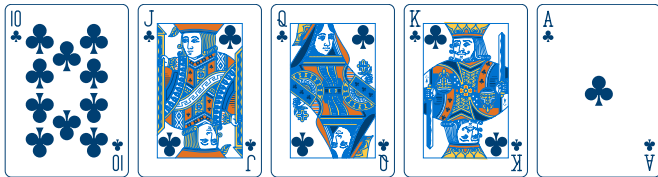


2P7: Probability & Statistics

Probability Fundamentals

Thierry Savin

Lent 2024



the *royal flush*, the best possible hand in poker, has a probability 0.000154%



Probability: mathematics of uncertain events

Statistics: science of collecting and analysing data

- ▶ Probability is logically self-contained
 - Few rules
 - Answers all follow logically from the rules
 - Computations can be tricky (but rarely messy)
 - **Example:** a fair coin is tossed 100 times, what is the probability of 60 or more heads?
- ▶ Statistics apply probability to draw conclusions from data
 - Computations can be messy (and tricky)
 - **Example:** an unknown coin is tossed 100 times and lands 60 heads, what can we conclude about its fairness?



Central questions of the student:

- ▶ **Why** do we need this?
 - Make inference about uncertain events
 - Test the strength of statistical evidence
 - Form the basis of many other theories
- ▶ **How** is it possible to say something about uncertain events?
How can we measure uncertainty?
- ▶ How do I get a good mark?

Examples of applications: failure analysis, design, risk assessment, reliability theory, environmental regulations, inventory theory, computing and simulation, mathematical finance, queueing theory, disease spread, clinical trials, quantum physics, telecommunication, traffic engineering, fitting and machine learning, neural dynamics, statistical mechanics, ...

Probability and statistics: important tools in all kinds of
Engineering!



Seven lectures (weeks 1-4) to cover the following:

1. Probability Fundamentals
2. Discrete Probability Distributions
3. Continuous Random Variables
4. Manipulating and Combining Distributions
5. Decision, Estimation and Hypothesis Testing



Introduction (*what we're doing now...*)

Foundations of Probability

Conditional Probability

Discrete random variables

Expectation and Entropy



- ▶ In classical **frequentist** statistics, the probability of an event is defined as “its long-run frequency in a repeatable experiment”.

Example: “the probability of rolling a 6 with a fair dice is $\frac{1}{6}$ ” because this is the relative frequency of this event as the number of experiments tends to infinity.

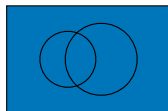
- ▶ However, some notions of chance don't lend themselves to a frequentist approach, and an interpretation of probability with a (subjective) degree of belief is possible; this is known as the **Bayesian** interpretation.

Example: “there is a 50% chance that the arctic polar ice cap will have melted by the year 2100”, it is not possible to define a repeatable experiment.

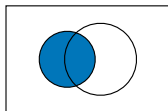
Both approaches can be treated using the same probability theory.



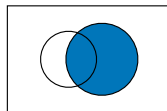
- ▶ A **set** Ω is a collection of **elements**.
- ▶ **Element**: We write $\omega \in \Omega$ to mean the element ω is in the set Ω .
- ▶ **Subset**: We say the set \mathcal{A} is a subset of Ω if all of its elements are in Ω . We write this as $\mathcal{A} \subset \Omega$ ($\mathcal{A} \subseteq \Omega$ if $\mathcal{A} = \Omega$ possible).
- ▶ **Complement**: The complement of \mathcal{A} in Ω is the set of elements of Ω that are not in \mathcal{A} . We write this as \mathcal{A}^c .
- ▶ **Union**: The union of \mathcal{A} and \mathcal{B} is the set of all elements in \mathcal{A} or \mathcal{B} or both. We write this as $\mathcal{A} \cup \mathcal{B}$.
- ▶ **Intersection**: The intersection of \mathcal{A} and \mathcal{B} is the set of all elements in both \mathcal{A} and \mathcal{B} . We write this as $\mathcal{A} \cap \mathcal{B}$.
- ▶ **Empty set**: The empty set is the set with no elements. We denote it \emptyset . If $\mathcal{A} \cap \mathcal{B} = \emptyset$, \mathcal{A} and \mathcal{B} are said to be *disjoint*.



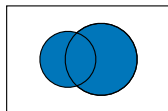
Ω



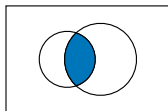
A



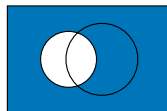
B



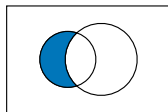
$A \cup B$



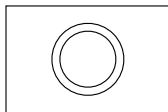
$A \cap B$



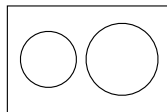
A^c



$A \cap B^c$



$A \subset B$



$A \cap B = \emptyset$

Examples: prove De Morgan's laws

$$\begin{aligned}(A \cup B)^c &= A^c \cap B^c \\ (A \cap B)^c &= A^c \cup B^c\end{aligned}$$



- ▶ **Experiment**: a repeatable procedure with well-defined¹ possible outcomes.
- ▶ **Sample space**: the *set* of all possible outcomes; noted Ω .
- ▶ **Event**: a *subset* of the sample space, $\mathcal{A} \subseteq \Omega$.

Examples:

- ▶ Experiment 1: toss a fair coin, report if it lands heads or tails.
 - $\Omega = \{H, T\}$.
 - $\mathcal{A} = \text{"heads"} = \{H\}$ (single *element*).
- ▶ Experiment 2: toss a fair coin 3 times, list the results.
 - $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$
 - $\mathcal{A} = \text{"exactly 2 heads"} = \{HHT, HTH, THH\}$.

¹each outcome is unique, and different outcomes are *mutually exclusive*



The probability \mathbb{P} is a *measure* that verifies the following:

- The probability of an event is a non-negative real number²

$$\mathbb{P}[\mathcal{A}] \in \mathbb{R} \quad \text{and} \quad \mathbb{P}[\mathcal{A}] \geq 0, \quad \text{for all } \mathcal{A} \subseteq \Omega$$

- The sample space (also called “certain event”) has unit probability

$$\mathbb{P}[\Omega] = 1$$

- Additivity for *incompatible* events (i.e. disjoint sets):

$$\mathbb{P}[\mathcal{A} \cup \mathcal{B}] = \mathbb{P}[\mathcal{A}] + \mathbb{P}[\mathcal{B}] \quad \text{if } \mathcal{A} \cap \mathcal{B} = \emptyset$$

Note: it is OK to write $\mathbb{P}[\mathcal{A}]$ as $\mathcal{P}[\mathcal{A}]$ in your handwritten notes/exam.

² \mathbb{R} denotes the set of real numbers.



These three axioms are sufficient to derive³ the following:

- ▶ Monotonicity

$$\text{if } \mathcal{A} \subseteq \mathcal{B} \text{ then } \mathbb{P}[\mathcal{A}] \leq \mathbb{P}[\mathcal{B}]$$

- ▶ Probability of the empty set

$$\mathbb{P}[\emptyset] = 0$$

- ▶ Complement rule

$$\mathbb{P}[\mathcal{A}^c] = 1 - \mathbb{P}[\mathcal{A}]$$

- ▶ Numeric bound

$$0 \leq \mathbb{P}[\mathcal{A}] \leq 1, \quad \text{for all } \mathcal{A} \subseteq \Omega$$

- ▶ Addition law

$$\mathbb{P}[\mathcal{A} \cup \mathcal{B}] = \mathbb{P}[\mathcal{A}] + \mathbb{P}[\mathcal{B}] - \mathbb{P}[\mathcal{A} \cap \mathcal{B}]$$

- ▶ Sum rule

$$\mathbb{P}[\mathcal{A} \cap \mathcal{B}] + \mathbb{P}[\mathcal{A} \cap \mathcal{B}^c] = \mathbb{P}[\mathcal{A}]$$

³proofs in the examples paper



Another consequence of the additivity axiom is

$$\mathbb{P}[\mathcal{A}] = \sum_{\omega \in \mathcal{A}} \mathbb{P}[\omega]$$

where $\{\omega \in \Omega\}$ are the mutually-disjoint *individual outcomes* (sometimes called “atomic events”) of the experiment.

- ▶ These calculations are often done with *combinatorics*;
- ▶ They inform on the “physical” origin of the probability of an event;
- ▶ We won’t do much of these here.

Conditional probability answers the question “**how does the probability of an event change if we have extra information?**”.

Example: Toss a fair coin 3 times.

- ▶ What is the probability of 3 heads?

$\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$, with all outcomes equally likely, so with $\mathcal{A} = \{HHH\}$

$$\mathbb{P}[\mathcal{A}] = 1/8$$

- ▶ **Given the 1st toss is heads**, what is the probability of 3 heads?

Reduced sample space to $\mathcal{B} = \{HHH, HHT, HTH, HTT\}$, with all outcomes equally likely, so

$$\mathbb{P}[\text{“ } \mathcal{A} \text{ occurs given } \mathcal{B} \text{ occurred ”}] = 1/4$$

The conditional probability of an event \mathcal{A} knowing that an event \mathcal{B} occurred is written:

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] \quad (\text{“ } \mathcal{A} \text{ given } \mathcal{B} \text{ ”})$$



The formal definition of conditional probability⁴ reads:

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] = \frac{\mathbb{P}[\mathcal{A} \cap \mathcal{B}]}{\mathbb{P}[\mathcal{B}]}, \quad \text{provided } \mathbb{P}[\mathcal{B}] \neq 0$$

Why? Frequency interpretation:

- ▶ Repeat experiment N times keeping track of events \mathcal{A} and \mathcal{B} .
- ▶ $N_{\mathcal{B}}$ number of times \mathcal{B} is realised:

$$\mathbb{P}[\mathcal{B}] \stackrel{N \rightarrow \infty}{\approx} \frac{N_{\mathcal{B}}}{N}$$

- ▶ $N_{\mathcal{A} \cap \mathcal{B}}$ number of times both \mathcal{A} and \mathcal{B} occur:

$$\mathbb{P}[\mathcal{A} \cap \mathcal{B}] \stackrel{N \rightarrow \infty}{\approx} \frac{N_{\mathcal{A} \cap \mathcal{B}}}{N}$$

- ▶ Conditional probability:

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] \stackrel{N_{\mathcal{B}} \rightarrow \infty}{\approx} \frac{N_{\mathcal{A} \cap \mathcal{B}}}{N_{\mathcal{B}}} \quad \text{“ proportion of } \mathcal{A} \text{ among the occurrences of } \mathcal{B} \text{ ”}$$

⁴one can show that the conditional probability follows the 3 axioms

- ▶ Product rule

$$\mathbb{P}[\mathcal{A} \cap \mathcal{B}] = \mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}]$$

Proof: from the definition of conditional probability □

- ▶ Law of Total Probability

$$\mathbb{P}[\mathcal{A}] = \mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}] + \mathbb{P}[\mathcal{A}|\mathcal{B}^c] \mathbb{P}[\mathcal{B}^c]$$

Proof: from the sum rule □

- ▶ More generally, with $\{\mathcal{B}_k : k = 1, 2, 3, \dots, n\}$ a set of *pairwise incompatible* events with⁵ $\bigcup_{k=1}^n \mathcal{B}_k = \Omega$,

$$\mathbb{P}[\mathcal{A}] = \sum_{k=1}^n \mathbb{P}[\mathcal{A}|\mathcal{B}_k] \mathbb{P}[\mathcal{B}_k]$$

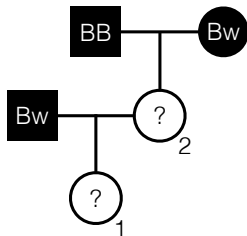
⁵We denote $\bigcup_{k=1}^n \mathcal{B}_k = \mathcal{B}_1 \cup \mathcal{B}_2 \cup \dots \cup \mathcal{B}_n$.

Conditional Probability

Example: Pedigree



UNIVERSITY OF
CAMBRIDGE
Department of Engineering



- ▶ Female (circle) and male (square) cats transmit each one colour gene to offspring
- ▶ Black (B) dominant, white (w) recessive
- ▶ What is the probability that cat 1 is heterozygous (carries both genes) if it is black?

\mathcal{B}_i : "cat i is black"

\mathcal{H}_j : "cat j is heterozygous"

What is $\mathbb{P}[\mathcal{H}_1|\mathcal{B}_1]$?

$$\begin{aligned}\mathbb{P}[\mathcal{B}_1] &= \mathbb{P}[\mathcal{B}_1|\mathcal{H}_2] \times \mathbb{P}[\mathcal{H}_2] + \mathbb{P}[\mathcal{B}_1|\mathcal{H}_2^c] \times \mathbb{P}[\mathcal{H}_2^c] \\ &= \frac{3}{4} \times \frac{1}{2} + 1 \times \frac{1}{2} = \frac{7}{8}\end{aligned}$$

$$\begin{aligned}\mathbb{P}[\mathcal{H}_1] &= \mathbb{P}[\mathcal{H}_1|\mathcal{H}_2] \times \mathbb{P}[\mathcal{H}_2] + \mathbb{P}[\mathcal{H}_1|\mathcal{H}_2^c] \times \mathbb{P}[\mathcal{H}_2^c] \\ &= \frac{1}{2} \times \frac{1}{2} + \frac{1}{2} \times \frac{1}{2} = \frac{1}{2}\end{aligned}$$

$$\mathbb{P}[\mathcal{H}_1|\mathcal{B}_1] = \frac{\mathbb{P}[\mathcal{H}_1 \cap \mathcal{B}_1]}{\mathbb{P}[\mathcal{B}_1]} = \frac{\mathbb{P}[\mathcal{H}_1]}{\mathbb{P}[\mathcal{B}_1]} = \frac{4}{7}$$

□



- Bayes' rule

$$\mathbb{P}[\mathcal{B}|\mathcal{A}] = \frac{\mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}]}{\mathbb{P}[\mathcal{A}]}$$

- Bayes' rule tells us how to **invert** conditional probabilities, i.e. to find $\mathbb{P}[\mathcal{B}|\mathcal{A}]$ from $\mathbb{P}[\mathcal{A}|\mathcal{B}]$.
- Proof via the *product rule*

$$\begin{aligned}\mathbb{P}[\mathcal{A} \cap \mathcal{B}] &= \mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}] \\ &= \mathbb{P}[\mathcal{B} \cap \mathcal{A}] \\ &= \mathbb{P}[\mathcal{B}|\mathcal{A}] \mathbb{P}[\mathcal{A}] \quad \square\end{aligned}$$

- Often used with the *law of total probability*

$$\mathbb{P}[\mathcal{B}|\mathcal{A}] = \frac{\mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}]}{\mathbb{P}[\mathcal{A}|\mathcal{B}] \mathbb{P}[\mathcal{B}] + \mathbb{P}[\mathcal{A}|\mathcal{B}^c] \mathbb{P}[\mathcal{B}^c]}$$

Conditional Probability

Example: Covid Test



UNIVERSITY OF
CAMBRIDGE
Department of Engineering

The rapid antigen test for covid has a sensitivity of 78% (true positive rate), and a specificity of 97% (true negative rate). The population sees 5% incidence of covid.

What is the probability that a person testing positive has covid?

\mathcal{T} : The test is positive

\mathcal{C} : The person has covid

What is $\mathbb{P}[\mathcal{C}|\mathcal{T}]$?

The data tell us:

$$\begin{aligned}\mathbb{P}[\mathcal{T}|\mathcal{C}] &= 0.78 & \mathbb{P}[\mathcal{T}|\mathcal{C}^c] &= 0.03 \\ \mathbb{P}[\mathcal{T}^c|\mathcal{C}] &= 0.22 & \mathbb{P}[\mathcal{T}^c|\mathcal{C}^c] &= 0.97 \\ \mathbb{P}[\mathcal{C}] &= 0.05 & \mathbb{P}[\mathcal{C}^c] &= 0.95\end{aligned}$$

Now we write:

$$\begin{aligned}\mathbb{P}[\mathcal{C}|\mathcal{T}] &= \frac{\mathbb{P}[\mathcal{T}|\mathcal{C}] \mathbb{P}[\mathcal{C}]}{\mathbb{P}[\mathcal{T}]} = \frac{\mathbb{P}[\mathcal{T}|\mathcal{C}] \mathbb{P}[\mathcal{C}]}{\mathbb{P}[\mathcal{T}|\mathcal{C}] \mathbb{P}[\mathcal{C}] + \mathbb{P}[\mathcal{T}|\mathcal{C}^c] \mathbb{P}[\mathcal{C}^c]} \\ &= \frac{0.78 \times 0.05}{0.78 \times 0.05 + 0.03 \times 0.95} = 58\%\end{aligned}$$

□



- ▶ Two events are *independent* if the knowledge that one occurred does not change the probability that the other occurs. **Independence is fundamental**, as we shall see later.
- ▶ In mathematical terms, that means that \mathcal{A} and \mathcal{B} are independent if $\mathbb{P}[\mathcal{A}|\mathcal{B}] = \mathbb{P}[\mathcal{A}]$. Using the product rule, we arrive at the ...
- ▶ Formal definition of independence:

$$\mathcal{A} \text{ and } \mathcal{B} \text{ independent} \quad \Leftrightarrow \quad \mathbb{P}[\mathcal{A} \cap \mathcal{B}] = \mathbb{P}[\mathcal{A}] \times \mathbb{P}[\mathcal{B}]$$

- ▶ Frequency interpretation:

$$\frac{N_{\mathcal{A}}}{N} \stackrel{N \rightarrow \infty}{\approx} \mathbb{P}[\mathcal{A}] = \frac{\mathbb{P}[\mathcal{A} \cap \mathcal{B}]}{\mathbb{P}[\mathcal{B}]} \stackrel{N \rightarrow \infty}{\approx} \frac{N_{\mathcal{A} \cap \mathcal{B}}}{N} \times \frac{N}{N_{\mathcal{B}}} = \frac{N_{\mathcal{A} \cap \mathcal{B}}}{N_{\mathcal{B}}}$$

Knowing when \mathcal{B} occurs is irrelevant to the probability of \mathcal{A} .



► Formal definition:

- A *discrete* random variable X is a function of the outcomes of a random experiment, $X : \Omega \rightarrow \mathbb{X}$, that takes a *discrete* set of scalar values forming \mathbb{X} (typically \mathbb{X} is a subset of \mathbb{R} , the set of real numbers).
- \mathbb{X} is called the *support* (or sometimes *alphabet*) of X .

► Example: game with 2 dice

- Roll a dice twice and record the outcomes (i, j) with i and j the result of the 1st and 2nd roll, respectively.
- $\Omega = \{(1, 1), (1, 2), (1, 3), \dots, (6, 6)\}$ (36 outcomes)
- We define the random variable $X(i, j) = \max\{i, j\}$.
- The event $X = 3$ is $\{(1, 3), (2, 3), (3, 3), (3, 2), (3, 1)\}$ and $\mathbb{P}[X = 3] = \frac{5}{36}$.
- We can build the table:

x	1	2	3	4	5	6	7...
$\mathbb{P}[X = x]$	$\frac{1}{36}$	$\frac{3}{36}$	$\frac{5}{36}$	$\frac{7}{36}$	$\frac{9}{36}$	$\frac{11}{36}$	0...



- ▶ The **probability mass function (PMF)** of a *discrete* random variable X is the function

$$P_X : \mathbb{X} \rightarrow [0, 1] \quad \text{with} \quad P_X(x) = \mathbb{P}[X = x]$$

- ▶ About the notations:
 - X designates the **random variable**. A random variable is often denoted by capital roman letters, such as $X, Y, Z, T \dots$
 - x is the **independent variable** (arbitrary input) and I could use whatever symbol for it: $P_X(a) = \mathbb{P}[X = a]$ does not change the definition of P_X .
 - the name of the **function** is P_X , with the “subscript X ” to remind us it is associated with the random variable X .
- ▶ Finding the PMF of a random variable comes from how the random variable originates from the outcomes.

Note: it is OK to write $P_X(x)$ as $\mathcal{P}_X(x)$ in your handwritten notes/exam.



- ▶ Note that the events “ $X = a$ ” and “ $X = b$ ” are disjoint if $a \neq b$. Hence $\sum_{x \in \mathbb{X}} P_X(x) = 1$.
- ▶ For $x \notin \mathbb{X}$, we can set $P_X(x) = 0$ (if x out of range, “ $X = x$ ” = \emptyset the empty event)
- ▶ It will be useful to introduce the **cumulative distribution function**:

$$F_X(x) = \mathbb{P}[X \leq x]$$

It is obtained by adding up the probabilities $P_X(\xi)$ as ξ runs from $-\infty$ to x , $F_X(x) = \sum_{\xi \leq x} P_X(\xi)$. It has the following properties:

- F_X is non-decreasing: $F_X(a) \leq F_X(b)$ if $a \leq b$
- $\lim_{x \rightarrow -\infty} F_X(x) = 0$ and $\lim_{x \rightarrow \infty} F_X(x) = 1$
- $\mathbb{P}[a < X \leq b] = F_X(b) - F_X(a)$
using $\{X \leq b\} = \{X \leq a\} \cup \{a < X \leq b\}$



Discrete random variables

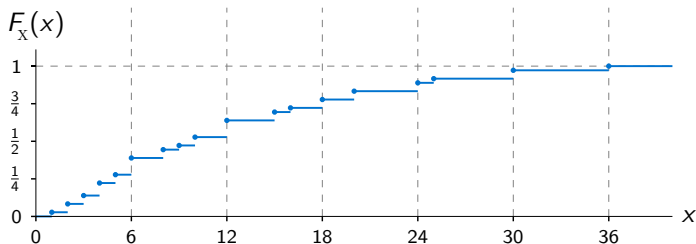
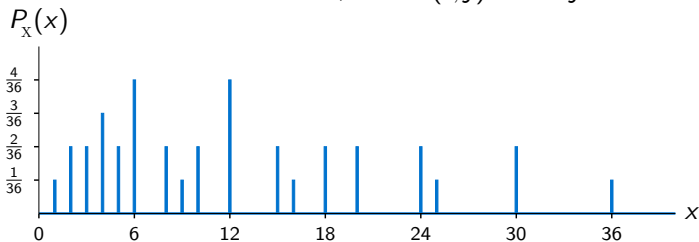
Probability mass function



UNIVERSITY OF
CAMBRIDGE

Department of Engineering

Game with 2 dice, and $X(i, j) = i \times j$





For two random variables $X : \Omega \rightarrow \mathbb{X}$ and $Y : \Omega \rightarrow \mathbb{Y}$ defined on the *same* sample space, we can introduce the **joint probability mass function**:

$$P_{XY} : \mathbb{X} \times \mathbb{Y} \rightarrow [0, 1] \quad \text{with} \quad P_{XY}(x, y) = \mathbb{P}[X = x \cap Y = y]$$

And we inherit the following ideas from our previous discussion:

- **Conditional probability**

$$P_{X|Y} : \mathbb{X} \times \mathbb{Y} \rightarrow [0, 1] \quad \text{with} \quad P_{X|Y}(x|y) = \frac{P_{XY}(x, y)}{P_Y(y)}$$

- **Law of total probability**, also known as **marginalisation**

$$P_X(x) = \sum_{y \in \mathbb{Y}} P_{XY}(x, y)$$

- **Bayes' rule**

$$P_{Y|X}(y|x) = \frac{P_{X|Y}(x|y)P_Y(y)}{P_X(x)} = \frac{P_{X|Y}(x|y)P_Y(y)}{\sum_{\xi \in \mathbb{Y}} P_{X|Y}(x|\xi)P_Y(\xi)}$$

Note: OK to write $P_{XY}(x, y)$ as $\mathbb{P}_{XY}(x, y)$ and $P_{X|Y}(x|y)$ as $\mathbb{P}_{X|Y}(x|y)$.



- Two random variables X and Y are **independent** iff *all the events* corresponding to values of X are independent of *all the events* corresponding to values of Y :

$$\begin{aligned}P_{XY}(x, y) &= \mathbb{P}[X = x \cap Y = y] \\ &= \mathbb{P}[X = x] \times \mathbb{P}[Y = y]\end{aligned}$$

$$P_{XY}(x, y) = P_X(x)P_Y(y) \quad \text{for all } x, y \in \mathbb{X} \times \mathbb{Y}$$

- For more than two random variables, we can also define a joint multivariate probability mass function $P_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n)$. The variables are *mutually independent* iff

$$P_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n) = P_{X_1}(x_1)P_{X_2}(x_2) \dots P_{X_n}(x_n)$$

for all $x_1, x_2, \dots, x_n \in \mathbb{X}_1 \times \mathbb{X}_2 \times \dots \mathbb{X}_n$.

This is different from *pairwise independence*, which only requires

$$P_{X_i X_j}(x_i, x_j) = P_{X_i}(x_i)P_{X_j}(x_j) \quad \text{for all } i, j \text{ and all } x_i, x_j \in \mathbb{X}_i \times \mathbb{X}_j$$

The XOR gate

- ▶ Consider two independent binary random variables X and Y with $P_X(x) = \frac{1}{2}$ if $x = 0$ or $x = 1$, and $P_X(x) = 0$ otherwise. X and Y are *identically distributed*, so $P_X = P_Y$.
- ▶ A third random variable Z is obtained by $Z = X \text{ XOR } Y$ ($a \text{ XOR } b = 0$ when $a = b$, 1 otherwise).
- ▶ Show that X, Y, Z are *pairwise*, but not *mutually*, independent.

The joint distribution of X, Y, Z :

$$P_{XYZ}(0, 0, 0) = \frac{1}{4}$$

$$P_{XYZ}(0, 0, 1) = 0$$

$$P_{XYZ}(0, 1, 0) = 0$$

$$P_{XYZ}(0, 1, 1) = \frac{1}{4}$$

$$P_{XYZ}(1, 0, 0) = 0$$

$$P_{XYZ}(1, 0, 1) = \frac{1}{4}$$

$$P_{XYZ}(1, 1, 0) = \frac{1}{4}$$

$$P_{XYZ}(1, 1, 1) = 0$$

By marginalisation over Y :

$$P_{XZ}(0, 0) = P_{XYZ}(0, 0, 0) + P_{XYZ}(0, 1, 0) = \frac{1}{4}$$

$$P_{XZ}(0, 1) = P_{XYZ}(0, 0, 1) + P_{XYZ}(0, 1, 1) = \frac{1}{4}$$

$$P_{XZ}(1, 0) = P_{XYZ}(1, 0, 0) + P_{XYZ}(1, 1, 0) = \frac{1}{4}$$

$$P_{XZ}(1, 1) = P_{XYZ}(1, 0, 1) + P_{XYZ}(1, 1, 1) = \frac{1}{4}$$

and $P_{XZ} = P_{YZ}$ by symmetry. By further marginalisation over X :

$$P_Z(0) = P_{XZ}(0, 0) + P_{XZ}(1, 0) = \frac{1}{2}$$

$$P_Z(1) = P_{XZ}(0, 1) + P_{XZ}(1, 1) = \frac{1}{2}$$

Verify that $P_{XZ} = P_X P_Z$, $P_{YZ} = P_Y P_Z$ and $P_{XY} = P_X P_Y$, but $P_{XYZ} \neq P_X P_Y P_Z$.



- ▶ We define the **expectation** of a random variable as the “centre of mass” of its distribution:

$$\mathbb{E}[X] = \sum_{x \in \mathbb{X}} x P_X(x) \quad (\text{OK to write } \mathbb{E}[x])$$

- ▶ More generally, we can define the expectation of *any function* of *any number* of random variables; for example

$$\mathbb{E}[g(X)] = \sum_{x \in \mathbb{X}} g(x) P_X(x) \quad \text{or} \quad \mathbb{E}[f(X, Y)] = \sum_{x, y \in \mathbb{X} \times \mathbb{Y}} f(x, y) P_{XY}(x, y)$$

- ▶ The expectation is **linear**:

$$\mathbb{E}[aX + bY] = a \mathbb{E}[X] + b \mathbb{E}[Y] \quad \text{for all } a, b \in \mathbb{R} \times \mathbb{R}$$

proof:
$$\begin{aligned} \mathbb{E}[aX + bY] &= \sum_{x, y \in \mathbb{X} \times \mathbb{Y}} (ax + by) P_{XY}(x, y) \\ &= a \sum_{x \in \mathbb{X}} x \sum_{y \in \mathbb{Y}} P_{XY}(x, y) + b \sum_{y \in \mathbb{Y}} y \sum_{x \in \mathbb{X}} P_{XY}(x, y) \\ &= a \sum_{x \in \mathbb{X}} x P_X(x) + b \sum_{y \in \mathbb{Y}} y P_Y(y) \quad \text{by marginalisation} \\ &= a \mathbb{E}[X] + b \mathbb{E}[Y] \end{aligned}$$





- ▶ For two *independent* random variables X and Y :

$$\begin{aligned}\mathbb{E}[X Y] &= \sum_{x,y \in \mathbb{X} \times \mathbb{Y}} x y P_{XY}(x, y) = \sum_{x,y \in \mathbb{X} \times \mathbb{Y}} x y P_X(x) P_Y(y) \\ &= \sum_{x \in \mathbb{X}} x P_X(x) \sum_{y \in \mathbb{Y}} y P_Y(y) \\ \mathbb{E}[X Y] &= \mathbb{E}[X] \mathbb{E}[Y]\end{aligned}$$

- ▶ This is not true in general of two random variables that are not independent.
- ▶ Two random variables X and Y for which $\mathbb{E}[X Y] = \mathbb{E}[X] \mathbb{E}[Y]$ are said to be *uncorrelated*.
- ▶ X and Y independent \Rightarrow X and Y uncorrelated but *the inverse is not true*.

The expectation can be used to further characterise a probability function by calculating its **central moments**.

- ▶ The central *second* moment, or **variance**, characterises the “spread” of the distribution as

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

the averaged squared-difference to the mean.

- ▶ The central *third* moment $\mathbb{E}[(X - \mathbb{E}[X])^3]$, characterises the “asymmetry” of the distribution.
- ▶ The central *fourth* moment $\mathbb{E}[(X - \mathbb{E}[X])^4]$, characterises the “tailedness” of the distribution.
- ▶ The central *n*-th moment $\mathbb{E}[(X - \mathbb{E}[X])^n]$, characterises higher-order features of the shape of the distribution.

We can rewrite the variance as $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

$$\begin{aligned}\text{Var}[X] &= \mathbb{E}[X^2] - 2\mathbb{E}[X \mathbb{E}[X]] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]^2 + \mathbb{E}[X]^2 \\ \text{Var}[X] &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$





- ▶ The idea of entropy⁶ in **information theory** was introduced by Claude Shannon (1916-2001) to convey the “**surprise**” of an event:
 - an event with 100% probability is perfectly unsurprising (and yields no information);
 - the less probable an event is, the more surprising it is (and the more information it yields).

- ▶ We define the *information content* of a random variable X as

$$I_x(x) = -\log_2 P_x(x)$$

- ▶ The **entropy** (in *bit*) is the average of the information content:

$$\mathbb{E}[I_x(X)] = \mathbb{H}[X] = - \sum_{x \in \mathbb{X}} P_x(x) \log_2 P_x(x)$$

(OK to write $H[X]$)

- ▶ We will discuss this, with examples, in the next part.

⁶we're not talking about thermodynamics here, don't panic

You can attempt Problems 1 to 5 of Examples Paper 5