

Probability Theory

Kexing Ying

February 2, 2022

Contents

1	Review of Measure Theory	2
2	Random Variables	6
2.1	Inequalities	9
2.2	Transformation of Random Variables	11
2.3	Independence	12
2.4	Weak Law of Large Numbers	13

1 Review of Measure Theory

Modern probability theory is based on measure theory and we will in this section recall some notions from measure theory.

Definition 1.1 (Algebra). Given a set Ω , a set of subsets \mathcal{A} of Ω is an algebra if $\Omega \in \mathcal{A}$ and \mathcal{A} is closed under finite union and complements.

It follows straight away that an algebra is also closed under finite intersections.

Definition 1.2 (Finitely Additive Measure). A function $\mu : \mathcal{A} \rightarrow [0, \infty]$ where \mathcal{A} is an algebra, is a finitely additive measure if for any disjoint sets $A, B \in \mathcal{A}$,

$$\mu(A \cup B) = \mu(A) + \mu(B).$$

Definition 1.3 (σ -Algebra). A σ -algebra \mathcal{F} is an algebra that is closed under countable unions.

Similarly, it follows that \mathcal{F} is closed under countable intersections.

Definition 1.4 (Measure). A function $\mu : \mathcal{F} \rightarrow [0, \infty]$ where \mathcal{F} is a σ -algebra, is a σ -additive measure (or simply measure) if given a sequence of pairwise disjoint sets A_1, A_2, \dots of \mathcal{F} , we have

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i).$$

We call a measure a probability measure if $\mu(\Omega) = 1$.

Definition 1.5 (σ -Finite Measure). A measure μ is said to be σ -finite if there exists a sequence of pairwise disjoint sets A_1, A_2, \dots of \mathcal{F} , such that $\bigcup_{i=1}^{\infty} A_i = \Omega$ and for all i , $\mu(A_i) < \infty$.

Definition 1.6 (Probability Space). A probability space is the triple $(\Omega, \mathcal{F}, \mathbb{P})$ consisting of a set Ω , a σ -algebra \mathcal{F} on Ω and \mathbb{P} a probability measure on \mathcal{F} .

We call elements of \mathcal{F} (i.e. a \mathcal{F} -measurable set) an event.

Proposition 1.1 (Continuity of Measures). Let $(A_n)_{n \in \mathbb{N}} \subseteq \mathcal{F}$, then

- (continuity from below) if (A_n) is increasing, then

$$\mathbb{P}\left(\bigcup_{n=1}^{\infty} A_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n).$$

- (continuity from above) if (A_n) is decreasing, then

$$\mathbb{P}\left(\bigcap_{n=1}^{\infty} A_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n).$$

We recall the the finiteness of the measure is vital for continuity from below while continuity from above is also valid for general measures.

Proof. Exercise. □

Proposition 1.2. A finitely additive probability measure on the σ -algebra \mathcal{F} is a probability measure if and only if it is continuous at 0.

Proof. The forward direction follows from above so we will prove the reverse. Suppose μ is finitely additive and for any decreasing $(A_n) \subseteq \mathcal{F}$ with $\bigcap A_n = \emptyset$, we have $\lim_{n \rightarrow \infty} \mu(A_n) = 0$. Then, μ is continuous from below, and so for any sequence of disjoint sets (B_n) , we have $(C_n) := (\bigcup_{i=1}^n B_i)$ is a sequence of increasing sets and thus,

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} B_i\right) = \mathbb{P}\left(\bigcup_{i=1}^{\infty} C_i\right) = \lim_{n \rightarrow \infty} \mathbb{P}(C_n) = \lim_{n \rightarrow \infty} \mathbb{P}\left(\bigcup_{i=1}^n B_i\right) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{P}(B_i)$$

implying μ is σ -additive and so, μ is a measure. □

Proposition 1.3. Given a collection $\{\mathcal{F}_i\}_{i \in I}$ σ -algebras of Ω , $\bigcap_{i \in I} \mathcal{F}_i$ is also a σ -algebra on Ω .

Definition 1.7 (σ -Algebra Generated By Sets). Given a collection of subsets S of Ω , the σ -algebra generated by S is

$$\sigma(S) := \bigcap \{\mathcal{F} \text{ a } \sigma\text{-algebra} \mid S \subseteq \mathcal{F}\}.$$

Definition 1.8 (Borel σ -Algebra). Given a topological space (X, \mathcal{T}) , the Borel σ -algebra on X is $\mathcal{B}(X) := \sigma(\mathcal{T})$.

Definition 1.9 (Product σ -Algebra). Given measurable spaces $(\Omega_1, \mathcal{F}_1), (\Omega_2, \mathcal{F}_2)$, the product σ -algebra on $\Omega_1 \times \Omega_2$ is

$$\mathcal{F}_1 \otimes \mathcal{F}_2 := \sigma(\mathcal{F}_1 \times \mathcal{F}_2) = \sigma(\{A_1 \times A_2 \mid A_1 \in \mathcal{F}_1, A_2 \in \mathcal{F}_2\}).$$

Definition 1.10 (Cylindrical σ -Algebra). A set $C \subseteq \mathbb{R}^\infty$ is said to be cylindrical if is of the form

$$C = \{x \in \mathbb{R}^\infty \mid (x_1, \dots, x_n) \in C_n\}$$

where $C_n \in \mathcal{B}(\mathbb{R}^n)$. The set of cylindrical sets $\mathcal{B}(\mathbb{R}^\infty)$ form a σ -algebra on \mathbb{R}^∞ and is called the cylindrical σ -algebra.

Definition 1.11 (Consistent). The sequence of measures \mathbb{P}_n on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ is said to be consistent if for all $n \in \mathbb{N}$, $\mathbb{P}_{n+1}(B_n \times \mathbb{R}) = \mathbb{P}_n(B_n)$ for all $B_n \in \mathcal{B}(\mathbb{R}^n)$.

Theorem 1 (Kolmogorov). Given any consistent sequence of measures \mathbb{P}_n , there exists a unique probability measure \mathbb{P} on $(\mathbb{R}^\infty, \mathcal{B}(\mathbb{R}^\infty))$ such that,

$$\mathbb{P}(\{x \in \mathbb{R}^\infty \mid (x_1, \dots, x_n) \in B_n\}) = \mathbb{P}_n(B_n)$$

for all $n \geq 1$, $B_n \in \mathcal{B}(\mathbb{R}^n)$.

Proof. Simply define the inner measure on the generating sets as claimed and use the Caratheodory extension (which provides both existence and uniqueness). □

Recall that a nondecreasing function g on \mathbb{R} is continuous up to possibly countably many discontinuities of the first kind. Furthermore, the derivative g' exists λ -a.e. (where λ is the Lebesgue measure on \mathbb{R}).

Proposition 1.4. Let $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbb{P})$ be a probability space. Defining $F(x) := \mathbb{P}(-\infty, x]$, we have

- F is nondecreasing;
- $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$;
- F is continuous on the right.

Proof. Clear by the monotonicity, continuity of measures (from above). \square

Definition 1.12 (Distribution Function). Any function $F : \mathbb{R} \rightarrow [0, 1]$ satisfying the above three properties is said to be a distribution function on \mathbb{R} .

It is clear that any probability measure induces a distribution. On the other hand the converse is also true.

Proposition 1.5. Given a distribution function F , there exists a unique probability measure \mathbb{P} on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ such that $F(x) = \mathbb{P}(-\infty, x]$ for all $x \in \mathbb{R}$.

Proof. Use Caratheodory extension theorem on the algebra $\{(-\infty, x] \mid x \in \mathbb{R}\}$ mapping $(-\infty, x] \mapsto F(x)$. The uniqueness of the probability measure follows by the uniqueness of the Caratheodory extension. \square

Definition 1.13 (Null-set). Given a measure μ , a set $S \subseteq \Omega$ is a null-set if there exists some measurable set $N \subseteq \Omega$ with measure 0 such that $S \subseteq N$.

Definition 1.14 (Complete Measure). A measure μ is complete if every μ -null set is measurable.

If a measure on the σ -algebra Σ is not complete, we may complete the σ -algebra by extending Σ to

$$\overline{\Sigma} := \sigma(\Sigma \cup \{N \mid N \text{ is a null-set}\}).$$

Clearly, the null-sets will have measure 0.

We note that the probability space $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathbb{P})$ is not complete as there exists subsets of a Borel null-set which are not Borel. With this in mind, we denote the completion of $\mathcal{B}(\mathbb{R})$ by $\mathcal{M}(\mathbb{R})$ and we say \mathbb{P} is the Lebesgue-Stieltjes measure.

Recall, that in elementary probability theory, we considered three types of distributions, namely discrete, absolutely continuous and singular continuous. Let us now consider them again in a formal measure theoretic setting.

- (Discrete) A random variable $X : \Omega \rightarrow \mathbb{R}$ is said to have discrete distribution if there exists some countable (including finite) set $A \subseteq \mathbb{R}$, such that for all $E \in \mathcal{B}(\mathbb{R})$, the push-forward measure satisfies

$$X_*\mathbb{P}(E) = \sum_{x \in A} p(x)\delta_x(E)$$

where $p(x) = X_*\mathbb{P}(\{x\})$ and δ_x is the Dirac measure at x .

We note that a distribution function F corresponds to a discrete random variable if and only if for all $x_0 \in \mathbb{R}$,

$$F(x_0) = \sum_{x \in A \cap \{\leq x_0\}} p(x).$$

It is clear that $\sum_A p(x) = 1$ since $\sum_A p(x) = X_*\mathbb{P}(\mathbb{R}) = 1$.

- (Absolutely continuous) A random variable $X : \Omega \rightarrow \mathbb{R}$ is said to be absolutely continuous if $X_*\mathbb{P} = f\lambda$ for some Lebesgue integrable function f and λ denotes Lebesgue measure. Thus, the distribution function corresponding to X satisfies

$$F(x) = \int_{(-\infty, x]} f d\lambda.$$

In particular, recalling the Radon-Nikodym theorem, we have X is absolutely continuous if and only if $X_*\mathbb{P} \ll \lambda$ (hence the name “absolutely continuous”).

Before introducing the last type of distribution, let us consider the following definition.

Definition 1.15 (Concentrated). A measure μ on the measurable space X is said to be concentrated on a measurable set A if $\mu(E) = 0$ for all $E \subseteq X \setminus A$.

- (Singular continuous) A random variable $X : \Omega \rightarrow \mathbb{R}$ is singular continuous if its distribution function F is continuous and $X_*\mathbb{P}$ is concentrated on a set A of Lebesgue measure 0 for which $F'(x) = 0$ for all $x \in A$ almost everywhere.

We note that, since $\lambda(A) = 0$, $X_*\mathbb{P} \perp \lambda$ by the set A (hence the name “singular”). Moreover, by continuity, $X_*\mathbb{P}(\{x\}) = 0$ for all $x \in \mathbb{R}$ in contrast to the discrete measure.

Analogous to the Lebesgue decomposition of measures, we may decompose any distribution function into a discrete, absolutely continuous and singular continuous distribution.

Theorem 2 (Hahn Decomposition for Distributions). Given a distribution function F , there exists $a_1 + a_2 + a_3 = 1$ and $F_{\text{disc}}, F_{\text{ac}}, F_{\text{sc}}$ discrete, absolutely continuous and singular continuous distribution functions respectively, such that

$$F = a_1 F_{\text{disc}} + a_2 F_{\text{ac}} + a_3 F_{\text{sc}}.$$

Proof. Recalling the refinement of the Lebesgue decomposition where we may decompose a measure μ with

$$\mu = \mu_d + \mu_a + \mu_s$$

where μ_d is a discrete measure, $\mu_a \ll \lambda$ and μ_s is singular continuous (i.e. mutually singular with respect to the Lebesgue measure and $\mu_s\{x\}$ for all x). Thus, by simply taking the decomposition of the measure corresponding to F (i.e. $X_*\mathbb{P}$), we obtain the required decomposition after normalization. \square

2 Random Variables

We will continue to let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space.

Definition 2.1 (Random Variable). A function $\xi : \Omega \rightarrow \mathbb{R}$ is said to be a random variable if it is \mathcal{F} -measurable (i.e. for any $B \in \mathcal{B}(\mathbb{R})$, we have $\xi^{-1}(B) \in \mathcal{F}$).

While we have already introduced the notion of a distribution within the previous section, we will present it here again for organization.

Definition 2.2 (Distribution of a Random Variable). Given a random variable ξ , the distribution of ξ is the push-forward measure of \mathbb{P} along ξ . Furthermore, the distribution function corresponding to ξ is

$$F(x) := X_*\mathbb{P}(-\infty, x].$$

Definition 2.3. Given a random variable ξ , we define $\mathcal{F}_\xi \subseteq \mathcal{F}$ to be the σ -algebra

$$\mathcal{F}_\xi := \{\xi^{-1}(B) \mid B \in \mathcal{B}(\mathbb{R})\}.$$

This is the least σ -algebra for which ξ is measurable.

We will recall some standard results about measurable functions. All proofs are left as exercises and can be found in the second year measure theory notes.

Lemma 2.1. If $\mathcal{B}(\mathbb{R}) = \sigma(\mathcal{D})$ for a collection of sets \mathcal{D} , ξ is a random variable if $\xi^{-1}(D) \in \mathcal{F}$ for all $D \in \mathcal{D}$.

Lemma 2.2. Given random variables f, g and $c \in \mathbb{R}$, $f + g, f - g, c \cdot f, |f|, fg, \max(f, g)$, and $\min(f, g)$ are all random variables. Furthermore, if $g(x) \neq 0$ for all x , then f/g is also a random variable.

Lemma 2.3. If (f_n) is a sequence of random variables, then

$$\sup_n f_n, \inf_n f_n, \lim_n f_n$$

are random variables if they exist.

Lemma 2.4. If ξ is a random variable and $f : \mathbb{R} \rightarrow \mathbb{R}$ is continuous, then $f(\xi)$ is a random variable.

Definition 2.4 (Simple Function). A random variable ξ is simple if there exists a partition of Ω , D_1, \dots, D_n such that

$$\xi(\omega) = \sum_{i=1}^n x_i 1_{D_i}(\omega)$$

for some x_1, \dots, x_n for all $\omega \in \Omega$.

Lemma 2.5. For any non-negative random variable ξ , there exists a sequence of nondecreasing simple random variables (ξ_n) such that for all $\omega \in \Omega$,

$$\xi_n(\omega) \uparrow \xi(\omega).$$

Definition 2.5 (Random Vector). A function $\xi : \Omega \rightarrow \mathbb{R}^n$ is a random vector if it is measurable. Again, we define its distribution to be its push-forward measure.

Lemma 2.6. $\xi : \Omega \rightarrow \mathbb{R}^n$ is a random vector if and only if $\xi_i := \text{pr}_i \circ \xi$ is a random variable for all $i = 1, \dots, n$ (where $\text{pr}_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is the i -th projection function).

Definition 2.6 (Independent Random Variables). Two random variables $\xi, \eta : \Omega \rightarrow \mathbb{R}$ are said to be independent if

$$(\xi, \eta)_* \mathbb{P} = \xi_* \mathbb{P} \otimes \eta_* \mathbb{P}.$$

Since, to check that two measures are equal, it suffices to check equality on generating sets, ξ, η are independent if for all $A, B \in \mathcal{B}(\mathbb{R})$,

$$\mathbb{P}(\xi \in A, \eta \in B) = \mathbb{P}(\xi \in A) \mathbb{P}(\eta \in B).$$

Let us quickly recall the construction of the Lebesgue integral.

1. Define the Lebesgue integral for simple functions.
2. Define the Lebesgue integral for non-negative functions by taking the limit of the Lebesgue integral of the monotone sequence of simple functions which converge to the said function.
3. Define the Lebesgue integral for general real-valued functions f by taking $\int f = \int f^+ - \int f^-$ if $\int |f| < \infty$.

Definition 2.7 (Expectation). Given a random variable $\xi : \Omega \rightarrow \mathbb{R}$, the expectation of ξ is simply

$$\mathbb{E}(\xi) := \int \xi d\mathbb{P}$$

if it exists. Furthermore, we say ξ is integrable if $\mathbb{E}(|\xi|) < \infty$.

Proposition 2.1. Let ξ, η be integrable random variables and let $c \in \mathbb{R}$, then

- $\mathbb{E}(c) = c$;
- $\mathbb{E}(\xi + \eta) = \mathbb{E}(\xi) + \mathbb{E}(\eta)$;
- $\xi \leq \eta$ a.e. implies $\mathbb{E}(\xi) \leq \mathbb{E}(\eta)$;
- $\xi = \eta$ a.e. implies $\mathbb{E}(\xi) = \mathbb{E}(\eta)$;
- $\xi \geq 0$ a.e. and $\mathbb{E}(\xi) = 0$ implies $\xi = 0$ a.e.

Proof. Follows directly from the properties of the Lebesgue integral. □

Let us recall some convergence theorems for the Lebesgue integral.

Theorem 3 (Dominated Convergence Theorem). Let (ξ_n) be a sequence of random variables such that $\xi_n \rightarrow \xi$ almost everywhere. If there exists some integrable η such that $|\xi_n| \leq \eta$ for all n , then, ξ is integrable and

$$\lim_{n \rightarrow \infty} \mathbb{E}(\xi_n) = \mathbb{E}(\xi).$$

Theorem 4 (Monotone Convergence Theorem). Let (ξ_n) be a sequence of non-negative increasing random variables. Then,

$$\lim_{n \rightarrow \infty} \mathbb{E}(\xi_n) = \mathbb{E} \lim_{n \rightarrow \infty} \xi_n.$$

We note that the right hand side limit always exists since for all $\omega \in \Omega$, $\xi_n(\omega)$ is increasing any bounded by ∞ .

We remark that the monotone convergence theorem applies if there exists some random variable η such that $\mathbb{E}(\eta) > -\infty$ such that $\eta \leq \xi_n$ for all n by considering $\xi_n - \eta$.

Corollary 4.1. If (η_n) is a sequence of non-negative random variables, then

$$\sum_{i=1}^{\infty} \mathbb{E}(\eta_i) = \mathbb{E} \left(\sum_{i=1}^{\infty} \eta_i \right).$$

Corollary 4.2 (Fatou's lemma). Let ξ_n be a sequence of non-negative random variables. Then,

$$\mathbb{E}(\liminf_n \xi_n) \leq \liminf_n \mathbb{E}\xi_n.$$

Proof. Apply the monotone convergence theorem to $\lambda_n := \inf_{k>n} \xi_k$. \square

Again, the non-negative condition can be replaced by the existence of some random variable η such that $\mathbb{E}(\eta) > -\infty$ and $\eta \leq \xi_n$ for all n . On the other hand, if $\mathbb{E}(\eta) < \infty$ and $\xi_n \leq \eta$, the theorem holds with limit supremum instead.

We note that in all above theorems, the statement still holds by replacing Ω with any measurable set by restricting the measure onto that set.

Theorem 5 (Change of Variables). Given a random variable ξ , a measurable function $g : \mathbb{R} \rightarrow \mathbb{R}$ and a measurable set A , we have

$$\int_A g d\xi_* \mathbb{P} = \int_{\xi^{-1}(A)} g \circ \xi d\mathbb{P},$$

where both integrals either exist or not exist simultaneously.

Proof. Apply usual method where one first prove the statement for indicator functions. Then, it follows that it holds for simple functions by the linearity of the integral. Finally, for any non-negative measurable function, we take a sequence of monotonically increasing simple functions, and apply the monotone convergence theorem. For arbitrary functions, the result follows by taking $f = f^+ - f^-$. \square

Corollary 5.1 (Law of the Unconscious Statistician). Given a random variable ξ and a measurable function $g : \mathbb{R} \rightarrow \mathbb{R}$,

$$\mathbb{E}(g(\xi)) = \int_{\mathbb{R}} g d\xi_* \mathbb{P}.$$

Corollary 5.2. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be measurable, then, if ξ be a discrete random variable,

$$\mathbb{E}(g(\xi)) = \sum_{x \in A} g(x)p(x).$$

On the other hand, if ξ is absolutely continuous, i.e. there exists some f such that $f\lambda = \xi_* \mathbb{P}$, then

$$\mathbb{E}(g(\xi)) = \int_{\mathbb{R}} g(x)f(x)\lambda(dx).$$

Proof. In the discrete case, we have

$$\mathbb{E}(g(\xi)) = \int g d \left(\sum_{x \in A} p(x) \delta_x \right) = \sum_{x \in A} p(x) \int g d \delta_x.$$

By considering $\int g d \delta_x = \int_{\{x\}} g d \delta_x + \int_{\mathbb{R} \setminus \{x\}} g d \delta_x = \delta_x(\{x\})g(x) + 0 = g(x)$. We have

$$\mathbb{E}(g(\xi)) = \sum_{x \in A} g(x)p(x),$$

as required.

On the other hand, if $f\lambda = \xi_*\mathbb{P}$, we have

$$\mathbb{E}(g(\xi)) = \int_{\mathbb{R}} g d(f\lambda) = \int_{\mathbb{R}} g(x)f(x)\lambda(dx)$$

as required. \square

Theorem 6 (Fubini's Theorem). Let $(E_1, \Sigma_1, \mu_1), (E_2, \Sigma_2, \mu_2)$ be σ -finite measure spaces. Then, for any $\Sigma_1 \otimes \Sigma_2$ -measurable functions $g : E_1 \times E_2 \rightarrow \mathbb{R}$, $g(\cdot, y_0)$ is Σ_1 -measurable for all $y_0 \in E_2$, and $g(x_0, \cdot)$ is Σ_2 -measurable for all $x_0 \in E_1$. Furthermore, $\int_{E_1} g d\mu_1, \int_{E_2} g d\mu_2$ are Σ_2 and Σ_1 -measurable respectively. Finally, if $\int |g| d\mu_1 \otimes \mu_2 < \infty$, then,

$$\int g d\mu_1 \otimes \mu_2 = \int \left(\int g(x, y) \mu_2(dy) \right) \mu_1(dx) = \int \left(\int g(x, y) \mu_1(dx) \right) \mu_2(dy).$$

2.1 Inequalities

Lemma 2.7 (Jensen's Inequality). Let ξ be an integrable random variable and let $g : \mathbb{R} \rightarrow \mathbb{R}$ be a measurable, convex function, then,

$$g(\mathbb{E}\xi) \leq \mathbb{E}g(\xi).$$

Proof. Recall that the function g is convex if for all $x_0 \in \mathbb{R}$, there exists some λ such that $g(x) \geq g(x_0) + (x - x_0)\lambda$ (graphically, λ is the slope (more accurately, a subderivative) of g at x_0 and so, the inequality is saying that the graph lies above the tangent line).

Setting $x = \xi$ and $x_0 = \mathbb{E}\xi$. Then, the above inequality becomes

$$g(\xi) \geq g(\mathbb{E}\xi) - (\xi - \mathbb{E}\xi)\lambda.$$

Thus, applying the expectation to both sides results in the required inequality by the linearity of the integral. \square

Corollary 6.1 (Lyapunov's Inequality). Let ξ be a random variable and let $0 < s < t$ be real numbers, then

$$\mathbb{E}(|\xi|^s)^{1/s} \leq \mathbb{E}(|\xi|^t)^{1/t}.$$

Proof. Use Jensen's inequality with $g(x) = |x|^{t/s}$.

Alternatively, setting $\eta = |\xi|^s$, by Hölder's inequality, we have

$$\|\xi^s\|_1 = \|\eta\|_1 \leq \|\eta\|_{t/s} = \|\xi\|_t^s.$$

Thus, taking both sides to the power of $1/s$, we obtain $\|\xi\|_s = (\|\xi^s\|_1)^{1/s} \leq (\|\xi\|_t^s)^{1/s} = \|\xi\|_t$ as required. \square

Proposition 2.2 (Markov Inequality). Let $\xi \geq 0$ be an integrable random variable and let $c > 0$. Then

$$\mathbb{P}(\xi \geq c) \leq \frac{\mathbb{E}(\xi)}{c}.$$

Proof. $\mathbb{E}(\xi) \geq \mathbb{E}(\xi \mathbf{1}_{\xi \geq c}) \geq c \mathbb{P}(\xi \geq c) = c \mathbb{P}(\xi \geq c)$. \square

Definition 2.8 (Variance). The variance (or dispersion) of a random variable ξ is defined to be

$$V_\xi := \mathbb{E}[(\xi - \mathbb{E}\xi)^2]$$

and we define $\sigma := \sqrt{V_\xi}$ the standard deviation of V_ξ .

By expanding the definition, we find $V_\xi = \mathbb{E}\xi^2 - (\mathbb{E}\xi)^2$.

Definition 2.9 (Covariance). The covariance of random variables ξ and η is defined to be

$$\text{cov}(\xi, \eta) := \mathbb{E}[(\xi - \mathbb{E}\xi)(\eta - \mathbb{E}\eta)].$$

Proposition 2.3. For random variables ξ, η , we have

- $V_{\xi+\eta} = V_\xi + V_\eta + 2\text{cov}(\xi, \eta)$;
- if $\text{cov}(\xi, \eta) = 0$, then $V_{\xi+\eta} = V_\xi + V_\eta$.

Proof. Clear. \square

Proposition 2.4 (Chebyshev's Inequality). Let ξ be a integrable random variable. Then, for all $\epsilon > 0$,

$$\mathbb{P}(|\xi - \mathbb{E}\xi| \geq \epsilon) \leq \frac{V_\xi}{\epsilon^2}.$$

Proof. By Markov inequality, for $\xi \geq 0$, we have

$$\mathbb{P}(\xi \geq \epsilon) = \mathbb{P}(\xi^2 \geq \epsilon^2) \leq \frac{\mathbb{E}\xi^2}{\epsilon^2}.$$

Thus, by replacing ξ by $|\xi - \mathbb{E}\xi|$, we have the required inequality. \square

Proposition 2.5 (Exponential Chebyshev's Inequality). Let $\xi \geq 0$ be a random variable and let $\epsilon, t > 0$ such that $\xi, e^{t\xi}$ are integrable. Then,

$$\mathbb{P}(\xi \geq \epsilon) \leq e^{-t\epsilon} \mathbb{E}(e^{t\xi}).$$

Proof. We observe, by Markov inequality

$$\mathbb{P}(\xi \geq \epsilon) = \mathbb{P}(e^{t\xi} \geq e^{t\epsilon}) \leq \frac{E(e^{t\xi})}{e^{t\epsilon}}.$$

□

Proposition 2.6 (Tail Probability). Let $\xi \geq 0$ be an integrable random variable. Then,

$$\mathbb{E}(\xi) = \int_{(0,\infty)} \mathbb{P}(\xi \geq x) \lambda(dx).$$

Proof. By change of variable, we have,

$$\mathbb{E}\xi = \int_{\Omega} \xi d\mathbb{P} = \int_{(0,\infty)} x(\xi_*\mathbb{P})(dx) = \int_{(0,\infty)} \int_{[0,x]} \lambda(dt)(\xi_*\mathbb{P})(dx).$$

Then, by Fubini's theorem to the function $g : (t, x) \mapsto \mathbf{1}_{[0,x]}(t)$, we have

$$\int_{(0,\infty)} \int_{[0,x]} \lambda(dt)(\xi_*\mathbb{P})(dx) = \int_{(0,\infty)^2} g(t, x) \lambda(dt)(\xi_*\mathbb{P})(dx) = \int_{(0,\infty)} \mathbb{P}(\xi \geq x) \lambda(dx)$$

as required. □

Definition 2.10 (Normal Random Variable). A random variable ξ is said to be norm if $\xi_*\mathbb{P} = f\lambda$ where

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

for some $m \in \mathbb{R}, \sigma > 0$. We denote this by $\xi \sim \mathcal{N}(m, \sigma^2)$.

Proposition 2.7. Let $\xi \sim \mathcal{N}(m, \sigma^2)$. Then, $\mathbb{E}\xi = m$ and $V_\xi = \sigma^2$, and so, a normal random variable is fully determined by its mean and variance.

Proof. Exercise. □

Definition 2.11 (Moment). Given a random variable ξ , we define the k -th moment of ξ to be $\mathbb{E}(\xi^k)$.

2.2 Transformation of Random Variables

Let $F_\xi(x)$ be a distribution function of a random variable ξ . Then, if ϕ is a real valued continuous function, we would like to consider the distribution of η where $\eta = \phi(\xi)$. An easy observation is that

$$F_\eta(y) = \mathbb{P}(\eta \leq y) = \mathbb{P}(\eta \in \phi^{-1}(-\infty, y]) = \int_{\phi^{-1}(-\infty, y]} d(\xi_*\mathbb{P}).$$

As one might expect, elementary methods from first year probability are sufficient for most cases we encounter (consider the case ϕ is linear or quadratic).

Suppose now that ξ is absolutely continuous (and so, has a density by Radon-Nikodym), we would like to find the density of $\eta := \phi(\xi)$. Assume first that $\xi(\Omega) \in I$ where I is a finite or infinite open interval and let ϕ be continuously differentiable and strictly increasing on I . Denoting $h(y) = \phi^{-1}(\{y\})$ which is well-defined and differentiable, for all $y \in \phi(I)$,

$$F_\eta(y) = \mathbb{P}(\eta \leq y) = \mathbb{P}(\xi \leq \phi^{-1}(\{y\})) = \int_{(-\infty, h(y)]} f_\xi d\lambda = \int_{(-\infty, y]} f_\xi(h(z))h'(z)\lambda(dz)$$

where f_ξ is the density of ξ . Hence, the density of η is $f_\eta(h(y))h'(y) = f_\xi(h(y))|h'(y)|$. Similarly, if ϕ is strictly decreasing, η remain to have the density $f_\xi(h(y))|h'(y)|$. With this in mind, we may obtain the density for a large class of transformations by de compositing the density into a strictly increasing and decreasing parts.

In the case that (ξ, η) is a random vector with joint distribution F and let ϕ be a continuous function. Then, $\phi(\xi, \eta)$ has distribution

$$F_{\phi(\xi, \eta)}(z) = \int_{\phi^{-1}(-\infty, z]} d((\xi, \eta)_* \mathbb{P}).$$

2.3 Independence

We recall that two random variables ξ, η are said to be independent if $(\xi, \eta)_* \mathbb{P} = \xi_* \mathbb{P} \otimes \eta_* \mathbb{P}$. Thus, if F_ξ, F_η are distributions of ξ and η respectively, then, $F_{(\xi, \eta)}(x, y) = F_\xi(x)F_\eta(y)$ for all $x, y \in \mathbb{R}$ where $F_{(\xi, \eta)}$ is a distribution of the random vector (ξ, η) .

Proposition 2.8. If ξ, η are independent random variables, then the distribution of $\xi + \eta$ is

$$F_{\xi+\eta}(z) = \int_{\mathbb{R}} F_\eta(z-x) \xi_* \mathbb{P}(dx) = \int_{\mathbb{R}} F_\xi(z-y) \eta_* \mathbb{P}(dy).$$

Proof. By taking $\phi : \mathbb{R}^2 \rightarrow \mathbb{R} : (x, y) \mapsto x + y$, we have

$$F_{\xi+\eta}(z) = \int_{\phi^{-1}(-\infty, z]} d((\xi, \eta)_* \mathbb{P}) = \int_{\phi^{-1}(-\infty, z]} d(\xi_* \mathbb{P} \otimes \eta_* \mathbb{P})$$

by independence. By considering that $(x, y) \in \phi^{-1}(-\infty, z]$ if and only if $x + y \leq z$, we have by Fubini's theorem,

$$\begin{aligned} \int_{\phi^{-1}(-\infty, z]} d(\xi_* \mathbb{P} \otimes \eta_* \mathbb{P}) &= \int_{\mathbb{R}^2} \mathbf{1}_{x+y \leq z} d(\xi_* \mathbb{P} \otimes \eta_* \mathbb{P}) \\ &= \int_{\mathbb{R}} \left(\int_{\mathbb{R}} \mathbf{1}_{x+y \leq z} \eta_* \mathbb{P}(dy) \right) \xi_* \mathbb{P}(dx) \\ &= \int_{\mathbb{R}} F_\eta(z-x) \xi_* \mathbb{P}(dx) = \int_{\mathbb{R}} F_\xi(z-y) \eta_* \mathbb{P}(dy). \end{aligned}$$

□

Recalling the definition of the convolution of a function, we may reformulate the above as the following corollaries.

Definition 2.12 (Convolution). Given two real-valued functions $f, g : \Omega \rightarrow \mathbb{R}$, we define the convolution of f with g by

$$f * g := t \mapsto \int_{\mathbb{R}} f(x)g(t-x)\mu(dx).$$

Corollary 6.2. The distribution function of the sum of two independent random variables is the convolution of their distribution functions.

Corollary 6.3. If ξ, η are independent absolutely continuous random variables, then, the density of $\xi + \eta$ is the convolution of their densities.

Proposition 2.9. Let ξ, η be independent integrable random variables. Then, $\xi \cdot \eta$ is integrable and $\mathbb{E}(\xi \cdot \eta) = \mathbb{E}(\xi)\mathbb{E}(\eta)$.

Proof. It is clearly true for indicator functions and so, we may extend to simple function by the linearity of expectation. Hence, by monotone convergence, the statement is true for non-negative random variables and hence true for arbitrary random variables by taking $\xi = \xi^+ - \xi^-$ and $\eta = \eta^+ - \eta^-$. \square

Definition 2.13. Random variables ξ, η are said to be uncorrelated if $\text{cov}(\xi, \eta) = 0$.

Proposition 2.10. Independent random variables are uncorrelated.

Proof. Clear since $\text{cov}(\xi, \eta) = \mathbb{E}(\xi \cdot \eta) - \mathbb{E}\xi\mathbb{E}\eta$. \square

We note that the converse is not true. Namely, uncorrelated does not imply independence.

2.4 Weak Law of Large Numbers

Consider $(\Omega_n, \mathcal{A}, \mathbb{P}_n)$ as a (finite) probability space such that

$$\Omega_n := \{\omega \mid \omega = (a_1, \dots, a_n), a_i \in \{0, 1\}\}, \mathcal{A} = \mathcal{P}(\Omega_n)$$

and

$$\mathbb{P}_n(\{\omega\}) = p(\omega) = p^{\sum_{i=1}^n a_i} q^{n - \sum_{i=1}^n a_i}$$

for some $0 < p < 1$ and $q = 1 - p$. Let $\xi_1, \dots, \xi_n : \Omega_n \rightarrow \{0, 1\}$ be random variables such that $\xi_i(\omega) = a_i$. It is easy to check that ξ_i are independent and identically distributed (iid.). Indeed,

$$(\xi_i)_* \mathbb{P}_n(\{1\}) = \mathbb{P}_n(\{\omega \mid a_j(\omega) = 1\}) = p \sum_{k=0}^{n-1} \binom{n-1}{k} p^k q^{n-k} = p$$

and $(\xi_i)_* \mathbb{P}_n(\{0\}) = (\xi_i)_* \mathbb{P}_n(\{1\}^c) = 1 - p = q$.

Now defining $S_n := \sum_{i=1}^n \xi_i$, we observe that

$$\mathbb{E}S_n = \sum_{i=1}^n \mathbb{E}\xi_i = \sum_{i=1}^n p = np.$$

Thus, the expectation of $\frac{1}{n}S_n$ is simply p . We now ask what is $|\frac{1}{n}S_n(\omega) - p|$. Immediately, we observe that $|\frac{1}{n}S_n(\omega) - p|$ cannot tends to 0 point-wise since $S_n(0) = 0$ for all n . Nonetheless, we observe that $\mathbb{P}_n(S_n = 0) = q^n \rightarrow 0$ as $n \rightarrow \infty$.

Recalling that the Kolmogorov extension theorem, as $\{\mathbb{P}_n\}$ is consistent, there exists a unique probability measure \mathbb{P} on $(\mathbb{R}^\infty, \mathcal{B}(\mathbb{R}^\infty))$ such that $\mathbb{P}(\xi \in \mathbb{R}^\infty \mid (\xi_1, \dots, \xi_n) \in B_n) = \mathbb{P}_n((\xi_1, \dots, \xi_n) \in B_n)$. With this measure in mind, using Chebyshev's inequality, and the fact that they are independent and hence, uncorrelated, we obtain

$$\mathbb{P}\left(\left|\frac{1}{n}S_n - p\right| \geq \epsilon\right) \leq \frac{V(\frac{1}{n}S_n)}{\epsilon^2} = \frac{1}{\epsilon^2} \sum_{j=1}^n V\left(\frac{1}{n}\xi_j\right) = \frac{1}{n^2\epsilon^2} \sum_{i=1}^n V_{\xi_i} = \frac{npq}{n^2\epsilon^2} \rightarrow 0$$

as $n \rightarrow \infty$. Thus, $\frac{1}{n}S_n$ converges to p in measure (probability).

This fact is known as Bernoulli's law of large numbers. By noting that we only used the uncorrelated fact, we obtain a more general theorem.

In general (with arbitrary ξ_i), we call the quantity $\frac{1}{n}S_n$ the time average and the (weak) law of large numbers tells us the time average converges in measure to the space average $\mathbb{E}\xi_i$.

Theorem 7 (Weak Law of Large Numbers). let ξ_1, ξ_2, \dots be integrable random variables. Defining $S_n^{(c)} = \sum_{i=1}^n (\xi_i - \mathbb{E}\xi_i)$ (we note that $\mathbb{E}S_n^{(c)} = 0$). Then, if ξ_1, ξ_2, \dots are uncorrelated and for all i , $V_{\xi_i} \leq C$ for some $C > 0$, for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\left|\frac{1}{n}S_n^{(c)}\right| \geq \epsilon\right) = 0$$

i.e. $\frac{1}{n}S_n^{(c)} \rightarrow 0$ in measure.

Proof. By Chebyshev's inequality, as ξ_i are uncorrelated,

$$\mathbb{P}\left(\left|\frac{1}{n}S_n^{(c)}\right| \geq \epsilon\right) \leq \frac{V(\frac{S_n^{(c)}}{n})}{\epsilon^2} \leq \frac{c}{n\epsilon^2} \rightarrow 0$$

as $n \rightarrow \infty$. □

Corollary 7.1. If ξ_1, ξ_2, \dots integrable iid. random variables such that $V_{\xi_i} < \infty$, then, $\frac{1}{n} \sum_{i=1}^n \xi_i$ converges to $\mathbb{E}\xi_i$ in measure.

We would also like to consider the limiting behaviour of the distribution. Let us first recall the big and small-O notation.

Given sequences of functions $(f_n), (g_n)$, we denote $g_n = O_{n \rightarrow \infty}(|f_n|)$ if $|g_n/f_n|$ is eventually bounded, i.e. there exists some c, N such that for all $n \geq N$, $|g_n| \leq c|f_n|$. On the other hand, we denote $g_n = o_{n \rightarrow \infty}(|f_n|)$ if $\lim_{n \rightarrow \infty} |g_n/f_n| = 0$.

Returning to our example of the Bernoulli random variables, we have the following result.

Theorem 8 (Local Limit Theorem). For any $0 < p < 1$,

$$\max_{0 \leq k \leq n} \left| \mathbb{P}(S_n = k) - \frac{1}{\sqrt{2\pi np(1-p)}} e^{-\frac{x^2}{2p(1-p)}} \right| = o\left(\frac{1}{\sqrt{n}}\right),$$

where $x = x_{k,n} := \frac{k-np}{\sqrt{n}}$.

Proof. Let $A_n > 0$ such that $A_n = o(n)$ (e.g. $A_n = n^\epsilon$ for some $0 < \epsilon < 1$). Then, let $k \in \mathbb{N}$ such that $|x_{k,n}| \leq A_n/\sqrt{n}$ and so,

$$np - A_n \leq k \leq np + A_n.$$

Then, $k, n - k \rightarrow \infty$ as $n \rightarrow \infty$ and using Stirling's formula,

$$\begin{aligned} \mathbb{P}(S_n = k) &= \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} \\ &= \frac{\sqrt{2\pi n}}{\sqrt{2\pi k} \sqrt{2\pi(n-k)}} e^{n \log n - k \log k - (n-k) \log(n-k)} p^k (1-p)^{n-k} \left(1 + O\left(\frac{1}{n}\right)\right). \end{aligned}$$

Then, by noting that $np + x\sqrt{n} = k = np(1 + O(A_n/n))$ and $n(1-p) - x\sqrt{n} = n - k = n(1-p)(1 + O(A_n/n))$, we have

$$\begin{aligned} \mathbb{P}(S_n = k) &= \sqrt{\frac{n}{np2\pi n(1-p)}} \left(1 + O\left(\frac{A_n}{n}\right)\right) e^{n \log n - (np+x\sqrt{n})(\log(np) + \frac{x}{p\sqrt{n}} - \frac{x^2}{2p^2\sqrt{n}} + O(\frac{x}{\sqrt{n}})^3)} \\ &\quad e^{-(n(1-p)-x\sqrt{n})(\log(n(1-p) - \frac{x}{(1-p)\sqrt{n}} - \frac{x^2}{2(1-p)^2n} + O(\frac{x}{\sqrt{n}})^3))} e^{k \log p + (n-k) \log(1-p)} \\ &= \sqrt{\frac{1}{2\pi np(1-p)}} \left(1 + O\left(\frac{A_n}{n}\right)\right) e^{-(np+x\sqrt{n})(\frac{x}{p\sqrt{n}} - \frac{x^2}{2p^2n} + O(\frac{x}{\sqrt{n}})^3)} \\ &\quad e^{(n(1-p)-x\sqrt{n})(\frac{x}{(1-p)\sqrt{n}} - \frac{x^2}{2(1-p)^2n} + O(\frac{x}{\sqrt{n}})^3)} \\ &= \sqrt{\frac{1}{2\pi np(1-p)}} \left(1 + O\left(\frac{A_n}{n}\right)\right) \\ &\quad e^{-\sqrt{n}x + \frac{x^2}{2p} - \frac{x^2}{p} + \frac{x^3}{2p^2\sqrt{n}}} e^{O(\frac{x}{\sqrt{n}})^3 + \sqrt{n}x + \frac{x^2}{2(1-p)} - \frac{x^2}{1-p} - \frac{x^3}{2(1-p)^2\sqrt{n}}} \\ &= \frac{1 + O\left(\frac{A_n}{n}\right)}{\sqrt{2\pi np(1-p)}} e^{-\frac{x^2}{2}(\frac{1}{p} + \frac{1}{1-p}) + O(\frac{A_n^3}{\sqrt{n}^3\sqrt{n}})} \\ &= \frac{1}{\sqrt{2\pi p(1-p)n}} e^{-\frac{x^2}{2p(1-p)}} \left(1 + O\left(\frac{A_n^3}{n^2}\right) + O\left(\frac{A_n}{n}\right)\right). \end{aligned}$$

Hence, choosing $A_n = n^{7/12}$, the result follows for $np - A_n \leq k \leq np + A_n$. Finally, by monotonicity of $P(S_n = k)$, we have the result follows for all $k \leq n$. \square

Theorem 9 (de Moivre-Laplace Central Limit Theorem). For any $0 < p < 1$ and $x \in \mathbb{R}$,

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{S_n - np}{\sqrt{np(1-p)}} \leq x\right) = \Phi(x),$$

where $\Phi(x)$ is the distribution of the normal random variable with parameters $m = 0, \sigma^2 = 1$.

