### I - Probability

### Probability spaces

A probability space is an ordered triple  $(\Omega, \mathcal{F}, P)$  in which

- $\star$  the universe  $\Omega$  is a non-empty set containing elementary outcomes  $\omega$
- \*  $\mathcal{F}$  is a set of subsets of  $\Omega$  that is (i) non-empty, (i) contains complements of its members, and (iii) contains countable unions of its members
- \* Probabilities  $P(A) \in [0,1]$  are defined for all  $A \in \mathcal{F}$ , (i)  $P(\Omega) = 1$  and (ii) the probability of a countable union of pairwise disjoint elements of  $\mathcal{F}$  equals the sum of their individual probabilities.
- $\star$  If  $A \in \mathcal{F}$  then  $P(A^c) = 1 P(A)$ , so  $P(\emptyset) = 0$  because  $P(\Omega) = 1$
- \* inclusion-exclusion:  $P(A \cup B) = P(A) + P(B) P(A \cap B)$ , more generally

$$P(\cup_{i=1}^{n} A_{i}) = \sum_{r=1}^{n} (-1)^{r+1} \sum_{1 \leq i_{1} < \dots < i_{r} \leq n} P(A_{i_{1}} \cap \dots \cap A_{i_{r}})$$

- $\star$  Boole's inequality:  $P(\cup_{i\in\mathbb{N}}A_i)\leq \sum_{i\in\mathbb{N}}P(A_i)$
- $\star \text{ Continuity of } P \text{ as a set function: } \lim_{n \to \infty} P\left(\cup_{i=1}^n A_i\right) = P\left(\lim_{n \to \infty} \cup_{i=1}^n A_i\right) \text{ and } \lim_{n \to \infty} P\left(\cap_{i=1}^n A_i\right) = P\left(\lim_{n \to \infty} \cap_{i=1}^n A_i\right)$

# Conditional probabilities and partitions of $\Omega$

- \* Conditional probability of A given B:  $P(A \mid B) = \frac{P(A \cap B)}{P(B)}$
- $\star P(A \cap B) = P(A)P(B) \iff A \text{ and B are independent}$
- \* If  $\{B_i\}_{i\in\mathbb{N}}$  partition  $\Omega$  (i.e.,  $B_i\cap B_j=\emptyset$  when  $i\neq j$  and  $\Omega=\cup_i B_i$ ), then (law of total probability)  $P(A)=\sum_i P(A\mid B_i)P(B_i)$ , giving (Bayes' theorem)

$$P(B_j \mid A) = \frac{P(A \mid B_j)P(B_j)}{\sum_i P(A \mid B_i)P(B_i)}.$$

#### Random variable

A random variable is a function  $X:\Omega\to\mathbb{R}$ . The support of X is  $S_X:=\{x\in\mathbb{R}:\exists\ \omega\in\Omega\ \text{s.t.}\ X(\omega)=x\}$ . Then

- $\star F_X(x) = P(X \le x)$
- $\star F_X \nearrow \text{on } \mathbb{R}$
- $\star F_X(b) F_X(a) = P(a < X < b)$
- $\star 1 F_X(x) = P(X > x)$
- \* The p quantile is  $x_p = \inf \{x \in \mathbb{R} : F_X(x) \ge p\}$ , for  $p \in (0, 1)$ .

### Discrete random variables: $S_X$ is finite or countable

- \* Probability mass function (PMF):  $f_X(x) = P(X = x)$
- $\star P(X \in B) = \sum_{x \in B} f_X(x)$
- \*  $f_{X|X \in B}(x \mid X \in B) = P(X = x \mid X \in B) = \frac{P(X = x \cap X \in B)}{P(X \in B)}$
- $\star \ \operatorname{E}\left\{g(X)\right\} = \sum_{x \in S_{|Y|}} g(x) f_X(x), \operatorname{Var}(X) = \operatorname{E}[\{X \operatorname{E}(X)\}^2] = \operatorname{E}(X^2) \operatorname{E}(X)^2$

# Continuous random variables: $S_X$ is uncountable

- \* Probability density function (PDF):  $f_X(x) = \frac{\mathrm{d}F_X}{\mathrm{d}x}(x) \neq P(X=x)$
- $\star P(X \in B) = \int_{x \in \mathbb{R}} f_X(x) I_B(x) \, \mathrm{d}x = \int_{x \in B} f_X(x) \, \mathrm{d}x$
- $\star f_{X|X \in B}(x) = \frac{f_X(x)I_B(x)}{P(X \in B)}$
- $\star \ \mathrm{E}\left\{g(X)\right\} = \int_{x \in \mathbb{R}} g(x) f_X(x) \ \mathrm{d}x, \mathrm{Var}(X) = \mathrm{E}(X^2) \mathrm{E}(X)^2$

#### Law of total expectation

 $\mathbb{E}\left\{g(X)\right\} = \sum_{i \in \mathbb{N}} \mathbb{E}\left\{g(X) \mid B_i\right\} P(B_i) \text{ where } \left\{B_i\right\}_{i \in \mathbb{N}} \text{ partitions } \Omega$ 

# Generating functions

- $\star$  The moment generating function of X is  $M_X(t) = \mathrm{E}\left(e^{tX}\right)$
- \* The cumulant generating function of X is  $K_X(t) = \ln M_X(t)$

### Notable variables and properties

- \* **Bernoulli variable** For  $A \in \Omega$  we have  $I_A(y) = 1$  if  $y \in A$  and  $I_A(y) = 0$  if  $y \notin A$ . We have  $f_X(x) = P(A)^x (1 P(A))^{1-x}$ ,  $E(I_A) = P(A)$ ,  $Var(I_A) = P(A) (1 P(A))$  and  $S_X = \{0, 1\}$  Can also be defined with a probability p,  $I_p$ , with  $f_X(x) = p^x (1 p)^{1-x}$ ,  $E(I_p) = p$  and  $Var(I_p) = p (1 p)$ .
- \* **Binomial** For  $Z_1, \ldots, Z_n \stackrel{\text{iid}}{\sim} I_p$  we define  $X = Z_1 + \cdots + Z_n$  and write  $X \sim B(n,p), n \in \mathbb{N}, p \in (0,1).$  We have  $f_X(x) = \binom{n}{x} p^x (1-p)^{n-x},$  E(X) = np, Var(X) = np(1-p) and  $S_X = \{0, 1, \ldots, n\}$
- \* Negative binomial If  $X \sim \text{NegBin}(n,p)$ ,  $n \in \mathbb{N}$ ,  $p \in (0,1)$ . We have  $f_X(x) = {x-1 \choose n-1} p^n (1-p)^{x-n}$ ,  $E(X) = \frac{n}{p}$ ,  $Var(X) = \frac{n(1-p)}{p^2}$  and  $S_X = \{n, n+1, \ldots, \infty\}$ . If n=1 we write  $X \sim \text{Geom}(p)$ .
- Hypergeometric If  $X \sim \text{HypGeom}(k,b,n)$ ,  $k,b,n \in \mathbb{N}$ . We have  $f_X(x) = \frac{\binom{b}{x}\binom{n}{n-x}}{\binom{k+b}{n}}$ ,  $\mathrm{E}(X) = n\frac{b}{b+k}$ ,  $\mathrm{Var}(X) = n\frac{bk(b+k-n)}{(b+k)^2(b+k-1)}$  and  $S_X = \{\max(0,b+k-n),\ldots,\min(b,n)\}$
- \* **Poisson** If  $X \sim \operatorname{Poiss}(\lambda)$ ,  $\lambda \in \mathbb{R}^{+*}$ . We have  $f_X(x) = e^{-\lambda} \frac{\lambda^x}{x!}$ ,  $\operatorname{E}(X) = \lambda$ ,  $\operatorname{Var}(X) = \lambda$  and  $S_X = \{0, 1, \dots, \infty\}$ . Note that for  $X_n \sim B(n, p_n)$  with  $\lim_{n \to \infty} np_n = \lambda$ , then  $X_n \stackrel{n \to \infty}{\to} X \sim \operatorname{Poiss}(\lambda)$
- \* Discrete uniform If  $X \sim \text{Unif}\{1,..,n\}, n \in \mathbb{N}$ . We have  $f_X(x) = \frac{1}{n}$ ,  $E(X) = \frac{n+1}{2}$ ,  $Var(X) = \frac{n^2-1}{12}$  and  $S_X = \{1,...,n\}$
- \* Continous Uniform If  $X \sim U(a,b)$ ,  $a < b \in \mathbb{R}$ . We have  $f_X(x) = \frac{1}{b-a}$ ,  $F_X(x) = \frac{x-a}{b-a}$ ,  $E(X) = \frac{b+a}{2}$ ,  $Var(X) = \frac{(b-a)^2}{12}$  and  $S_X = (a,b)$
- \* Exponential If  $X \sim \exp(\lambda)$ ,  $\lambda \in \mathbb{R}^{+*}$ . We have  $f_X(x) = \lambda e^{-\lambda x}$ ,  $F_X(x) = 1 e^{-\lambda x}$ ,  $E(X) = \frac{1}{\lambda}$ ,  $Var(X) = \frac{1}{\lambda^2}$  and  $S_X = (0, +\infty)$
- $\star \text{ Pareto If } X \sim \operatorname{Pareto}(\alpha,\beta), \ \alpha,\beta \in \mathbb{R}^+. \quad \text{We have } f_X(x) = \frac{\alpha\beta^\alpha}{x^{\alpha+1}},$   $F_X(x) = 1 \left(\frac{\beta}{x}\right)^\alpha, \quad \operatorname{E}(X) = \frac{\alpha\beta}{\alpha-1} \quad \text{defined only for } \alpha > 1,$   $\operatorname{Var}(X) = \frac{\beta^2\alpha}{(\alpha-1)^2(\alpha-2)} \quad \text{defined only for } \alpha > 2 \text{ and } S_X = (\beta,+\infty)$
- $\begin{array}{l} \star \text{ Laplace If } X \sim \operatorname{Laplace}(\lambda,\eta), \, \lambda \in \mathbb{R}^+, \, \eta \in \mathbb{R}. \, \text{ We have } f_X(x) = \frac{\lambda}{2} e^{-\lambda |x-\eta|}, \\ F_X(x) = \frac{1}{2} e^{\lambda (x-\eta)} \, \text{ for } \, x \leq \, \, \eta \, \text{ and } \, F_X(x) = 1 \frac{1}{2} e^{-\lambda (x-\eta)} \, \text{ for } \, x \, > \, \eta, \\ \operatorname{E}(X) = \eta, \operatorname{Var}(X) = \frac{2}{\lambda^2} \, \operatorname{and} \, S_X = \mathbb{R} \end{array}$
- $\star \ \, \mathbf{Gamma} \ \, \mathbf{If} \ \, X \ \, \sim \ \, \mathbf{Gamma}(\alpha,\beta), \ \, \alpha,\beta \ \, \in \ \, \mathbb{R}^{+*}, \ \, \mathbf{then} \ \, f_X(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, \\ F_X(x) = \frac{\int_0^{\beta x} t^{\alpha-1} e^{-t} \, \mathrm{d}t}{\Gamma(\alpha)}, \\ \mathbf{E}(X) = \frac{\alpha}{\beta}, \ \, \mathbf{Var}(X) = \frac{\alpha}{\beta^2} \ \, \mathbf{and} \ \, S_X = (0,+\infty)$

- \* Gaussian If  $X \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma \in \mathbb{R}^{+*}$ . We have  $f_X(x) = \frac{e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma^2}\right)^2}}{\sqrt{2\pi\sigma^2}}$ ,  $F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$ ,  $E(X) = \mu$ ,  $Var(X) = \sigma^2$  and  $S_X = \mathbb{R}$
- $\star \ \chi^2 \ \text{For} \ Z_1, \dots, Z_{\nu} \overset{\text{iid}}{\sim} \mathcal{N}(0,1) \ \text{we define} \ X = Z_1^2 + \dots + Z_{\nu}^2 \ \text{and write} \ X \sim \chi_{\nu}^2$  for  $\nu \in \mathbb{N}$ . We have  $f_X(x) = \frac{x^{\nu/2-1}e^{-x/2}}{2^{\nu/2}\Gamma(\nu/2)}, \ F_X(x) = \frac{\int_0^{x/2} t^{\nu/2-1}e^{-t} \ \text{d}t}{\Gamma(\nu/2)}$  for x > 0, and  $\operatorname{E}(X) = \nu$ ,  $\operatorname{Var}(X) = 2\nu$  and  $S_X = \mathbb{R}^{+*}$ .

#### Random vectors

- $\star \mathbb{R}^n \ni X = (X_1, \dots, X_n)$  is a random vector.
- $\star f_X(x) = f_{X_1, \dots, X_n}(x) : \mathbb{R}^n \to \mathbb{R}$  is its joint density function.
- \* If  $A\subset\{1,\ldots,n\}$  satisfies |A|=p, B is its complement, and we write  $x=(x_A,x_B)$ , then  $X_A=X_{i:i\in A}$  has **marginal density function**  $f_{X_A}(x_A)=\int_{x_B\in\mathbb{R}^{n-p}}f_X(x)\,\mathrm{d}x_B=\int_{x_B\in\mathbb{R}^{n-p}}f_X(x_A,x_B)\,\mathrm{d}x_B$  and the **conditional density function** of  $X_B$  given that  $X_A=x_A$  is

$$f_{X_B\mid X_A}(x_B\mid x_A) = \frac{f_X(x_A,x_B)}{f_{X_A}(x_A)}, \quad x_A\in\mathbb{R}^p, \quad x_B\in\mathbb{R}^{n-p}.$$

- $\text{ If } f_{X_B \mid X_A}(x_B \mid x_A) = f_{X_B}(x_B) \text{ for all possible values of } x_A \text{ and } x_B, \text{ then } X_A \text{ and } X_B \text{ are independent and we can write } f_X(x) = f_{X_A}(x_A) f_{X_B}(x_B).$
- \* We define  $\mathbb{E}\left\{g(X)\right\} = \int_{\mathbb{R}^n} g(x) f_X(x) \, \mathrm{d}x$ , where  $g(x) : \mathbb{R}^n \to \mathbb{R}$ .
- $\star$  We define the **mean vector** as  $E(X) = \mu = (E(X_1), \dots, E(X_n))^T \in \mathbb{R}^n$ .
- \* For  $1 \leq k, l < n$ , the covariance is  $Cov(X_k, X_l) = E(X_k X_l) E(X_k)E(X_l)$ . Note that  $X_k \perp \!\!\! \perp X_l \implies Cov(X_k, X_l) = 0$ .
- \* The **correlation** between  $X_k$  and  $X_l$  is  $\operatorname{Corr}(X_k, X_l) := \frac{\operatorname{Cov}(X_k, X_l)}{\sqrt{\operatorname{Var}(X_k)\operatorname{Var}(X_l)}}$ . Note that  $\operatorname{Corr}(X_k, X_l) \in [-1, 1]$
- \* The covariance matrix of X is  $Cov(X) = \Omega \in \mathbf{R}^{n \times n}$ , where  $\Omega_{k,l} = Cov(X_k, X_l)$ . As  $Cov(X_k, X_l) = Cov(X_l, X_k)$ ,  $\Omega = \Omega^T$ , and  $\Omega$  is symmetric positive semi-definite.

### Multivariate Gaussian distribution

- $\text{* For } X = (X_1, \dots, X_n)^T \in \mathbb{R}^{n \times 1} \text{ with } \mathbf{E}(X) = \mu \text{ and } \mathbf{Cov}(X) = \Omega, \text{ if } \forall u \in \mathbb{R}^n, \\ u^T X \sim \mathcal{N}\left(u^T \mu, u^T \Omega u\right), \text{ we write } X \sim \mathcal{N}_n\left(\mu, \Omega\right).$
- \* If Rank( $\Omega$ ) = n then  $f_X(x) = \{(2\pi)^n |\Omega|\}^{-1/2} e^{-(x-\mu)^T \Omega^{-1} (x-\mu)/2}, x \in \mathbb{R}^n$ .
- \* With  $X \sim \mathcal{N}_n(\mu,\Omega)$  and  $A, B, X_A \in \mathbb{R}^p$  and  $X_B \in \mathbb{R}^{n-p}$  defined as above, we define  $\mu_A = \mathrm{E}(X_A), \, \mu_B = \mathrm{E}(X_B), \, \Omega_{AA} = \mathrm{Cov}(X_A) \in \mathbb{R}^{p \times p}, \, \Omega_{BB} = \mathrm{Cov}(X_B) \in \mathbb{R}^{(n-p) \times (n-p)}, \, \mathrm{Cov}(X_A, X_B) = \Omega_{AB} \in \mathbb{R}^{p \times (n-p)}$  and  $\Omega_{BA} = \Omega_{AB}^T$ . Then  $X_A \sim \mathcal{N}_p(\mu_A, \Omega_{AA})$  and

$$X_{\mathcal{A}} \mid X_{\mathcal{B}} = x_{\mathcal{B}} \sim \mathcal{N}_{p} \left( \mu_{\mathcal{A}} + \Omega_{\mathcal{A}\mathcal{B}} \Omega_{\mathcal{B}\mathcal{B}}^{-1} (x_{\mathcal{B}} - \mu_{\mathcal{B}}), \Omega_{\mathcal{A}\mathcal{A}} - \Omega_{\mathcal{A}\mathcal{B}} \Omega_{\mathcal{B}\mathcal{B}}^{-1} \Omega_{\mathcal{B}\mathcal{A}} \right).$$

### Transformation of variables

For  $X=(X_1,\ldots,X_n)\in\mathbb{R}^n$  and a function  $g(x):\mathbb{R}^n\to\mathbb{R}$ , the distribution of Y:=g(X) is  $P(Y\leq y)=F_Y(y)=\int_{x:g(x)\leq y}f_X(x)\,\mathrm{d}x$ .

The case n=1 with  $g \nearrow$  and  $g^{-1}(y)=x$  bijective on  $S_X$  gives  $F_X(x)=P(X \le x)=P(g(X) \le g(x))=F_Y(g(x))$   $\Longrightarrow F_Y(y)=F_X(g^{-1}(y)).$ 

A similar calculation applies for  $g \searrow \text{ and } g^{-1}(y) = x \text{ bijective on } S_X$ .

# II - Data fitting and statistics

### Approximations and convergence

- \* Below a > 0, h(x) > 0 for all  $x \in \mathbb{R}$  and g(x) is a convex function on  $\mathbb{R}$ .
- \* Basic inequality:  $P\{h(X) \ge a\} \le E\{h(X)\}/a$ .
- \* Markov's inequality  $P(|X| \ge a) \le E(|X|)/a$ .
- \* Chebyshev's inequality:  $P(|X| \ge a) \le E(X^2)/a^2$  or  $P\{|X E(X)| \ge a\} \le \frac{\text{Var}(X)}{a^2}$
- $\star$  Jensen's inequality:  $g\{E(X)\} \leq E\{g(X)\}.$
- \* Quadratic mean convergence  $X_n \stackrel{2}{\to} X \iff \lim_{n \to \infty} \mathbb{E}\left((X_n X)^2\right) = 0$  with  $E(X_n^2), E(X^2) < \infty$
- $\star \ \, \text{Convergence in probability:} \ \, X_n \overset{P}{\to} X \ \iff \forall \varepsilon > 0, P\left(\lim_{n \to \infty} |X_n X| > \varepsilon\right) = 0$
- $\star \text{ Convergence in distribution/law: } X_n \overset{D}{\to} X \iff \lim_{n \to \infty} F_{X_n}(x) = F_X(x) \text{ at all } x$  where F(x) is continuous.
- $\star \ X_n \stackrel{2}{\to} X \implies X_n \stackrel{P}{\to} X \implies X_n \stackrel{D}{\to} X$
- $\star \ X_n \overset{D}{\to} x_0 \implies X_n \overset{P}{\to} x_0 \text{ for constant } x_0$
- $\star \ X_n \overset{P}{\to} x_0 \implies g(X_n) \overset{P}{\to} g(x_0)$  with g(x) continuous at  $x_0$ .
- \* Slutsky's lemma:  $X_n \stackrel{D}{\to} X, Y_n \stackrel{P}{\to} y_0 \implies X_n + Y_n \stackrel{D}{\to} X + y_0$  and  $X_n Y_n \stackrel{D}{\to} X y_0$
- \* Law of small numbers: If  $X_n \sim B(n, p_n)$  with  $\lim_{n \to \infty} np_n = \lambda$ , then  $X_n \stackrel{D}{\to} X$  where  $X \sim \text{Poiss}(\lambda)$ .
- \* Weak law of large numbers:  $X_1, \ldots, X_n \overset{\text{iid}}{\sim} F$  where  $E(X_i) = \mu < \infty$   $\Longrightarrow \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \overset{P}{\rightarrow} \mu \left( \Longleftrightarrow \forall \varepsilon > 0, P\left( \lim_{n \to \infty} |X_n \mu| > \varepsilon \right) = 0 \right)$
- \* Strong law of large numbers:  $X_1,\ldots,X_n \overset{\text{iid}}{\sim} F$  where  $E(X_i)=\mu<\infty$   $\Longrightarrow P\left(\lim_{n\to\infty}\overline{X}_n=\mu\right)=1$

### Central limit theorem

- \* CLT:  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} F$  such that  $\mathrm{E}(X_i) = \mu$  and  $0 < \mathrm{Var}(X_i) = \sigma^2 < \infty$   $\implies Z_n = \frac{\overline{X}_n \mu}{\sqrt{\sigma^2/n}} \stackrel{D}{\rightarrow} Z \sim \mathcal{N}(0, 1)$
- $\star \ \, \text{ Delta method: } X_1, \dots, X_n \overset{\text{iid}}{\sim} F \text{ such that } \\ E(X_i) = \mu \text{ and } \\ 0 < \text{Var}(X_i) = \sigma^2 < \infty \\ \text{with } g(x) \text{ such as } g'(\mu) \neq 0 \implies Z_n = \frac{g\left(\overline{X}_n\right) g(\mu)}{\sqrt{g'(\mu)^2\sigma^2/n}} \overset{D}{\rightarrow} Z \sim \mathcal{N}(0,1)$
- \* Quantiles:  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} F, p \in (0, 1)$  where  $x_p = F^{-1}(p)$  and  $f(x_p) > 0$   $\implies Z_n = \frac{X_{(\lceil np \rceil)} x_p}{\sqrt{p(1-p)/n} f(x_p)^2} \stackrel{D}{\rightarrow} Z \sim \mathcal{N}(0, 1)$

# Statistics

- \* A statistic G is a function that depends only on the data  $y = (y_1, \dots, y_n)$ : G = g(y)
- $\star$  A random sample is a set of independent identically distributed data  $Y_1,\ldots,Y_n\stackrel{\mathrm{iid}}{\sim} F$
- \* From the order statistics  $y_{(i) \leq y_{(j)}}$ ,  $1 \leq i \leq j \leq n$ , the empirical quantiles/quantiles of the sample  $y_{(\lceil np \rceil)}$  for  $p \in (0,1)$  can be defined
- $\star$  The breakdown point of a statistic is  $p\times 100\%$ , where  $p\in (0,1)$  is the (asymptotically as  $n\to\infty)$  smallest value such that sending  $x_1,\dots,x_{\lceil np\rceil}\to\pm\infty$  sends the statistic to  $\pm\infty.$

### Notable statistics

- \* Summaries of location: the average (arithmetic mean)  $\overline{y} := n^{-1} \sum_{i=1}^{n} y_i$  and the sample median  $y_{(\lceil n/2 \rceil)}$ .
- \* Summaries of scale / dispersion: the inter-quartile range  $\ \mathrm{IQR} := y_{\lceil \lceil 3n/4 \rceil)} y_{\lceil \lceil n/4 \rceil)}$ , the range  $y_{(n)} y_{(1)}$ , and the sample standard deviation

$$s(y) := \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (y_i - \overline{y})^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (y_i^2 - n\overline{y}^2)}$$

\* Summary of dependence for  $(x_1, y_1), \ldots, (x_n, y_n)$ : the sample correlation coefficient

$$r_{xy} = \frac{\sum_{i=1}^{n} (x_i - \overline{x}) (y_i - \overline{y})}{\left\{\sum_{i=1}^{n} (x_i - \overline{x})^2 \sum_{i=1}^{n} (y_i - \overline{y})^2\right\}^{1/2}}$$

### Useful plots

- **\* Boxplot**: defined by five numbers, the central line  $y(\lceil n/2 \rceil)$ , the limits of the box  $y(\lceil n/4 \rceil)$  and  $y(\lceil 3n/4 \rceil)$ , the limits of the "whiskers", at the  $y_i$  that are most extreme but inside  $y(\lceil n/4 \rceil) 1.5$  IQR and  $y(\lceil 3n/4 \rceil) + 1.5$  IQR. Points outside the whiskers are shown individually.
- **Q-Q plot**: Assuming that  $y_1,\ldots,y_n \stackrel{\text{iid}}{\sim} F$ , the order statistics are plotted against the corresponding quantiles of F, i.e., we plot  $(F^{-1}(i/(n+1)),y_{(i)})$ . A line close to x=y suggests that the data come from F. If F is parametric, a modified plot (depending on F) can be used to estimate some parameters.

### Hypothesis testing

- \* Hypothesis testing is 'proof by stochastic contradiction': we suppose that a null hypothesis H<sub>0</sub> about reality is true and attempt to disprove H<sub>0</sub> using data. Distributions of data Y under H<sub>0</sub> are denoted by subscript 0; these are 'null distributions'.
- \* A test requires a test statistic T = t(Y), large values of which suggest that  $H_0$  is false.
- $\star$  The observed value of  $T,t_{\rm obs},$  is used to compute a p-value  $p_{\rm obs}=P_0$   $(T\geq t_{\rm obs}),$  small values of which cast doubt on  $H_0.$
- If a clear decision is required, a 'significance level'  $\alpha \in (0,1)$  is chosen (e.g.,  $\alpha = 0.05, 0.01, 0.001$ ), and  $H_0$  is rejected iff  $p_{\rm obs} < \alpha$ , or equivalently if  $t_{\rm obs} > t_{1-\alpha}$ , where  $t_{1-\alpha}$  is the  $1-\alpha$  quantile of the null distribution of T.

In a decision setting the possible outcomes of a statistical test are:

	State of nature	
Decision on $H_0$	$H_0$ is true	$H_0$ is false
Not rejected (negative)	True negative	False negative (type II error)
Rejected (positive)	False positive (type I error)	True positive

- $\star$  We may have a clearly-specified alternative/counter hypothesis  $H_1$  that is true when  $H_0$  is false.
- \* With  $H_1$  and  $\alpha$  specified, the probability of a true positive is  $P(\text{rejecting } H_0 \text{ at significance level } \alpha \text{ when } H_1 \text{ is true}) = P_1(T \geq t_{1-\alpha}), \text{ where } \alpha = P_0(T \geq t_{1-\alpha}) \text{ is the probability of a false positive.}$
- $\star \alpha$  is called the size and  $P_1(T \geq t_{1-\alpha}) =: \beta(\alpha)$  the power of the test.
- $\star$  A test is said to be optimal if it maximizes  $\beta(\alpha)$  for all  $\alpha \in (0,1)$
- \* Pearson's statistic When data follow a multinomial distribution (under a certain hypothesis) with denominator n and k categories (it models an experiment with k possible outcomes repeated independently n times, generalising the binomial law), then Pearson's statistic

$$T=\sum_{i=1}^k\frac{(O_i-E_i)^2}{E_i} \text{ follows a } \chi^2_{k-1} \text{ if } \sum E_i/k \geq 5. \text{ This is widely used for tests of fit.}$$

#### Point estimation

- \* An estimator of the parameter  $\theta$  of a parametric model is a function of the data T=t(Y) that estimates  $\theta$ . An estimate is a specific value t=t(y) of T=t(Y).
- \* The bias of an estimator  $\tilde{\theta}$  is  $b_{\tilde{\theta}(\theta)} = E(\tilde{\theta}) \theta$ .
- \* The mean square error of an estimator  $\tilde{\theta}$  is  $MSE_{\tilde{\theta}(\theta)} = E\left(\tilde{\theta} \theta\right)^2 = b_{\tilde{\theta}(\theta)}^2 + Var(\tilde{\theta})$
- $\star$  For two unbiased estimators of  $\theta$ ,  $\tilde{\theta}_1$  and  $\tilde{\theta}_2$ , we say that  $\tilde{\theta}_1$  is more efficient than  $\tilde{\theta}_2$  if  $Var(\tilde{\theta}_1) \leq Var(\tilde{\theta}_2)$

# Types of estimators

- \* Moment estimator: For  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta \in \mathbb{R}^p$ , and moments  $\mathrm{E}(Y_j^r) = \mu_r(\theta)$  for  $r \leq p$ :  $\frac{1}{n} \sum_{i=1}^n Y_i^r \stackrel{n \to \infty}{\to} \mu_r(\theta)$ . This gives a set of p equations in  $\theta$  whose solution gives an estimator for  $\theta$ .
- \* Maximum likelihood: For  $Y_1, \ldots, Y_n \overset{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta \in \mathbb{R}^p$ , and  $Y = (Y_1, \ldots, Y_n) \sim f_Y(\theta)$  the likelihood function  $L(\theta) = f_Y(y, \theta) = \prod_{i=1}^n f(y_i, \theta)$  and the log-likelihood  $l(\theta) = \log L(\theta)$  functions can be defined.

The value  $\hat{\theta}$  such that  $L(\hat{\theta}) \geq L(\theta)$  (or  $l(\hat{\theta}) \geq l(\theta)$ )  $\forall \theta$  is the maximum likelihood estimator. From this, the observed information  $J(\theta) = -\frac{\mathrm{d}^2 l(\theta)}{\mathrm{d}\theta^2}$  and expected / Fisher information  $I(\theta) = \mathrm{E}(J(\theta))$  can be defined for later use.

#### Interval estimation

- \* If  $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta \in \mathbb{R}^p$ , a confidence interval for  $\theta$  is a statistic in the form of an interval that contains  $\theta$  with a given probability (called the confidence level of the interval).
- \* An interval of the form (L,U) is called bilateral and an interval of the form  $(-\infty,U)$  or  $(L,+\infty)$  is called unilateral.
- \* If  $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta \in \mathbb{R}^p$  and  $\tilde{\theta}$  is an estimator of  $\theta$  with V an estimator of  $\text{Var}(\tilde{\theta})$ . Then  $V^{1/2}$  is called a standard deviation of  $\tilde{\theta}$

# Construction of an interval

- \* Using the CLT: If  $Y_1,\ldots,Y_n\stackrel{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta\in\mathbb{R}^p$  and if  $\tilde{\theta}$  is an estimator of  $\theta$  with a standard deviation  $V^{1/2}$  with  $\tilde{\theta}\sim\mathcal{N}(\theta,V)$ , then  $(L,U)=\left(\tilde{\theta}-V^{1/2}z_{1-\alpha_L},\tilde{\theta}-V^{1/2}z_{\alpha_U}\right)$  is a confidence interval with approximate confidence level  $1-\alpha_L-\alpha_U$ .
  - For a bilateral interval we usually chose  $\alpha_L=\alpha_U=\alpha/2$  to have a symmetrical interval. For a unilateral interval we chose  $\alpha_L=\alpha_U=\alpha$  and replace the unwanted limit by  $\pm\infty$ .
- **Limit law of the MLE**: If  $Y_1, \ldots, Y_n \overset{\text{iid}}{\sim} f_{Y_1}(\theta)$  where  $\theta \in \mathbb{R}^p$  and if  $\hat{\theta}$  is the maximum likelihood estimator of  $\theta$ , then under mild regularity conditions,  $J(\hat{\theta})^{1/2}(\hat{\theta}-\theta) \overset{D}{\to} \mathcal{N}_p(0,I_p)$ . We then use the method above.

By Adam Avedissian Content based on the lecture slides of Professor Anthony Davison Format based on the template added by Fingal Mychkine Nagel Persoud