Artificial Intelligence

A Course About Foundations



Probabilistic Reasoning: Representation & Inference

Marco Piastra

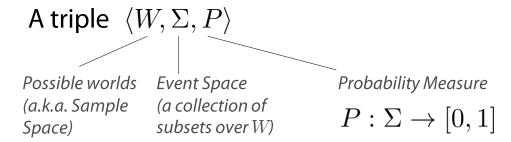
Artificial Intelligence 2024–2025 Probabilistic Reasoning [1]

Probability Space

Artificial Intelligence 2024–2025 Probabilistic Reasoning [2]

Probability Space (preliminary definition)

Probability space



The intuitive definition is simple enough, its mathematical translation ... not so much

Artificial Intelligence 2024–2025 Probabilistic Reasoning [3]

Event Space: a collection of subsets of possible worlds

Boolean algebra

A non-empty collection of subsets Σ of a set W such that:

- 1) $A, B \in \Sigma \implies A \cup B \in \Sigma$
- 2) $A \in \Sigma \implies A^c \in \Sigma$
- 3) $\varnothing \in \Sigma$

Corollary:

The sets \varnothing e W belong to any Boolean algebra generated on W Σ is also closed under <u>binary</u> intersection

• σ -algebra

A non-empty collection of subsets Σ of a set W such that:

- 1) $A_k \in \Sigma, \ \forall k \in \mathbb{N}^+ \implies (\bigcup_{k=1}^{\infty} A_k) \in \Sigma$
- 2) $A \in \Sigma \implies A^c \in \Sigma$
- 3) $\varnothing \in \Sigma$

Corollary:

This is a stronger requirement: closeness under <u>countable</u> union Hence a σ -algebra is a boolean algebra but not vice-versa

The sets \varnothing and W belong to any σ - algebra generated on W Σ is also closed under *countable intersection*

Probability Measure

• σ -algebra (*Event Space*)

A non-empty collection of subsets Σ of a set W such that:

- 1) $A_k \in \Sigma, \ \forall k \in \mathbb{N}^+ \implies (\bigcup_{k=1}^{\infty} A_k) \in \Sigma$
- 2) $A \in \Sigma \implies A^c \in \Sigma$
- 3) $\varnothing \in \Sigma$
- Probability <u>measure</u> over a σ -algebra (i.e., over the events)

A function $P:\Sigma \to [0,1]$

i.e. P assigns a measure (i.e. a real number) to each elements of a σ -algebra Σ of subsets of W

- 1) $\forall A \in \Sigma, P(A) \geq 0$
- 2) $A_1, A_2 \in \Sigma$ are $\underline{disjoint} \implies P(A_1 \cup A_2) = P(A_1) + P(A_2)$ Finite additivity $A_k \in \Sigma, \ \forall k \in \mathbb{N}^+$ are all $\underline{disjoint} \implies P\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} P(A_k)$ Countably infinite additivity

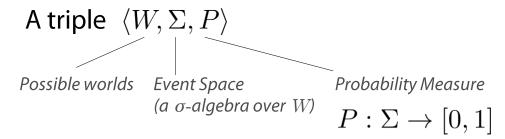
3) $P(\varnothing) = 0$

4) $P(A^c) = 1 - P(A)$ (which implies P(W) = 1)

Artificial Intelligence 2024–2025 Probabilistic Reasoning [5]

Probability Space

Probability space



Why bothering so much with these (very) technical definitions?

Rationale (just a few hints)

Closure w.r.t. countable unions of a σ -algebra (as well as countable additivity of P) is required for dealing with <u>infinite sequences</u> of events

A σ -algebra is <u>included</u> in the *power set* of W (i.e., the collection of all its possible subsets): requiring closure on <u>countable</u> union is a <u>restriction</u>, to ensure <u>measurability</u>

(see the so-called Banach-Tarski paradox for counterexamples)

Artificial Intelligence 2024–2025 Probabilistic Reasoning [6]

An Aside: Probability is Systemic

In general

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

It follows from the additivity property

If $A \cap B = \emptyset$ then events A and B are <u>disjoint</u>

$$P(A \cup B) = P(A) + P(B)$$

(*) Note that $A\cap B=\varnothing\implies P(A\cap B)=0$ but not vice-versa: as an event can have zero probability without being empty

(**) Unlike in propositional logic, knowing P(A) and P(B) is <u>not</u> sufficient for determining $P(A \cup B)$

Namely, probability is not compositional ...

Discrete Probability

Artificial Intelligence 2024–2025 Probabilistic Reasoning [8]

Studying basic properties*: (*) in a finitary setting

A simpler setting that allows a more intuitive definition of fundamental properties

Basic assumption: the set of possible worlds $\,W\,$ is $\,$ finite $\,$

Finite event space

 Σ is a *finite* collection of subsets

In such setting boolean algebra $\equiv \sigma$ -algebra

Events could also be defined via propositional logic (à la de Finetti, 1937)

Finitely additive probability measure

Just summations, no integrals

Computability will be always guaranteed

Artificial Intelligence 2024–2025 Probabilistic Reasoning [9]

Random Variables*

(*) In a finitary setting

Artificial Intelligence 2024–2025 Probabilistic Reasoning [10]

Partitions, random variables*

Partition

A <u>finite</u> collection A_i of <u>disjoint</u> subsets (i.e. <u>events</u>) such that

$$\bigcup_{i} A_i = W$$

A σ -algebra can be generated from a *partition* by taking its closure under *union* and *complement*

Artificial Intelligence 2024–2025 Probabilistic Reasoning [11]

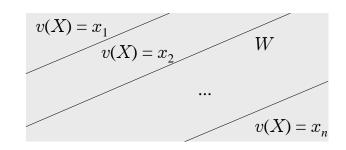
Partitions, random variables*

Random Variable (i.e. a convenient way to define a σ -algebra)

Let X be a variable having a <u>finite</u> set of possible values $\{x_1, x_2, \dots, x_n\}$ In each possible world, the variable X is assigned a specific value x_i

- The set of possible assignments $\{X = x_1, X = x_2, ... X = x_n\}$ defines a <u>partition</u> over W
- A σ -algebra can be obtained by taking the closure of the partition under union and complement
- $X = x_i$ defines an <u>event</u> (i.e. a subset of W)
- $X=x_i$ and $X=x_j$ are <u>disjoint events</u>, whenever $i \neq j$ $P(X=x_i \cup X=x_j) = P(X=x_i) + P(X=x_j)$

Random variables having binary values are also said to be <u>binomial</u> (also Bernoullian) Random variables with multiple values are also said to be <u>multinomial</u>



Random variables, joint distribution*

(*) In a finitary setting

Multiple random variables

In practice, in a probabilistic representation, there will be multiple random variables

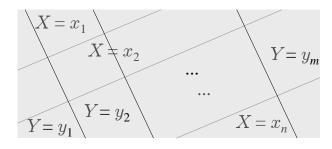
Example:

 X_i occurrence of a word i in the body of an email (binomial)

Y classification of that email as *spam* (binomial)

The intersection of two or more σ -algebras is a σ -algebra

Together, a collection of random variables defines a partition of $\,W\,$



Joint Probability Distribution

for a given set of random variables, e.g. X, Y, Z

It is a <u>function</u> that associates a value in [0, 1] to each individual combination of values

$$P(X = x, Y = y, Z = z)$$

Given that *X*, *Y* and *Z* define each a *partition* over *W*:

$$\sum_{x} \sum_{y} \sum_{z} P(X = x, Y = y, Z = z) = 1$$

Random variables: notation*

(*) In a finitary setting

• Random variables, events and σ -algebras

Sometimes the notation can be ambiguous

Examples:

This is the probability measure over the σ -algebra generated by the random variable X

$$P(X=x)$$

This the probability (i.e. a value in [0,1]) associated to the <u>event</u> X=x

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

This is the probability measure over the $\underline{\sigma}$ -algebra generated by the random variable X in the subspace of W corresponding to the event Y=y

Fundamental Operations*

(*) In a finitary setting

Artificial Intelligence 2024–2025 Probabilistic Reasoning [15]

Marginalization*

(*) In a finitary setting

Removing a random variable from a joint distribution

Given a joint probability distribution

The <u>marginal probability</u> P(X) is obtained via a summation:

$$P(X) := \sum_{y \in \mathcal{Y}} P(X, Y = y)$$

A marginal probability can be a joint probability as well ...

$$P(X,Y) := \sum_{z \in \mathcal{Z}} P(X,Y,Z=z)$$

Marginal probability, shorthand notation with values of Y omitted:

$$P(X) = \sum_{Y} P(X,Y)$$
Shorthand notation

Artificial Intelligence 2024–2025 Probabilistic Reasoning [16]

Conditionalization*

(*) In a finitary setting

Exploring 'what if' something becomes known

Given a joint probability distribution

The <u>conditional probability</u> $P(X \mid Y = y)$ is defined as:

$$P(X \mid Y = y) := \frac{P(X, Y = y)}{P(Y = y)}$$

A conditional probability can be a joint probability as well ...

$$P(X, Y \mid Z = z) := \frac{P(X, Y, Z = z)}{P(Z = z)}$$

Conditional probability, more general notation:

$$P(X\mid Y) \ := \ \frac{P(X,Y)}{P(Y)} \qquad \qquad \text{Conditional probabilities for the } \underbrace{\text{whole } \sigma\text{-algebra} \text{ generated by } Y}$$

(it represents a family of probability measures)

Conditionalization*

(*) In a finitary setting

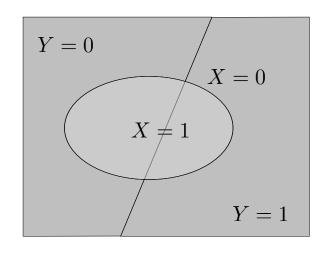
Exploring 'what if' something becomes known Conditional probability, more general notation:

$$P(X \mid Y) := \frac{P(X,Y)}{P(Y)}$$

Assume both variables are binary $X,Y\in\{0,1\}$

$$P(X \mid Y = 0) := \frac{P(X, Y = 0)}{P(Y = 0)}$$

$$P(X \mid Y = 1) := \frac{P(X, Y = 1)}{P(Y = 1)}$$



Each value of the conditioning variable defines a distinct event (sub)space

Chain Rule

Starting point: a joint distribution

conditioning on X:

$$P(Y,Z|X) = \frac{P(X,Y,Z)}{P(X)}$$

which implies:

$$P(X, Y, Z) = P(X) P(Y, Z|X)$$

then, conditioning on Y:

$$P(Z|X,Y) = \frac{P(Y,Z|X)}{P(Y|X)}$$

which implies:

$$P(Y,Z|X) = P(Y|X) P(Z|X,Y)$$

hence and finally:

$$P(X, Y, Z) = P(X) P(Y|X) P(Z|X, Y)$$

<u>Univariate factorization</u> of the joint distribution

Any other sequence of conditionalization would do

Inference (without *learning*)

Artificial Intelligence 2024–2025 Probabilistic Reasoning [20]

Probabilistic Inference* (general structure)

(*) In a finitary setting

General structure of probabilistic inference problems

The starting point is a fully-specified joint probability distribution

$$P(X_1, X_2, \ldots, X_n)$$

In an *inference* problem, the set of random variables $\{X_1, X_2, \dots, X_n\}$ is divided into three categories:

- 1) Observed variables $\{X_o\}$ having a definite (and supposedly known) value
- 2) Irrelevant variables $\{X_i\}$ which are neither observed nor directly part of the answer
- 3) Relevant variables $\{X_r\}$ which are part of the answer we seek

In general, the problem is finding:

$$P(\{X_r\}|\{X_o\}) = \sum_{\{X_i\}} P(\{X_r\}, \{X_i\}|\{X_o\})$$

- "Decidability" ("computability") is guaranteed (* in a finitary setting)
 Given that the joint probability distribution is completely specified
- Computational efficiency can be a problem

The number of value combinations to be considered in the summation grows exponentially with the number of random variables in $\{X_r\} \cup \{X_i\}$

Bayes' Theorem* (T. Bayes, 1764)

(*) In a finitary setting

Definition

A relation between conditional and marginal probabilities

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

P(Y|X) is also called the *likelihood* L(X|Y)



The theorem follows from the definition of conditional probability (chain rule)

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

Furthermore, given the definition of marginalization:

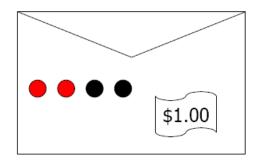
$$P(Y) = \sum_{X} P(X,Y) = \sum_{X} P(Y|X)P(X) \qquad \qquad \qquad \text{Also called 'law of total probability'}$$

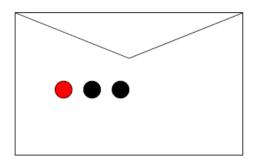
it follows an alternative formulation of the Bayes' theorem:

$$P(X|Y) = \frac{P(Y|X)P(X)}{\sum_{X} P(Y|X)P(X)}$$

Artificial Intelligence 2024–2025 Probabilistic Reasoning [22]

Example: information and bets





Two envelopes, only one is extracted

One envelope contains two red tokens and two black tokens, it is worth \$1.00 One envelope contains one red token and two black tokens, it is valueless

The envelope has been extracted.

Before posing you bet, you are allowed to extract on token from it

- a) The token is black. How much do you bet?
- b) The token is red. How much do you bet?

Purpose: showing that Bayes' Theorem makes the representation easier

Artificial Intelligence 2024–2025 Probabilistic Reasoning [23]

Independence

Artificial Intelligence 2024–2025 Probabilistic Reasoning [24]

Independence, conditional independence

■ **Independence** (also marginal independence)

Two variables are <u>independent</u> iff their joint probability can be factorized into the product of *marginals*

$$< X \perp Y> \qquad \Leftrightarrow \qquad P(X,Y) = P(X)P(Y) \qquad \Leftrightarrow \qquad < Y \perp X>$$
 $\Rightarrow \qquad P(X|Y) = \frac{P(X,Y)}{P(Y)} = \frac{P(X)P(Y)}{P(Y)} = P(X)$

Conditional independence

Two variables are <u>conditionally independent</u> given a third variable, iff their joint conditional probability can be factorized into the product of *conditional marginals*

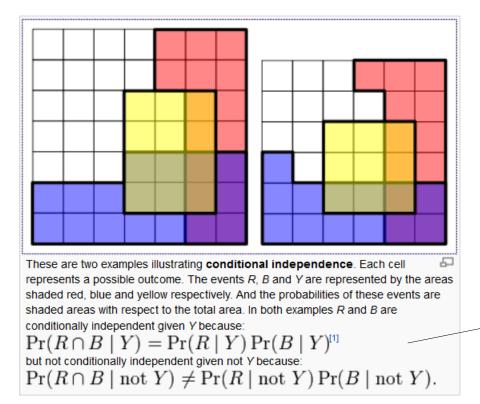
$$< X \perp Y \mid Z> \qquad \Leftrightarrow \qquad P(X,Y|Z) = P(X|Z)P(Y|Z) \qquad \Leftrightarrow \qquad < Y \perp X \mid Z>$$

$$\Rightarrow \qquad P(X|Y,Z) = \frac{P(X,Y|Z)}{P(Y|Z)} = \frac{P(X|Z)P(Y|Z)}{P(Y|Z)} = P(X|Z)$$

CAUTION: the two forms of <u>independence</u> are distinct!

$$< X \perp Y> \implies < X \perp Y \mid Z>$$
 $< X \perp Y \mid Z> \implies < X \perp Y>$

Independence, conditional independence



[from Wikipedia, "Conditional Independence"]

R, *B* and *Y* here are subsets, i.e. <u>events</u>, not random variables

The example above shows that (marginal or conditional) independence of two specific <u>events</u> does NOT imply (marginal or conditional) independence of the whole σ -algebras

Artificial Intelligence 2024–2025 Probabilistic Reasoning [26]

Continuous Random Variables

Artificial Intelligence 2024–2025 Probabilistic Reasoning [27]

Continuous random variables (hints)

Although intuitively similar, dealing with continuous random variables is technically difficult

X=x does <u>not</u> describe an *event* in a continuous setting For technical reasons (of *measurability*), a point must have probability <u>zero</u>

Events need to be *subsets*, or better, <u>intervals</u>:

$$X \leq a \ , X \leq b \ , \quad a < X \leq b$$
 . Assuming $a < b$

Probability measures these subsets

$$P(X \le b) = P(X \le a) + P(a < X \le b)$$
These two events are disjoint

$$P(a < X \le b) = P(X \le b) - P(X \le a)$$

$$P_X(x) = P(X \le x) \qquad \qquad \qquad \text{Shorthand notation}$$

Density and Cumulative Distribution

Probability Density Function (pdf)

It is defined as the derivative

$$p_X(x) := \frac{dP_X(x)}{dx}$$

provided that it <u>exists</u> everywhere in \mathcal{X} and is non-negative: $p_X(x) \geq 0$

Probability Measure as Cumulative Distribution Function (CDF)

Cumulative Distribution Function (CDF)

$$P(a < X \leq b) := \int_a^b p_X(x) \; dx$$

As a probability measure, it must integrate to unity

$$\int_{x \in \mathcal{X}} p_X(x) \ dx = 1$$

Note that p(x) may well be above 1 (its integral over the value space $\mathcal X$ will be equal to one)

Artificial Intelligence 2024–2025 Probabilistic Reasoning [29]

Expected value of a random variable

(also: *Expectation*)

Basic definition*

$$\mathbb{E}_X[X] := \sum_{x \in \mathcal{X}} x \ P(X = x)$$

*More concise notation**

$$\mathbb{E}[X] := \sum_{x \in \mathcal{X}} x \, P_X(x)$$

Continuous case

$$\mathbb{E}[X] := \int_{x \in \mathcal{X}} x \ p_X(x) dx$$

Also denoted as: μ_X

Expectation is a linear operator

$$\mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$
$$\mathbb{E}[cX] = c\mathbb{E}[X]$$

Conditional expectation*

$$\mathbb{E}_X[X|Y=y] = \mathbb{E}[X|Y=y] := \sum_{x \in \mathcal{X}} x \ P(X=x|Y=y)$$

(*) In a finitary setting

Variance of a random variable

Basic definition

$$\mathrm{Var}(X) := \mathbb{E}_X[(X-\mathbb{E}_X[X])^2] = \mathbb{E}_X[(X-\mu_X)^2]$$
 ______ Also denoted as: σ_X^2

where*:

$$Var(X) := \sum_{x \in \mathcal{X}} P(X = x) (x - \mu)^2$$

 $\sigma_X := \sqrt{\operatorname{Var}(X)}$ Standard Deviation

Variance is <u>not</u> a linear operator

Conditional Variance

$$Var(X|Y=y) := \mathbb{E}_X[(X - \mathbb{E}_X[X|Y=y])^2 | Y=y]$$

Variance lemma

$$Var(X) = \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}[X^2] - 2\mu_X \mathbb{E}[X] + \mu_X^2$$
$$= \mathbb{E}[X^2] - 2\mu_X^2 + \mu_X^2 = \mathbb{E}[X^2] - \mu_X^2$$

$$\mathbb{E}[X^2] = \mu_X^2 + \sigma_X^2$$

(*) In a finitary setting