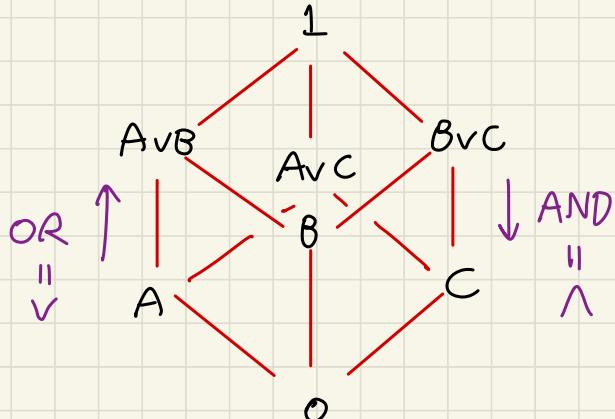


# Probability & Statistics



Logical operation	Set-theoretic
OR	$\cup$
AND	$\cap$
NOT	Complement

Hasse diagram of Boolean algebra with three elementary events  $A, B, C$

## Probability axioms

- ①  $\Pr(A) \geq 0$
- ②  $A$  is certain  $\Leftrightarrow \Pr(A) = 1$  Mutual exclusive
- ③  $\Pr(A \text{ OR } B) = \Pr(A) + \Pr(B)$  if  $A \cap B = \emptyset$

Random variable  $X$  has values  $x_1, x_2, \dots, x_N$   
 $\Pr(X=x) = p_X(x) = p(x)$  Elementary events  
 /alternatives

## Joint distribution

$$\Pr(X=x, Y=y) = p(x, y)$$

## Marginal distribution

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

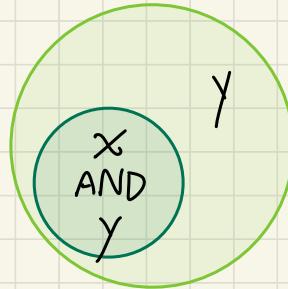
$$p(y) = \sum_x p(x, y)$$

## Conditional probability

(Definition)  $\Pr(X=x|Y=y) = p(x|y) = \frac{P(x,y)}{P(y)}$

Probability that  $X$  takes the value  $x$  given that  $Y$  takes the value  $y$ .

$$p(x|y) + p(\text{NOT } x|y) = 1$$



## Bayes theorem

By the equality  $p(x|y)p(y) = p(x,y) = p(y|x)p(x)$

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

## "Law of total probability"

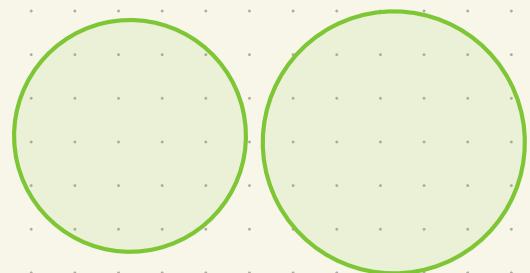
$$p(x) = \sum_y p(x|y)p(y)$$

Mutually exclusive  $\Leftrightarrow E_1 \wedge E_2 = \emptyset$

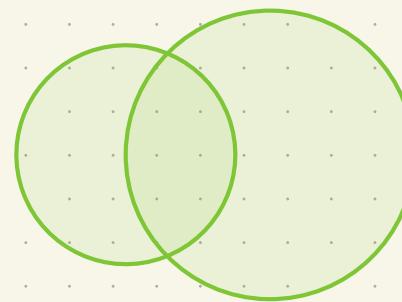
Independent  $\Leftrightarrow \Pr(E_1 \wedge E_2) = \Pr(E_1) \Pr(E_2)$

### Remarks

- ① Mutually exclusive events can't be independent; if a coin flip comes up head, I know that it did not come up tail.
- $E_1 \wedge E_2 = \emptyset \Rightarrow \Pr(E_1 \wedge E_2) = 0$ . But  $\Pr(E_1 \wedge E_2) = \Pr(E_1) \Pr(E_2)$  can't be 0 if  $\Pr(E_1), \Pr(E_2) > 0$ .  $\square$
- ② Independence can't be easily visualized on a Venn diagram, in contrast to mutual exclusivity, since an area in a Venn diagram doesn't conventionally represent probability



Mutually exclusive



Merely overlapping doesn't guarantee independence

How do we assign probabilities in the real world (and stat. mech.)?

There are two main camps of probabilists.

① **Frequentists** interpret probabilities to be whatever the frequencies  $f_x = N_x/N$  of obtaining outcome  $x$  in the long-running limit  $N \rightarrow \infty$ .

"Objective prob." It can be proved that  $f_x$  is likely to be close to the theoretical probability  $p_x$  in the large  $N$  limit, but typically  $f_x \neq p_x$  for any finite  $N$ .

Statements such as  $\Pr(\text{Saturn's mass} = 5.68 \times 10^{26} \text{ kg}) = 0.9$  is also meaningless for frequentists since there is no multiple copies of Saturns with different masses.

② **Bayesians** view probabilities as degrees of belief that can be used for making judgements/decisions. In this view, probabilities are not directly related to observed frequencies (data) can (and should) be used to update relevant probability assignments.

"Subjective prob." We will adopt a pragmatic stance that we assign subjective probabilities based on the best available information and then experimentally check the validity of our probability assignments

$$\underline{\text{Moments}} \quad \langle x^n \rangle = \int dx x^n p(x)$$

Not clear yet that knowing all the moments is equivalent to knowing the PDF

Characteristic function is just the Fourier transform of the PDF

$$\varphi_x(k) = \langle e^{-ikx} \rangle = \int dx e^{-ikx} p(x)$$

Inverse FT gives the PDF back

$$p(x) = \frac{1}{2\pi} \int dk e^{ikx} \varphi_x(k)$$

The characteristic function "generates" all the moments of  $p(x)$  (up to the factor  $(-i)^n$ ), meaning that it is a polynomial (here in  $k$ ) whose coefficients are the moments:

$$\begin{aligned} \varphi_x(k) &= \langle 1 - ikx + \dots + \frac{(-ikx)^n}{n!} + \dots \rangle \\ &= 1 - i\langle x \rangle k - \frac{\langle x^2 \rangle k^2}{2} + \dots + (-i)^n \frac{\langle x^n \rangle k^n}{n!} + \dots \end{aligned}$$

↑  
Normalization    ↑  
Mean

Moments of  $p(x)$  can be computed by differentiating w.r.t.  $k$  and then setting  $k=0$ .

$$\begin{aligned} \frac{d^n}{dk^n} \varphi_x(k) &= \int dx p(x) \left. \frac{d^n}{dk^n} e^{-ikx} \right|_{k=0} \\ &= (-i)^n \int dx x^n p(x) = (-i)^n \langle x^n \rangle \end{aligned}$$

## Cumulants

The cumulant generating function is the natural log of the characteristic function.

$$K_X(k) = \ln \varphi(k) = \sum_{n=1}^{\infty} \frac{(-ik)^n}{n} \langle x^n \rangle_c$$

Cumulant (Defined implicitly)

Beware! Log and integration can't be interchanged

$$\ln \left[ \int_0^{\infty} dx e^{-x} \right] = \ln 1 = 0 \neq \int_0^{\infty} dx \ln e^{-x} = - \int_0^{\infty} dx x = -\infty$$

$$\ln(1+\epsilon) = - \sum_{n=1}^{\infty} \frac{(-1)^n \epsilon^n}{n} = \epsilon - \frac{\epsilon^2}{2} + \frac{\epsilon^3}{3} + \dots$$

$$\begin{aligned} &\Rightarrow \ln \left( 1 - i \langle x \rangle k - \frac{\langle x^2 \rangle k^2}{2} + \dots \right) \\ &= - \sum_{n=1}^{\infty} \frac{(-1)^n}{n} \left( -i \langle x \rangle k - \frac{\langle x^2 \rangle k^2}{2} + \dots \right)^n \end{aligned}$$

First-order term:  $-i \langle x \rangle k \Rightarrow \langle x \rangle_c = \langle x \rangle$

Second-order term:  $-\frac{\langle x^2 \rangle k^2}{2} + \underbrace{\frac{1}{2} \langle x \rangle^2 k^2}_{\text{From } n=2}$

$\underbrace{\phantom{0}}$  From  $n=1$        $\underbrace{\phantom{0}}$  From  $n=2$

Compare with  $-\langle x^2 \rangle k^2 / 2 \Rightarrow \langle x^2 \rangle_c = \langle x^2 \rangle - \langle x \rangle^2 = \sigma_x^2$

This is the part of the 2nd moment that doesn't come from the 1st moment.

## Laws of large number (LLN)

We will derive Markov  $\rightarrow$  Chebyshev  $\rightarrow$  Weak LLN

Moments  $\langle x^n \rangle = \int dx x^n p(x)$

Mean  $\langle x \rangle = \int dx x p(x)$

Variance  $\sigma_x^2 := \langle (x - \langle x \rangle)^2 \rangle$   
 $= \langle x^2 \rangle - \langle x \rangle^2$

Deviation

$$\Delta x = x - \langle x \rangle$$

$$\langle \Delta x \rangle = 0$$

For a non-negative random variable  $X$

$$p(X \geq a) \leq \frac{\langle x \rangle}{a} \quad (\text{Markov})$$

$$\bullet \langle x \rangle = \int_0^\infty dx x p(x) \geq \int_a^\infty dx x p(x)$$

$$\geq \int_a^\infty dx a p(x) = a p(X \geq a)$$

Ineq. saturated if  $\text{Supp}(p) \subseteq \{0, a\}$  (*Delta functions*)  
(at  $x=0$  and  $a$ )

Applying this to the  $\Delta x^2$  to obtain

$$p(\Delta x^2 \geq a^2) \leq \frac{\langle \Delta x^2 \rangle}{a^2} \Leftrightarrow p(|x - \langle x \rangle| \geq a) \leq \frac{\sigma_x^2}{a^2} \quad (\text{Chebyshev})$$

Let  $X_1, X_2, \dots, X_N$  be independent, identically distributed (i.i.d.) random variables with finite mean and variance.

Sample mean  $S = \frac{1}{N} \sum_{k=1}^N X_k$  from many trials

Mean of sample mean  $\langle S \rangle = \frac{1}{N} \sum_{k=1}^N \langle X_k \rangle = \langle X \rangle$

$$\langle S^2 \rangle = \frac{1}{N^2} \sum_{k,l} \langle X_k X_l \rangle$$

$$\begin{aligned} & \langle (\langle X \rangle + \Delta X_k)(\langle X \rangle + \Delta X_l) \rangle \\ &= \langle X \rangle^2 + 2 \cancel{\langle \Delta X_k \rangle} \cancel{\langle \Delta X_l \rangle} + \langle \Delta X_k \Delta X_l \rangle \\ &= \langle X \rangle^2 + \delta_{kl} \langle \Delta X^2 \rangle \end{aligned}$$

$$\therefore \langle S^2 \rangle = \langle X \rangle^2 + \frac{1}{N} \langle \Delta X^2 \rangle = \langle X \rangle^2 + \frac{\sigma_X^2}{N}$$

$$\text{Thus, } \langle \Delta S^2 \rangle = \langle S^2 \rangle - \langle S \rangle^2 = \frac{\sigma_X^2}{N}$$

Now we are equipped to prove the weak LCL.

$$P(|S - \langle S \rangle| \geq \epsilon) \leq \frac{\langle \Delta S^2 \rangle}{\epsilon^2}$$

$$P(|S - \langle X \rangle| \geq \epsilon) \leq \frac{\sigma_X^2}{\epsilon^2 N}$$

Taking the limit  $N \rightarrow \infty$ , we see that

$$\lim_{N \rightarrow \infty} P(|S - \langle X \rangle| \leq \epsilon) \rightarrow 1 \quad (\text{Weak LCL})$$

This provides a link between the theoretical mean and the sample mean. But be careful about what the weak law of large number doesn't say. It doesn't say that  $S$  becomes  $\langle X \rangle$  identically. It only says that  $S$  takes the value  $\langle X \rangle$  with overwhelming probability (convergence in probability)

How to obtain this relation at the level of frequencies?

⇒ Next

$$\langle x \rangle_c = \langle x \rangle$$

$$\langle x^2 \rangle_c = \langle x^2 \rangle - \langle x \rangle^2$$

$$\langle x^3 \rangle_c = \langle x^3 \rangle - 3\langle x^2 \rangle \langle x \rangle + 2\langle x \rangle^3$$

$$\langle x^4 \rangle_c = \langle x^4 \rangle - 4\langle x^3 \rangle \langle x \rangle - 3\langle x^2 \rangle^2 + 12\langle x^2 \rangle \langle x \rangle^2 - 6\langle x \rangle^4$$

### Normal (Gaussian) distribution

Characteristic function

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$$\tilde{\varphi}(k) = \int dx p(x) e^{-ikx} = \exp\left(-ik\mu - \frac{k^2\sigma^2}{2}\right)$$

$$\Rightarrow \ln \tilde{\varphi}(k) = -ik\mu - \frac{k^2\sigma^2}{2} \quad \text{Only the first and second cumulants are non zero!}$$

$$\langle x \rangle_c = \mu, \langle x^2 \rangle_c = \sigma^2$$

$$\Rightarrow \langle x^3 \rangle = 3\sigma^2\mu + \mu^3$$

$$\langle x^4 \rangle = 3\sigma^4 + 6\sigma^2\mu^2 + \mu^4$$

Multinomial

### Binomial distribution

A occurs with prob.  $p$

$$B \longrightarrow \underbrace{\dots}_{q=1-p}$$

In  $N$  trials, the prob. that A occurs exactly  $n$  times is

$$\frac{N!}{n!(N-n)!} p^n q^{N-n}$$

$$p(n_1, \dots, n_d) = \frac{N!}{n_1! \dots n_d!} p_1^{n_1} \dots p_d^{n_d}$$

## Binomial & multinomial distributions

Let  $X$  be a discrete random variable that can take on 2 values  $x_1, x_2$  with probability  $p_1$  and  $p_2$  respectively.

Head  $\xrightarrow{\quad}$  Tail  $\xrightarrow{\quad}$    
 $p \quad \quad \quad q$

In  $N$  trials ( $N$  independent and identical coin tosses),

$$\left( \begin{array}{l} \text{Prob. for a} \\ \text{particular sequence} \\ \text{with } n \text{ heads} \end{array} \right) = p^n q^{N-n}$$

$$\left( \begin{array}{l} \text{Prob. for any} \\ \text{sequence with} \\ n \text{ heads} \end{array} \right) = \binom{N}{n} p^n q^{N-n} =: p(n)$$

Normalization  $\sum_{n=0}^N \binom{N}{n} p^n q^{N-n} \xrightarrow{\text{Binomial thm}} (p+q)^N = 1$

$$\langle n \rangle = \sum_{n=0}^N n p(n) = \sum_n n \binom{N}{n} p^n q^{N-n}$$

$$= p \frac{\partial}{\partial p} \sum \binom{N}{n} p^n q^{N-n} = p \frac{\partial}{\partial p} (p+q)^N$$

$$= N p (p+q)^{N-1} = \boxed{N p} \quad \boxed{\sigma_n^2 = p(1-p)N}$$

## d outcomes

$$P(n_1, \dots, n_d) = \frac{N!}{n_1! \cdots n_d!} p_1^{n_1} \cdots p_d^{n_d}$$

$$\sum_{j=1}^d n_j = N$$

Multinomial theorem

$$\sum_{\substack{n_1, \dots, n_d \\ \sum n_j = N}} P(n_1, \dots, n_d) = (p_1 + \cdots + p_d)^N = 1$$

$$\langle n_j \rangle = p_j \frac{\partial}{\partial p_j} \sum_{n_1, \dots, n_d} P(n_1, \dots, n_d)$$

$$= N_j p_j (p_1 + \cdots + p_d)^{N-1} = N_j p_j$$

## Correlation matrix

$$\langle n_j n_k \rangle = \sum_{n_1, \dots, n_d} n_j n_k P(n_1, \dots, n_d)$$

$$= p_j \frac{\partial}{\partial p_j} \left[ p_k \frac{\partial}{\partial p_k} \left( \sum_{n_1, \dots, n_d} P(n_1, \dots, n_d) \right) \right]$$

$$= p_j \frac{\partial}{\partial p_j} \left[ N p_k (p_1 + \cdots + p_d)^{N-1} \right]$$

$$= N p_j \delta_{jk} (p_1 + \cdots + p_d)^{N-1} + N(N-1) p_j p_k (p_1 + \cdots + p_d)^{N-2}$$

$$= N^2 p_j p_k + N p_j (\delta_{jk} - p_k)$$

$$\sum_{j=1}^d n_j = N$$

$$\begin{aligned}\langle \Delta n_j \Delta n_k \rangle &= \langle (n_j - \langle n_j \rangle)(n_k - \langle n_k \rangle) \rangle \\ &= \langle n_j n_k \rangle - \langle n_j \rangle \langle n_k \rangle \\ &= N p_j (\delta_{jk} - p_k)\end{aligned}$$

Variances  $\langle \Delta n_j^2 \rangle = N p_j (1 - p_j)$

Frequencies

$$f_j = \frac{n_j}{N} \Rightarrow$$



$$\langle f_j \rangle = \frac{\langle n_j \rangle}{N} = p_j$$

$$\langle f_j f_k \rangle = \frac{\langle n_j n_k \rangle}{N^2} = p_j p_k + \frac{p_j}{N} (\delta_{jk} - p_k)$$

$$\langle \Delta f_j \Delta f_k \rangle = \frac{p_j}{N} (\delta_{jk} - p_k) \quad \boxed{\text{Important}}$$

$$\langle \Delta f_j^2 \rangle = p_j (1 - p_j) \quad \text{Variance goes as } 1/N$$

Therefore, the weak LOL implies that  $\overrightarrow{p}$  within the sphere inside the hypercube

$$p(|\vec{f}_j - p_j| \geq \epsilon_j, \forall j) \geq p\left[\sum_j (\vec{f}_j - p_j)^2 \geq \epsilon_j^2 \forall j\right]$$

$\vec{f}$  is within a hypercube

$$= \sum_j \frac{\langle \Delta f_j^2 \rangle}{\epsilon_j^2} = \sum_j \frac{p_j}{\epsilon_j^2} - \frac{\sum_j p_j^2}{N \epsilon_j^2}$$

$$p(|x - \langle x \rangle| \geq a) \leq \frac{\sigma_x^2}{a^2} \quad (\text{Chebyshev})$$

Central limit theorem (CLT) is stronger

$$p(S) = \frac{1}{\sqrt{2\pi(\Delta S)^2}} \exp\left[-\frac{(S - \langle S \rangle)^2}{(\Delta S)^2}\right]$$

To show convergence to a Gaussian distribution (not rigorous), we only need to show that all the cumulants higher than the 3rd cumulants vanish in the  $N \rightarrow \infty$  limit.

Joint cumulant

$$\langle x_1^{n_1} * \cdots * x_N^{n_N} \rangle = (-i)^{n_1} \frac{\partial}{\partial k_1^{n_1}} \cdots (-i)^{n_N} \frac{\partial}{\partial k_N^{n_N}} \ln \tilde{\mathcal{Q}} \Big|_{\vec{k}=0}$$

$\langle x_1 * x_2 \rangle = 0$  if  $x_1$  and  $x_2$  are independent random variables

Independent random variables

$$\tilde{\mathcal{Q}}_{\sum X_j / \sqrt{N}}(k) = \langle e^{-i \sum k_j x_j / \sqrt{N}} \rangle = \prod_{j=1}^N \tilde{\mathcal{Q}}_{X_j} \left( \frac{k_j}{\sqrt{N}} \right)$$

$$\Rightarrow K_{\sum X_j / \sqrt{N}}(k) = \ln \tilde{\mathcal{Q}}_{\sum X_j / \sqrt{N}}(k) = \sum_{j=1}^N K_{X_j} \left( \frac{k_j}{\sqrt{N}} \right)$$

So the  $n$ th cumulant will be attached to the  $\frac{k^n}{N^{n/2}}$  term in the cumulant generating function.

$$\Rightarrow n\text{th cumulant} = \frac{\sum \langle x^n \rangle_c}{N^{n/2}} \stackrel{\text{Identical and cumulant bounded by a const. } C}{\leq} \frac{NC}{N^{n/2}} = N^{1-\frac{n}{2}} C$$

## Entropy

The (Shannon/Gibbs) entropy is a measure of the average information gained from observing an outcome of a random source

$$H(X) = H_X = - \sum_j p_j \ln p_j$$

*non-negative*

If all outcomes are equiprobable

$$p_j = \frac{1}{N} \Rightarrow H_X = \ln N$$

Boltzmann entropy  $S = k_B \ln N$

*Multiplicity  $\Sigma$  in stat. mech.*

Information is additive: suppose that  $X$  and  $Y$  are two independent random variables with  $N$  and  $M$  elementary events respectively.

$$H(X, Y) = - \sum_{x,y} p(x, y) \ln p(x, y)$$

$$= - \sum_x \left[ \sum_y p(x, y) \right] \ln p(x) - \sum_y \left[ \sum_x p(x, y) \right] \ln p(y)$$

$$= H(X) + H(Y) \text{ in particular } H(X, Y) = \ln N + \ln M \text{ if equiprobable}$$

The base of the log differs from field to field. For example, statistical mechanics use the natural log  $\ln$  while information theory often uses log base 2

In the limit  $N \rightarrow \infty$ , the number of times  $x_j$  appears in a sequence is  $\langle n_j \rangle = N p_j$  ("typical" sequence)

(Prob. of any given typical sequence)

$$\begin{aligned}
 &= p_1^{n_1} p_2^{n_2} \dots p_d^{n_d} = p_1^{N p_1} \dots p_d^{N p_d} \\
 &= e^{N p_1 \log p_1} \dots e^{N p_d \log p_d} \\
 &= e^{N(p_1 \log p_1 + \dots + p_d \log p_d)} \\
 &= e^{-N H(x)}
 \end{aligned}$$

$$x^N = e^{N \ln x}$$

Count the number of typical sequences

$$\ln \left( \frac{N!}{n_1! \dots n_d!} \right) = \ln N! - \sum_j \ln n_j!$$

Stirling

$$\ln N! \sim N \ln N - N$$

$$= N \ln N - N - \sum_j (n_j \ln n_j - n_j)$$

$$= N \ln N - \sum_j N p_j \ln(N p_j)$$

$$\sum_j N p_j (\ln N + \ln p_j)$$

$$N \ln N + N \sum_j p_j \ln p_j$$

$$= N \left( - \sum_j p_j \ln p_j \right) = N H(x)$$

Each typical sequence appears with prob.  $e^{-N H(x)}$  and there are  $e^{N H(x)}$  of them, so they take all the probability.