/6\}$

Probability Mass Function (PMF) \rightarrow Probability that a discrete random variable equals a specific value.

$$P(X=x) = P(x)$$

Probability Density Function (PDF) \rightarrow How dense probability is around values (in the continuous case, probability at an exact point is zero).

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

Expectations → Long run average outcome (over many repetitions). weighted average, where weights are probabilities.

$$E_p[f(x)] = \sum_{x=1}^k p(x) f(x) \quad (\text{discrete})$$
$$= \int_{-\infty}^{\infty} p(x) f(x) dx \quad (\text{continuous})$$

Joint distribution (in discrete space) & Marginalization:

Ex: Suppose we observe 2 random variables.

W = weather: Rainy, Sunny

U = Whether a person carries an umbrella: Yes, No.

Qn: What is the probability of a specific combination of weather & umbrella choice?

ex:- rainy AND carries an umbrella

- sunny AND does NOT carry an umbrella

let's say after observing many days, we get:

W	U	P(W, U)
Rainy	Yes	0.4
Rainy	No	0.1
Sunny	Yes	0.1
Sunny	No	0.1

← This table is the joint distribution $P(W, U)$

Marginalization

Suppose we ask what is the probability that it is rainy, regardless of the umbrellas?

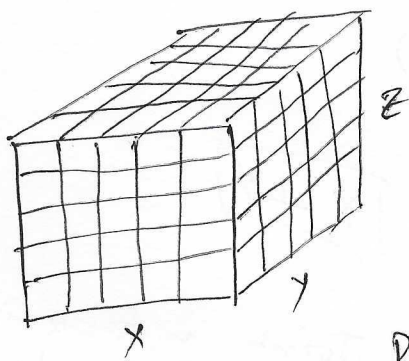
Then we marginalize out (ignore) the umbrella variable.

$$P(W = \text{Rainy}) = P(\text{Rainy, Yes}) + P(\text{Rainy, No}) = 0.4 + 0.1 = 0.5$$

Similar for $P(W = \text{sunny})$ as well.

This is called marginal distribution of weather.

②



$$p(x, y) \downarrow \text{marginalize out } y$$

$$p(x=x) = \sum_y p(x=x, y=y)$$

↑
Sum of all values where $x=x$ occurs with all possible values of y

$$p(x, y, z) \downarrow \text{we want to compute } p(x)$$

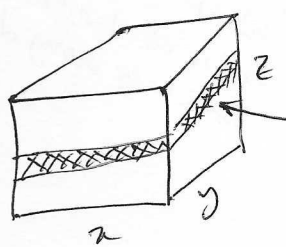
$$p(x, y) = \sum_z p(x, y, z)$$

$$p(x) = \sum_y p(x, y)$$

If there are n variables, we want marginalization of everything other than a single variable, we need $(n-1)$ summation.

Conditional probabilities:

$p(x=x | y=y)$: What do you believe about $y=y$, if I tell you $x=x$?



$$p(x, y | z) \text{ (conditioning on a variable is taking that slice out)}$$

$$p(x, y | z=z) = \frac{p(x, y, z)}{p(z)} \leftarrow \frac{\text{Joint}}{\text{Marginal}}$$

Chain rule

$$p(y=y | x=x) = \frac{p(y=y, x=x)}{p(x=x)}$$

$$\Rightarrow p(y=y, x=x) = p(y=y | x=x) p(x=x)$$

For d -dimensional case of joint distribution, we need to recursively apply the chain rule:

$$p(x_1=x_1, \dots, x_d=x_d) = p(x_2=x_2, \dots, x_d=x_d | x_1=x_1) p(x_1=x_1)$$

$$= p(x_3=x_3, \dots, x_d=x_d | x_2=x_2, x_1=x_1)$$

$$p(x_2=x_2 | x_1=x_1) p(x_1=x_1)$$

$$= \prod_{j=1}^d p(x_j=x_j | x_1=x_1, \dots, x_{j-1}=x_{j-1})$$

Conditional independence

$$P(Y=y, X=x) = P(Y=y | X=x) P(X=x)$$

India wins world cup

today is sunny

$$X \perp Y$$

$$= P(Y=y) P(X=x)$$

Joint is factorized into marginal.

Bayes's rule

from chain rule,

$$P(Y=y, X=x) = P(Y=y | X=x) P(X=x)$$

↓ swap

$$P(X=x, Y=y) = P(X=x | Y=y) P(Y=y)$$

$$\Rightarrow P(Y=y | X=x) = \frac{P(X=x | Y=y) P(Y=y)}{P(X=x)}$$

posterior
(what do you believe of the classes after you see the data)

Likelihood
(how much does a certain hypothesis explain the data)

Normalization

Prior
(what do you believe before seeing any data)

Entropy: Measures amount of uncertainty in a distribution.

$$H(X) = - \sum_{x \in X} p(x) \log p(x)$$



KL divergence: An asymmetric measure of distance between two distributions.

$$KL[P || Q] = \sum_x P(x) [\log P(x) - \log Q(x)]$$

$KL > 0$ unless $P=Q$ (in that case $KL=0$)