

Some authors use the term *σ -field* to refer to a σ -algebra.

with a discrete sample space, we will often work with the set of all subsets as the relevant σ -algebra. The set of all subsets is sometimes referred to as the **power set** and written 2^S . If the sample space is continuous (like the set of all real numbers),

probability space. With a discrete sample space, we will often work with the

if $P(A) \geq 0$ for all $A \in \mathcal{A}$; $P(\emptyset) = 0$; $P(S) = 1$; and if A_1, A_2, \dots is a series of **disjoint** sets ($A_i \cap A_j = \emptyset$ when $i \neq j$), then $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$. The triple (S, \mathcal{A}, P) is called a

this last requirement is what restricts \mathcal{A} to a σ -algebra.¹

A σ -algebra \mathcal{A} , is a class of subsets of S with three requirements: We require that (i) if $A \in \mathcal{A}$ then $A^c \in \mathcal{A}$; (ii) if $A \in \mathcal{A}$ and $B \in \mathcal{A}$ then $A \cup B$ is also in \mathcal{A} . This defines an algebra. We also require that (iii) if $A_1, A_2, \dots \in \mathcal{A}$, then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$.

Probability theory is defined in terms of set theoretic notions. We begin by defining a **sample space**, which is the set of all possible outcome of some event. Example, role of a die: $S = \{1, 2, 3, 4, 5, 6\}$. We assign probability to **events**, which are subset of S . The set of events that we can assign probabilities to is called a **σ -algebra**, which is a set of subsets of S that satisfies some properties, ensuring that probability measures are well-defined and consistent.

1 Probability

1. $P(A) = 1 - P(A^c)$.
2. $P(A) = P(A \cap B) + P(A \cap B^c)$.
3. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.
4. $P(A \cap B) = P(A)P(B|A)$.
5. $P(A \cap B) \leq \min\{P(A), P(B)\}$.
6. $E[EY|X]] = E[Y]$.
7. $E[a + bX] = a + bE[X]$.
8. $E[aX + bY] = aE[X] + bE[Y]$.
9. $E[g(X)] = \int g(x)f^X(x)dx$ when the density f^X exists.
10. $E[XY] = E[X]E[Y]$ when X and Y are independent.
11. $E[|X + Y|] \leq E[|X|] + E[|Y|]$.
12. $\text{var}[X] = E[X^2] - (E[X])^2$.
13. $\text{var}[\frac{X}{n}] = \frac{\text{var}[X]}{n}$.
14. $\text{cov}[X, Y] = E[XY] - E[X]E[Y]$.
15. $\text{cov}[X, X] = \text{var}[X]$.
16. $\text{var}[a + bX] = b^2 \text{var}[X]$.

Table 1: Standard normal distribution table

	s : combine with z to find cell entry $\Phi(z+s)$									
z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.500	0.504	0.508	0.512	0.516	0.520	0.524	0.528	0.532	0.536
0.1	0.540	0.544	0.548	0.552	0.556	0.560	0.564	0.567	0.571	0.575
0.2	0.579	0.583	0.587	0.591	0.595	0.599	0.603	0.606	0.610	0.614
0.3	0.618	0.622	0.626	0.629	0.633	0.637	0.641	0.644	0.648	0.652
0.4	0.655	0.659	0.663	0.666	0.670	0.674	0.677	0.681	0.684	0.688
0.5	0.691	0.695	0.698	0.702	0.705	0.709	0.712	0.716	0.719	0.722
0.6	0.726	0.729	0.732	0.736	0.739	0.742	0.745	0.749	0.752	0.755
0.7	0.758	0.761	0.764	0.767	0.770	0.773	0.776	0.779	0.782	0.785
0.8	0.788	0.791	0.794	0.797	0.800	0.802	0.805	0.808	0.811	0.813
0.9	0.816	0.819	0.821	0.824	0.826	0.829	0.831	0.834	0.836	0.839
1.0	0.841	0.844	0.846	0.848	0.851	0.853	0.855	0.858	0.860	0.862
1.1	0.864	0.867	0.869	0.871	0.873	0.875	0.877	0.879	0.881	0.883
1.2	0.885	0.887	0.889	0.891	0.893	0.894	0.896	0.898	0.900	0.901
1.3	0.903	0.905	0.907	0.908	0.910	0.911	0.913	0.915	0.916	0.918
1.4	0.919	0.921	0.922	0.924	0.925	0.926	0.928	0.929	0.931	0.932
1.5	0.933	0.934	0.936	0.937	0.938	0.939	0.941	0.942	0.943	0.944
1.6	0.945	0.946	0.947	0.948	0.949	0.951	0.952	0.953	0.954	0.954
1.7	0.955	0.956	0.957	0.958	0.959	0.960	0.961	0.962	0.962	0.963
1.8	0.964	0.965	0.966	0.966	0.967	0.968	0.969	0.969	0.970	0.971
1.9	0.971	0.972	0.973	0.973	0.974	0.974	0.975	0.976	0.976	0.977
2.0	0.977	0.978	0.978	0.979	0.979	0.980	0.980	0.981	0.981	0.982
2.1	0.982	0.983	0.983	0.983	0.984	0.984	0.985	0.985	0.985	0.986
2.2	0.986	0.986	0.987	0.987	0.987	0.988	0.988	0.988	0.989	0.989
2.3	0.989	0.990	0.990	0.990	0.990	0.991	0.991	0.991	0.991	0.992
2.4	0.992	0.992	0.992	0.992	0.993	0.993	0.993	0.993	0.993	0.994
2.5	0.994	0.994	0.994	0.994	0.994	0.995	0.995	0.995	0.995	0.995
2.6	0.995	0.995	0.996	0.996	0.996	0.996	0.996	0.996	0.996	0.996
2.7	0.997	0.997	0.997	0.997	0.997	0.997	0.997	0.997	0.997	0.997
2.8	0.997	0.998	0.998	0.998	0.998	0.998	0.998	0.998	0.998	0.998
2.9	0.998	0.998	0.998	0.998	0.998	0.998	0.998	0.999	0.999	0.999

Note: To calculate $\Phi(x)$ for $x \geq 0$, find the row using the first decimal in x , and find the second decimal to find the correct column. For $x < 0$, use the symmetry of the distribution to find $\Phi(x) = 1 - \Phi(-x)$. The table is calculated with the R-function `pnorm` (R Core Team, 2023).

Compact notes on probability and statistics

Erik Ø. Sørensen*

3rd edition, 2025

Contents

1	Probability	4
2	Random variables	5
2.1	Distribution	6
2.2	Expectation of a random variable	7
2.3	Expectation results	8
2.4	Common distributions	9

*FAIR, Department of Economics, NHH Norwegian School of Economics.
erik.sorensen@nhh.no

This compact set of notes is designed to support the August Method Camp in probability and statistics for incoming PhD students at NHH Norwegian School of Economics. It is not intended to replace a comprehensive textbook or serve as a resource for independent study.

The notes draw heavily on two key references that I have used in teaching over the past decade: Hogg et al. (2013), a solid general textbook in mathematical statistics, and Linton (2017), which is more concise and specifically tailored for training econometricians. For the statistics component, econometrics textbooks will likely suffice for most students. However, for those seeking a deeper theoretical understanding of probability, a good starting point is Rosenthal (2006), while the rigor of Billingsley (1995) should meet the needs of any reader.

Wikimedia can also be a surprisingly helpful resource for topics in probability and statistics.

Although we will not cover software during the Method Camp, I recommend Wickham and Grolemund (2017) as a strong introduction to modern R, particularly for data wrangling in applied settings.

2.5	Functions of random variables	11
3	Random vectors	13
4	Estimation	15
5	Hypothesis testing	18
6	Asymptotic theory	23
6.1	Convergence in probability	23
6.2	Convergence in distribution	24
6.3	The Central Limit Theorem	25
7	Selected maths facts	27
8	Some calculating rules for P, E, var and cov	28
	References	30

If (when) you find a typo, an error, or an inconsistency, please register an issue or a pull request at <https://github.com/ErikOSorensen/CompactNotes>.

- Linton, Oliver (2017). *Probability, Statistics and Econometrics*. London, UK: Academic Press.
- 62(1): 45–53.
- Leemis, Lawrence M. and Jacquelyn T. McQueston (2008). “Univariate distribution relationships.” *American Statistician*, 7th edition.
- Hogg, Robert V., Joseph W. McKean, and Allen T. Craig (2013). *Introduction to Mathematical Statistics*. Pearson Education.
- Haavelmo, Trygve (1944). “The probability approach in econometrics.” *Econometrica*, 12: iii–115.
- Billingsley, Patrick (1995). *Probability and Measure*. Wiley, 3rd edition.
- Benjamin, Daniel J., et al. (2017). “Redefine statistical significance.” *Nature Human Behaviour*, 2: 6–10.

References

17. $\text{var}[aX \pm bY] = a^2 \text{var}[X] + b^2 \text{var}[Y] \pm 2ab \text{cov}[X, Y]$.
18. $\text{cov}[X + Y, Z] = \text{cov}[X, Z] + \text{cov}[Y, Z]$.

- R Core Team (2023). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rosenthal, Jeffrey S. (2006). *A First Look at Rigorous Probability Theory*. World Scientific, 2nd edition.
- Wickham, Hadley and Garrett Grolemund (2017). *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data*. O’Reilly.

This can also be proven with a Taylor approximation. The Δ -rule is the default approach to calculating the distribution of derived statistics. If we can show that the CLT applies

$$\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow{D} N(0, (g'(\theta))^2 \sigma^2).$$

Theorem 3 If $\{X_n\}$ is a sequence of random variables such that $\sqrt{n}(X_n - \theta)$ converges to $N(0, \sigma^2)$ in distribution, g is a differentiable function at θ , and $g'(\theta) \neq 0$, then

says that

The Δ -rule is often useful in conjunction with the CLT. It function of $X - \mu$. second order Taylor-approximation of the moment generating where X has a moment generating function, the trick is to do a simple proof can be constructed for the subset of cases

distribution $N(0, 1)$.

converges in distribution to a random variable with a normal

$$Y_n = \sum_{i=1}^n X_i - n\mu = \sqrt{n}(X_n - \mu) / \sigma$$

random variable

Theorem 2 Let X_1, X_2, \dots, X_n denote a random sample from a distribution with mean μ and positive variance σ^2 . Then the

If the set of outcomes is bounded ($\min(X) > -\infty$ and $\max(X) < \infty$), the expectation always exists. If X is a discrete random variable with pmf p_X , $E[X] = \sum_Y Y p_Y(Y)$. Having different expressions for discrete and continuous random variables is awkward, since sometimes we need to develop theory for

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx.$$

If X is a continuous random variable with pdf $f(x)$ and $\int_{-\infty}^{\infty} |x| f(x) dx < \infty$, the expectation of X exists and is

2.2 Expectation of a random variable

the density function.

For discrete r.v., we define the **probability mass function** (pmf) as $p_X(s) = P(X \in \{s\})$ for $s \in S$. This is analogous to are not necessarily the same random variables.

$F_Y(s)$ for all s , X and Y are **identically distributed**, but they X and Y have distribution functions F_X and F_Y and $F_X(s) = F_Y(s)$ (and sometimes also the limits of integration): $\int f(s) ds = 1$. If 1. When it is clear from the context, we drop the subscripts $f_X(s) \geq 0$, $\int_a^b f_X(s) ds = P(X \in [a, b])$, and $\int_{-\infty}^{\infty} f_X(s) ds = 1$. When the derivative of the distribution function for X exists, it is called the **density** (pdf) of X , written $f_X(s) = F_X'(s)$. We then say X is distributed **continuously**. General properties:

technical requirements: X must be **measurable**:

$$\mathcal{A}_X = \{A \subset S : X(A) \in \mathcal{B}\} = \{X^{-1}(B) : B \in \mathcal{B}\} \subseteq \mathcal{A}.$$

Here \mathcal{B} is a σ -algebra on $S_X \subset \mathbb{R}$ and \mathcal{A} is the σ -algebra on S . Now the probability measure P_X is defined $P_X(B) = P(X^{-1}(B))$, and we have mapped the probability space (S, \mathcal{A}, P) to the probability space (S_X, \mathcal{B}, P_X) .

To recognize the importance of measurability, consider the sample space $S_e = \{a, b\}$, with the σ -algebra $\mathcal{A}_e = \{\emptyset, S_e\}$, and the function $X_e : S_e \rightarrow \mathbb{R}$, taking values $X_e(a) = 0$, $X_e(b) = 1$. This cannot be a random variable, since $X_e^{-1}(0) = \{a\}$ and $X_e^{-1}(1) = \{b\}$, and neither $\{a\}$ nor $\{b\}$ are in the σ -algebra \mathcal{A}_e . Measurability is a requirement on the combination of the function and the σ -algebra.

It is customary to use capital roman letters to refer to random variables and lower case letters to refer to particular values (numbers), such that $X(s) = x$ is the statement that at the point $s \in S$, the random variable X (a function) takes on the value x (a number, $x \in \mathbb{R}$).

2.1 Distribution

A random variable X has a **distribution function** (d.f.) F_X such that

$$F_X(s) = P_X((-\infty, s]) = P(\{a : X(a) \leq s\}).$$

to some moments or parameters, and we are interested in a function of these moments or parameters, we can use the Δ -rule to calculate the distribution of these functions. Statistical packages will often do this automatically (example: Stata's `testnl` command).

7 Selected maths facts

In this section, there are some mathematical facts that turn out to be useful for probability and statistics but are not part of statistics itself.

The **exponential function**, written e^s or $\exp(s)$ can be defined by

$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n,$$

or by the series expansion

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}.$$

Remember that $de^s / ds = e^s$.

If ψ is some function such that $\psi(n) \rightarrow 0$ as $n \rightarrow \infty$, then

$$\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n} + \frac{\psi(n)}{n}\right)^n = e^x,$$

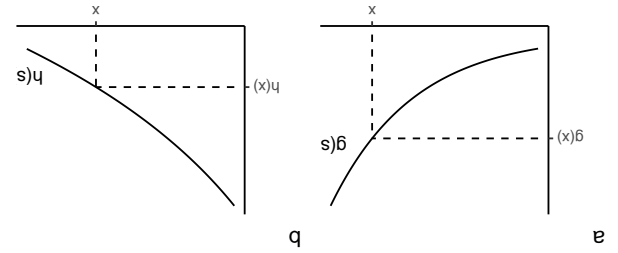
and calculate each term separately. See exercise 11 in the second set of exercises for an example of this.

$$F_Y(y) = P(Y \leq y) = P(Y \leq y, X \in A) + P(Y \leq y, X \in B),$$

probability.

For a function that is not one-to-one, but increasing on the subset A and decreasing on B , we can use the law of total derivative controls for whether g is increasing or decreasing. $f_Y(y) = f_X(g^{-1}(y)) \cdot |dg^{-1}(y)/dy|$. The absolute value of the

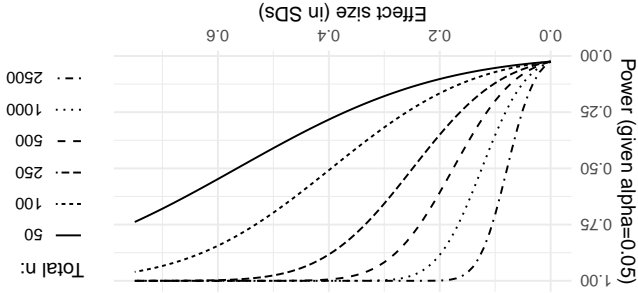
Figure 1: Functions of random variables



alternative hypothesis true, while $1 - p$ is our belief about the encode our prior knowledge as p , our belief that nature made the low power affects statistical inference. Consider Figure 3: we methods, we can uncover some unsettling consequences of how If we analyze classical hypothesis testing with Bayesian most consider sufficient (0.8 and above).

effect sizes require large sample sizes to have the power that size and at different total sample size. Measuring moderate test of $\mu_0 = 0$ at different alternative hypotheses about the effect In Figure 2 we see the power of a two-sided, two-sample t -

Figure 2: Power for a two-sided, two-sample t -test of $H_0 : \mu = 0$ as the effect size and sample size varies ($\alpha = 0.05$). Calculated with the R-command `power.t.test` (R Core Team, 2023).



1. Suppose $\text{plim } X_n = X$ and $\text{plim } Y_n = Y$. Then $\text{plim } (X_n + Y_n) = X + Y$.
2. Suppose $\text{plim } X_n = X$ and a is a constant. Then $\text{plim } aX_n = aX$.
3. Suppose $\text{plim } X_n = a$ and the function g is continuous at a . Then $\text{plim } g(X_n) = g(a)$.
4. Suppose $\text{plim } X_n = X$ and $\text{plim } Y_n = Y$, then $\text{plim } X_n Y_n = XY$.

The concept of probability limit is closely tied to the statistical concept of consistency. Let X be a r.v. with d.f. $F(s, \theta)$, for some $\theta \in \Omega$. Let X_1, X_2, \dots, X_n be a random sample on X , and let T_n be a statistic. T_n is a **consistent** estimator of θ if $\text{plim } T_n = \theta$.

6.2 Convergence in distribution

$\{X_n\}$ is a sequence of random variables and X is a random variable, let F_{X_n} and F_X be the respective distribution functions. Let $C(F_X)$ be the set of points at which F_X is continuous. Now X_n **converges in distribution** to X if $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for all $x \in C(F_X)$, and we sometimes write

$$X_n \xrightarrow{D} X.$$

$k < m$, then $E[X^k]$ exists.

2. Let $u(X)$ be a nonnegative function of the random variable X . If $E[u(X)]$ exists, then for every positive constant c ,

$$P(u(X) \geq c) \leq E[u(X)]/c$$

(Markov).

3. Let the random variable X have a distribution with finite variance σ^2 . Then for every $k > 0$,

$$P(|X - \mu| \geq k\sigma) \leq 1/k^2$$

(Chebyshev).

4. If ϕ is convex on an open interval I and X is a random variable with finite expectation and support in I , then $\phi(E[X]) \leq E[\phi(X)]$ (Jensen).

2.4 Common distributions

There are some distributions that we should recognize:

Uniform A uniform distribution on $[a, b]$ is written $U(a, b)$, has a density $1/(b-a)$ on $[a, b]$ and a d.f. $F(x) = x/(b-a)$ on $[a, b]$. Expectation is $(a+b)/2$, the variance is $(b-a)^2/12$.

In Figure 1b, we start with a different inequality: If $h(X) \leq y$, then $X \geq y_{-1}(y)$, and $f_Y(y) = P(X \geq y_{-1}(y)) = 1 - P(X \leq y_{-1}(y)) = 1 - F_X(y_{-1}(y))$. Either way, for increasing or decreasing transformations g , taking derivatives we find that

another random variable. In Figure 1a, $X = g(X)$ is an increasing function (one-to-one). Considering the distribution of X , the distribution function for X is $F_X(y) = P(X \leq y)$. We can substitute the definition of X into this: $F_X(y) = P(g(X) \leq y)$. We can now apply the inverse of g on both sides of the inequality in P , relying on g being increasing: $P(g(X) \leq y) = P(g^{-1}(g(X)) \leq g^{-1}(y)) = P(X \leq g^{-1}(y))$, and we can conclude that $F_X(y) = F_X(g^{-1}(y))$. The density follows by differ-

If X is a random variable and g is a function, $g(X)$ is

2.5 Functions of random variables

A more complete overview of distributions is given by Leemis and McQuestion (2008). There are also multi-volume handbooks written with properties of various distributions.

distribution function is not pretty. The moment generating function is $m(t) = \exp(\lambda(\exp(t) - 1))$. The expectation is λ and the variance is λ .

Standard normal A standard normal distribution is written $N(0, 1)$, the density is

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

The distribution function is often written Φ but there is no analytical expression for this. See Table 1 in the back for some tabulated values. The moment generating function is $m(t) = \exp(t^2/2)$. The expectation is 0 and the variance is 1. The standard normal distribution is **symmetric** around zero, $\Phi(s) = 1 - \Phi(-s)$ for all $s \in \mathbb{R}$.

Exponential An exponential distribution with parameter λ is written $\text{Exp}(\lambda)$, has density $f(x) = \lambda \exp(-\lambda x)$ on $[0, \infty)$, and distribution function $F(x) = 1 - \exp(-\lambda x)$. The expectation is $1/\lambda$ and the variance is $1/\lambda^2$.

Binomial If an event happens with probability p in a single (binary) trial, the distribution of the number of events k in n trials has probability mass function $p(k) = \binom{n}{k} p^k (1-p)^{n-k}$ (with the binomial coefficient $\binom{n}{k} = n! / k!(n-k)!$). The expectation is np and the variance is $np(1-p)$. The special case $n = 1$ is often very useful (**Bernoulli**).

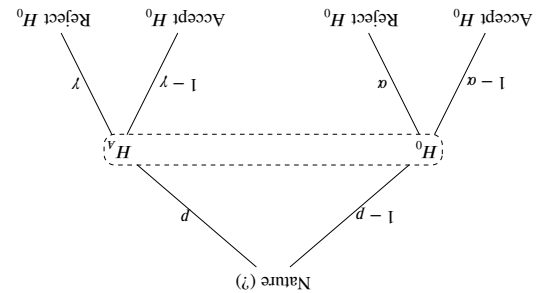
Poisson A Poisson distribution on the integers $\{0, 1, 2, \dots\}$ has probability mass function $p(x) = \lambda^x \exp(-\lambda)/x!$, the

Consider a surprising ($p = 0.1$) rejection of H_0 in an under-powered study ($\gamma = 0.5$), our posterior belief in H_A should not

$$P(H_A | \text{Reject } H_0) = \frac{P(\text{Reject } H_0 | H_A) \cdot P(H_A)}{P(\text{Reject } H_0)} = \frac{(1 - \alpha)d + \gamma}{\gamma}.$$

probability that H_0 is true. We cannot tell if we are at the H_0 or at the H_A node, but if we are at H_0 we reject with probability α , the level of significance for our test (a type-I-error). If we are at H_A , we reject with probability γ , the power of the test. Rejecting H_0 , the Bayesian posterior belief in H_A is

Figure 3: Statistical game tree



be very strong: $0.1 \cdot 0.5 / (0.9 \cdot 0.05 + 0.1 \cdot 0.5) = 0.53$. For this reason some people argue for a much stricter α : Benjamin et al. (2017) argue that we should use $\alpha = 0.005$ as a conventional level instead of 0.05 (in this example, if we had the surprising result at $\alpha = 0.005$, that would give us $P(H_A | \text{Reject } H_0) = 0.92$).

6 Asymptotic theory

6.1 Convergence in probability

$\{X_n\}$ is a sequence of random variables, and X is a random variable, both defined on the same sample space. Now X_n **converges in probability** to X if, for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0.$$

We say that $\text{plim } X_n = X$, or that $X_n \xrightarrow{P} X$.

Theorem 1 *Let $\{X_n\}$ be sequence of iid random variables with common mean μ and finite variance σ^2 . If $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$, then $\text{plim } \bar{X}_n = \mu$. (Weak law of large numbers.)*

There are a number of results about how probability limits work:

Our task is, based on a sample from the distribution of X , determine whether to keep insisting that H_0 is true or reject that in favor of H_A . What we want is often to fix the probability of Type I error ("significance level") to $\alpha = 0.01, 0.05, \dots$, and conditional on that, minimize probability of type II error. By **Type I error** we mean that a null hypothesis is rejected when the null is true; by **Type II error** we mean that the null hypothesis is not rejected even though it is false.

$$H_0 : \theta \in \omega_0 \quad \text{vs} \quad H_A : \theta \in \omega_A.$$

Assume that we have two alternative ways to think of what might be true. The distribution of X is $f(x; \theta)$, with parameter $\theta \in \Omega$, and our ideas about truth can be described by regions of Ω :

5 Hypothesis testing

Inference, which is a crucial input into hypothesis testing. Often we want our estimators to be **unbiased**, such that $E[\hat{\theta}] = \theta_0$, with θ_0 being the true value. We can define the **bias** as $b = E[\hat{\theta}] - \theta_0$. Even if estimators are biased (most nonlinear estimators are biased), they can be **consistent**. To define consistency we need some asymptotic theory (in Section 6).

4 Estimation

For estimation, we would like to have some rules to uncover empirical analogs to theoretically defined parameters. We will assume we have access to a random sample from some defined population (with dependency in data, we often start by a transforming data such that the transformed data is a random sample). The **analog principle** is to start with a theoretically defined parameter and set its empirical analog equal to the same in the

A vector $\mathbf{X} = (X_1, \dots, X_n)$ where each element is an *independent* draw from the same distribution as the random variable X is known as a **random sample** of size n of the variable X . integration with change of variables in calculus.

where $\mathbf{J}(y_1, y_2)$ is the Jacobian of $\mathbf{w}(y)$ and B is the image $B = \mathbf{u}(A)$. This is an application of the general formula for

$$= \iint_A \iint_{X_1, X_2} f_{X_1, X_2}(x_1, x_2) dx_1 dx_2 = \iint_B f_{X_1, X_2}(w_1(y_1, y_2), w_2(y_1, y_2)) |\det(\mathbf{J}(y_1, y_2))| dy_1 dy_2$$

let w_1 and w_2 be inverses. Then in general,

with), we require that for all Borel subsets $S \subset \mathbb{R}$,

$$\begin{aligned} E[P(A|X) \cdot 1_{X \in S}] &= P(A \cap \{X \in S\}), \\ E[E[Y|X] \cdot 1_{X \in S}] &= E[Y \cdot 1_{X \in S}]. \end{aligned}$$

These equations tie down both conditional probability and conditional expectation. Note that with the indicator function, it would be possible to *define* probability as $P(X \in A) = E[1_{X \in A}]$, this would push the problem of defining conditional probability into the theory of integration.

Let (X_1, X_2) have a joint distribution function $F(x_1, x_2)$ and let X_1 and X_2 have marginal distribution functions $F_1(x_1)$ and $F_2(x_2)$. Then X_1 and X_2 are **independent** if and only if

$$F(x_1, x_2) = F_1(x_1)F_2(x_2),$$

the joint distribution function being the product of marginal distribution functions. It follows that if the corresponding densities exist, X_1 and X_2 are independent if the joint density is the product of marginal densities.

If X_1 and X_2 are independent and $E[u(X_1)]$ and $E[v(X_2)]$ exist, then it is the case that $E[u(X_1)v(X_2)] = E[u(X_1)] \cdot E[v(X_2)]$.

If X and Y are random variables, $E[E[Y|X]] = E[Y]$, known as the law of **iterated expectations**.

For complete generality, let $Y_1 = u_1(X_1, X_2)$ and $Y_2 = u_2(X_1, X_2)$ be one-to-one transformed random variables, and

Strictly speaking, we decide on a regions of the outcome space, $(X_1, \dots, X_n) \in C$, in which to accept or reject H_0 , but in practice we define test-statistics $T_n(X_1, \dots, X_n)$ and decide on **critical regions** of the test statistics. Instead of reporting reject/accept with a given α -criterion (often indicated with stars in tables), papers often report p -values.

P-values are the probability of at least as extreme data given H_0 ,

$$P(\text{data}|H_0).$$

P -values do not address the likelihood of H_0 being true. In classical inference, probabilities are *never* attached to hypotheses: Hypotheses are true or false, and this is not a sampling issue. If we want to be Bayesian about hypotheses (have beliefs about them), we need to incorporate the priors:

$$P(H_0|\text{data}) = \frac{P(\text{data}|H_0)P(H_0)}{P(\text{data})}.$$

Very little of applied economics is explicitly Bayesian, most is anchored in **frequentist inference**. In frequentist philosophy, there is always a true (but unknown) value of any parameter, and we do not put probabilities on our beliefs about this parameter. Probabilities are restricted to the **sampling distribution** of our estimators *under some null hypothesis*.

If we assume that X is normally distributed, we could test