# Review on Probability

김지수 (Jisu KIM)

통계적 기계학습(Statistical Machine Learning), 2025 1st semester

## Probability Spaces

A probability space is a triple $(\Omega, \mathcal{F}, P)$ where $\Omega$ is a set of "outcomes," $\mathcal{F}$ is a set of "events," and $P : \mathcal{F} \to [0, 1]$ is a function that assigns probabilities to events.

**Definition.** Let $\Omega$ be a set. A nonempty collection $\mathcal{F}$ of subsets of $\Omega$ is called $\sigma$-algebra (or field) if

(i) if $A \in \mathcal{F}$ then $\Omega \backslash A \in \mathcal{F}$, and

(ii) if $A_1, A_2, \cdots \in \mathcal{F}$, then $\bigcup\limits_{i=1}^{\infty} A_i \in \mathcal{F}$.

**Example.** $\mathcal{F} = \{\phi, \Omega\}$ trivial $\sigma-$field

$\mathcal{F} = 2^{\Omega} = \{A| \ A \subset \Omega\}$ : power set $\Longrightarrow \sigma-$field

Without $P$, $(\Omega, \mathcal{F})$ is called a measurable space, i.e., it is a space on which we can put a measure.

**Definition.** A measure is a nonnegative countably additive set function; that is, for an $\sigma$-algebra $\mathcal{F}$, a function $\mu : \mathcal{F} \to [0, \infty]$ is a measure if

(i) $\mu(A) \geq \mu(\phi) = 0$ for all $A \in \mathcal{F}$, and

(iii) For $A_1, A_2, \cdots \in \mathcal{F}$ with $A_i \cap A_j = \phi$ for any $i \neq j$,

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

**Definition.** (1) $\mu(\Omega) < \infty \Longrightarrow$ finite measure

(2) $\mu(\Omega) = 1 \Longrightarrow$ probability measure

(3) $\exists$ a partition $A_1, A_2, \cdots$ with $\bigcup\limits_{i=1}^{\infty} A_i = \Omega$ and $\mu(A_i) < \infty \Longrightarrow \sigma-$finite measure

**Theorem** ([1, Theorem 1.1.4]). *Let $\mu$ be a measure on $(\Omega, \mathcal{F})$.*

*(i) Monotonicity. If $A \subset B$ then $\mu(A) \leq \mu(B)$.*

*(ii) Subadditivity. If $A \subset \bigcup\limits_{i=1}^{\infty} A_i$ then $\mu(A) \le \sum\limits_{i=1}^{\infty} \mu(A_i)$.*

*(iii) Continuity from below. $A_n \uparrow A$ ( i.e. $A_1 \subset A_2 \subset \cdots$ and $A = \bigcup\limits_{i=1}^{\infty} A_i$) then $\mu(A_i) \uparrow \mu(A)$.*

*(iv) Continuity from above. $A_n \downarrow A$ ( i.e. $A_1 \supset A_2 \supset \cdots$ and $A = \bigcap\limits_{i=1}^{\infty} A_i$) with $\mu(A_1) < \infty$ then $\mu(A_i) \downarrow \mu(A)$.*

**Definition.** Let $\mathcal{A}$ be a class of subsets of $\Omega$. Then $\sigma(\mathcal{A})$ denotes the smallest $\sigma-$algebra that contains $\mathcal{A}$.

For any any $\mathcal{A}$, such $\sigma(\mathcal{A})$ exists and is unique: [1, Exercise 1.1.1].

**Definition.** Borel $\sigma-$field on $\mathbb{R}^d$, denoted by $\mathcal{R}^d$, is the smallest $\sigma-$field containing all open sets.

**Theorem** ([1, Theorem 1.1.2]). *There is a unique measure $\mu$ on $(\mathbb{R}, \mathcal{R})$ with*

$$\mu((a, b]) = b - a.$$

*Such measure is called Lebesgue measure.*

**Example** ([1, Example 1.1.3]). Product space

$(\Omega_i, \mathcal{F}_i, \mathcal{P}_i)$ : sequence of probability spaces

Let $\Omega = \Omega_1 \times \cdots \times \Omega_n = \{(\omega_1, \cdots, \omega_n) | \omega_i \in \Omega_i\}$

$\mathcal{F} = \mathcal{F}_1 \times \cdots \times \mathcal{F}_n =$ the $\sigma-$field generated by $A_1 \times \cdots \times A_n$, where $A_i \in \mathcal{F}_i$

$P = P_1 \times \cdots \times P_n$ (i.e. $P(A_1 \times \cdots \times A_n) = P_1(A_1) \cdots P_n(A_n)$)

# Distribution and Random Variables

**Definition.** Let $(\Omega, \mathcal{F})$ and $(S, \mathcal{S})$ are measurable spaces. A mapping $X : \Omega \to S$ is a measurable map from $(\Omega, \mathcal{F})$ to $(S, \mathcal{S})$ if

$$\text{for all } B \in \mathcal{S}, \ X^{-1}(B) := \{\omega \in \Omega : \ X(\omega) \in B\} \in \mathcal{F}.$$

If $(S, \mathcal{S}) = (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ and $d > 1$ then $X$ is called a random vector. If $d = 1$, $X$ is called a random variable.

**Example.** A trivial but useful example of a random variable is indicator function $1_A$ of a set $A \in \mathcal{F}$:

$$1_A(\omega) = \begin{cases} 1 & \omega \in A, \\ 0 & \omega \notin A. \end{cases}$$

If $X$ is a random variable, then $X$ induces a probability measure on $\mathbb{R}$.

**Definition.** The probability measure $\mu$ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ defined as $\mu(A) = P(X \in A)$ for all $A \in \mathcal{B}(\mathbb{R})$ is called the distribution of $X$.

*Remark.* The distribution can be defined similarly for random vectors.

The distribution of a random variable $X$ is usually described by giving its distribution function.

**Definition.** The distribution function $F(x)$ of a random variable $X$ is defined as $F(x) = P(X \le x)$.

**Theorem** ([1, Theorem 1.2.1]). *Any distribution function $F$ has the following properties:*

(i) $F$ *is nondecreasing.*

(ii) $\lim\limits_{n \to \infty} F(x) = 1, \quad \lim\limits_{n \to -\infty} F(x) = 0.$

(iii) $F$ *is right continuous. i.e.* $\lim\limits_{y \downarrow x} F(y) = F(x)$.

(iv) $P(X < x) = F(x-) = \lim\limits_{y \uparrow x} F(x)$.

(v) $P(X = x) = F(x) - F(x-)$.

**Theorem** ([1, Theorem 1.2.2]). *If $F$ satisfies (i) (ii) (iii) in [1, Theorem 1.2.1], then it is the distribution function of some random variable. That is, there exists a triple $(\Omega, \mathcal{F}, P)$ and a random variable $X$ such that $F(x) = P(X \le x)$.*

**Theorem.** *If $F$ satisfies (i) (ii) (iii), then $\exists!$ probability measure $\mu$ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ such that for all $a < b$,*

$\mu((a, b]) = F(b) - F(a)$

**Definition.** If $X$ and $Y$ induce the same distribution $\mu$ on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$, we say $X$ and $Y$ are equal in distribution. We write

$$X \stackrel{d}{=} Y.$$

**Definition.** When the distribution function $F(x) = P(X \leq x)$ has the form $F(x) = \int_{-\infty}^{x} f(y) dy$, then we say $X$ has the density function $f$.

*Remark.* $f$ is not unique, but unique up to Lebesque measure 0.

**Theorem** ([1, Theorem 1.3.2]). *If $X : (\Omega, \mathcal{F}) \to (S, \mathcal{S})$ and $f : (S, \mathcal{S}) \to (T, \mathcal{T})$ are measurable maps, then $f(X)$ is measurable.*

**Theorem.** *$f : (S, \mathcal{S}) \to (T, \mathcal{T})$ and suppose $\mathcal{S} = \sigma(\text{open sets})$, $\mathcal{T} = \sigma(\text{open sets})$. Then, if $f$ is continuous then $f$ is measurable.*

**Theorem** ([1, Theorem 1.3.3]). *If $X_1, \cdots, X_n$ are random variables and $f : (\mathbb{R}^n, \mathcal{R}^n) \to (\mathbb{R}, \mathcal{R})$ is measurable, then $f(X_1, \cdots, X_n)$ is a random variable.*

**Theorem** ([1, Theorem 1.3.4]). *If $X_1, \cdots, X_n$ are random variables then $X_1 + \cdots + X_n$ is a random variable.*

*Remark.* If $X, Y$ are random variables, then

$$cX \ (c \text{ is scalar}), \ X \pm Y, \ XY, \ \sin(X), \ X^2, \ \cdots,$$

are all random variables.

**Theorem** ([1, Theorem 1.3.5]). *$\inf_n X_n, \ \sup_n X_n, \ \limsup_n X_n, \ \liminf_n X_n$ are random variables.*

# Integration

Let $\mu$ be a $\sigma$-finite measure on $(\Omega, \mathcal{F})$.

**Definition.** For any predicate $Q(\omega)$ defined on $\Omega$, we say $Q$ is true $(\mu-)$almost everywhere (or a.e.) if $\mu(\{\omega : Q(\omega) \ is \ false\}) = 0$

**Step 1.**

**Definition.** $\varphi$ is a simple function if $\varphi(\omega) = \sum\limits_{i=1}^{n} a_i 1_{A_i}$ with $A_i \in \mathcal{F}$

If $\varphi$ is a simple function and $\varphi \geq 0$, we let

$$\int \varphi d\mu = \sum\limits_{i=1}^{n} a_i \mu(A_i)$$

**Step 2.**

**Definition.** If $f$ is measurable and $f \geq 0$ then we let

$$\int f d\mu = \sup\{\int h d\mu : \ 0 \leq h \leq f \ and \ h \ simple\}$$

**Step 3.**

**Definition.** We say measurable $f$ is integrable if $\int |f| d\mu < \infty$

let $f^+(x) := f(x) \vee 0$, $f^-(x) := (-f)(x) \vee 0$ where $a \vee b = \max(a, b)$

We define the integral of $f$ by

$$\int f d\mu = \int f^+ d\mu - \int f^- d\mu$$

we can also define $\int f d\mu$ if $\int f^+ d\mu = \infty$ and $\int f^- d\mu < \infty$, or $\int f^+ d\mu < \infty$ and $\int f^- d\mu = \infty$

**Theorem.** *(1.4.7) Suppose $f$ and $g$ are integrable.*

*(i) If $f \geq 0$ a.e. then $\int f d\mu \geq 0$*

*(ii) $\forall a \in \mathbb{R}, \ \int a f d\mu = a \int f d\mu$*

*(iii) $\int f + g d\mu = \int f d\mu + \int g d\mu$*

*(iv) If $g \leq f$ a.e. then $\int g d\mu \leq \int f d\mu$*

*(v) If $g = f$ a.e. then $\int g d\mu = \int f d\mu$*

*(vi) $|\int f d\mu| \leq \int |f| d\mu$*

# Independence

**Definition.** Let $(\Omega, \mathcal{F}, P)$ be probability space. Two events $A, B \in \mathcal{F}$ are independent if

$P(A \cap B) = P(A) \times P(B)$

Two random variables $X$ and $Y$ are independent if

$\forall C, D \in \mathcal{R}, \ P(X \in C, \ Y \in D) = P(X \in C)P(Y \in D)$

Two $\sigma$-fields $\mathcal{F}_1$ and $\mathcal{F}_2 (\subset \mathcal{F})$ are independent if

$\forall A \in \mathcal{F}_1, \ \forall B \in \mathcal{F}_2, \ A$ and $B$ are independent.

*Remark.* An infinite collection of objects ($\sigma-$fields, random variables, or sets) is said to be independent if every finite subcollection is.

**Definition.** $\sigma-$fields $\mathcal{F}_1, \cdots, \mathcal{F}_n$ are independent if

$P(\bigcap\limits_{i=1}^{n} A_i) = \prod\limits_{i=1}^{n} P(A_i), \ \forall A_i \in \mathcal{F}_i$

random variables $X_1, \cdots, X_n$ are independent if

$P(\bigcap\limits_{i=1}^{n} \{X_i \in B_i\}) = \prod\limits_{i=1}^{n} P(X_i \in B_i), \ \forall B_i \in \mathcal{R}$

Sets $A_1, \cdots, A_n$ are independent if

$P(\bigcap\limits_{i \in I} A_i) = \prod\limits_{i \in I} P(A_i)$ for all $I \subset \{1, \cdots, n\}$

*Remark.* the definition of independent events is not enough to assume pairwise independent, which is $P(A_i \cap A_j) = P(A_i)P(A_j), \ i \neq j$. It is clear that indenendent events are pairwise independent, but converse is not true.

**Example.** Let $X_1, \ X_2, \ X_3$ be independent random variables with $P(X_i = 0) = P(X_i = 1) = \frac{1}{2}$

Let $A_1 = \{X_2 = X_3\}, \ A_2 = \{X_3 = X_1\}$ and $A_3 = \{X_1 = X_2\}$. These events are pairwise independent but not independent.

# Weak laws of large numbers

## Various modes of convergence

$\{X_n\}$ and $X$ are random variables defined on $(\Omega, \mathcal{F}, P)$

**Definition.** $X_n \to X$ almost surely (a.s.) ( with probability 1(w.p. 1), almost everywhere(a.e.) ) if $P\{\omega : X_n(\omega) \to X(\omega)\} = 1$

Equivalent definition : $\forall \epsilon,\ \lim_{m\to\infty} P\{\omega : |X_n(\omega) - X(\omega)| \leq \epsilon\ \forall n \geq m\} = 1$

or $\forall \epsilon,\ \lim_{m\to\infty} P\{\omega : |X_n(\omega) - X(\omega)| > \epsilon\ \forall n \geq m\} = 0$

**Definition.** $X_n \to X$ in probability (in pr, $\xrightarrow{p}$) if $\lim_{n\to\infty} P\{|X_n - X| > \epsilon\} = 0$

**Theorem.** $X_n \to X$ a.s. $\Longrightarrow X_n \xrightarrow{p} X$

*Remark.* $X_n \xrightarrow{p} X \nRightarrow X_n \to X$ a.s.

**Definition.** $X_n \to X$ in $L_p$, $0 < p < \infty$

if $\lim_{n\to\infty} E(|X_n - X|^p) = 0$ provided $E|X_n|^p < \infty$, $E|X|^p < \infty$.

**Theorem.** $X_n \to X$ *in* $L_p \implies X_n \xrightarrow{p} X$

**Theorem.** *(Chebyshev inequality)*

$P(|X| \geq \epsilon) \leq \frac{E|X|^p}{\epsilon^p}$

*Remark.* $X_n \xrightarrow{p} X \nRightarrow X_n \to X$ in $L_p$

**Example.** $\Omega = [0,1]$, $\mathcal{F} = \mathcal{B}[0,1]$, $P = Unif[0,1]$

$X(\omega) = 0$, $X_n(\omega) = nI(0 \leq \omega \leq \frac{1}{n})$

Then $P\{|X_n(\omega) - X(\omega)| > \epsilon\} = P\{0 \leq \omega \leq \frac{1}{n}\} = \frac{1}{n} \to 0$

But $E|X_n - X| = E|X_n| = 1$

**Theorem.** $X_n \xrightarrow{p} X$ *and there exists a random variables* $Z$ *s.t.*

$|X_n| \leq Z$ *and* $E|Z|^p < \infty$

*Then $X_n \to X$ in $L_p$.*

*Remark.* If $E|X| < \infty$, then

$$\lim_{n\to\infty} \int_{A_n} |X| dP \to 0 \text{ whenever } P(A_n) \to 0$$

## 2..2.1. $L_2$ weak law

**Theorem** ([1, Theorem 2.2.3]). *Let $X_1, X_2, \cdots$ be uncorrelated random variables with $EX_i = \mu$ and $Var(X_i) \leq C < \infty$*

*Let $S_n = \sum_{i=1}^{n} X_i$. Then*
$\frac{S_n}{n} \to \mu$ *in $L_2$ and so in pr.*

**Theorem** ([1, Theorem 2.2.9]). *Weak law of large numbers*

*Let $X_1, X_2, \cdots$ be i.i.d. random variables with $E|X_i| < \infty$.*

*Let $S_n = X_1 + \cdots + X_n$ and let $\mu = EX_1$.*

*Then $\frac{S_n}{n} \to \mu$ in pr.*

# Weak Convergence

**Definition.** A sequence of distribution function $F_n$ converges weakly to a limit $F$ ($F_n \Rightarrow F$, $F_n \xrightarrow{w} F$)

if $F_n(y) \to F(y)$ $\forall y$ that are continuity points of $F$.

**Definition.** A sequence of random variables $\{X_n\}$ converges weakly or converges in distribution to a limit $X$

($X_n \Rightarrow X$, $X_n \xrightarrow{w} X$, $X_n \xrightarrow{d} X$)

If the distribution function $F_n$ of $X_n$ converges weakly to the distribution of $X$.

**Example** ([1, Example 3.2.1]). Let $X_1, X_2, \cdots$ be iid with $P(X_1 = 1) = P(X_1 = -1) = \frac{1}{2}$.

Let $S_n = X_1 + \cdots + X_n$.

Then $F_n(y) = P(S_n/\sqrt{n} \leq y) \to \int_{-\infty}^{y} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$ $\forall y$

That is, $F_n \Rightarrow N(0, 1)$

**Example** ([1, Example 3.2.3]). Let $X \sim F$ and $X_n = X + \frac{1}{n}$

    Then $F_n(x) = P(X_n \leq x) = F(x - \frac{1}{n}) \to F(x-)$

    Hence $F_n(x) \to F(x)$ only when $F(x) = F(x-)$

    (i.e. $x$ is a continuity point of $F$)

    so $X_n \to X$

**Example** ([1, Example 3.2.4]). $X_p \sim Geo(p)$ (i.e. $P(X_p \geq m) = (1-p)^{m-1}$)

    Then $P(X_p > \frac{x}{p}) = (1-p)^{\frac{x}{p}} \to e^{-x}$ as $p \to 0$

# Central Limit Theorem

**Theorem** ([1, Theorem 3.4.1]). *Let $X_1, X_2, \cdots$ be iid with $EX_i = \mu$ and $Var(X_i) = \sigma^2 > 0$.*

    *If $S_n = X_1 + \cdots + X_n$, then*

    $(S_n - n\mu)/(\sqrt{n}\sigma) \xrightarrow{d} N(0,1)$

**Theorem** ([1, Theorem 3.4.9]). *Berry-Essen theorem*

    *Let $X_1, X_2, \cdots$ be i.i.d. with $EX_i = 0$, $EX_i^2 = \sigma^2$ and $E|X_1|^3 = \rho < \infty$*

    *Let $F_n(x)$ be the distribution function of $(X_1 + \cdots + X_n)/(\sigma\sqrt{n})$ and $\Phi(x)$ be the standard normal distribution.*

    *Then $\sup\limits_{x}|F_n(x) - \Phi(x)| \leq 3\rho/(\sigma^3\sqrt{n})$*

# Stochastic Order Notation

The classical order notation should be familiar to you already.

1. We say that a sequence $a_n = o(1)$ if $a_n \to 0$ as $n \to \infty$. Similarly, $a_n = o(b_n)$ if $a_n/b_n = o(1)$.

2. We say that a sequence $a_n = O(1)$ if the sequence is eventually bounded, i.e. for all $n$ large, $|a_n| \leq C$ for some constant $C \geq 0$. Similarly, $a_n = O(b_n)$ if $a_n/b_n = O(1)$.

3. If $a_n = O(b_n)$ and $b_n = O(a_n)$ then we use either $a_n = \Theta(b_n)$ or $a_n \asymp b_n$.

When we are dealing with random variables we use stochastic order notation.

1. We say that $X_n = o_P(1)$ if for every $\epsilon > 0$, as $n \to \infty$

$$\mathbb{P}\left(|X_n| \geq \epsilon\right) \to 0,$$

    i.e. $X_n$ converges to zero in probability.

2. We say that $X_n = O_P(1)$ if for every $\epsilon > 0$ there is a finite $C(\epsilon) > 0$ such that, for all $n$ large enough:

$$\mathbb{P}\left(|X_n| \geq C(\epsilon)\right) \leq \epsilon.$$

The typical use case: suppose we have $X_1, \ldots, X_n$ which are i.i.d. and have finite variance, and we define:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i.$$

1. $\hat{\mu} - \mu = o_P(1)$ (Weak Law of Large Number)

2. $\hat{\mu} - \mu = O_P(1/\sqrt{n})$ (Central Limit Theorem)

As with the classical order notation, we can do some simple "calculus" with stochastic order notation and observe that for instance: $o_P(1) + O_P(1) = O_P(1)$, $o_P(1)O_P(1) = o_P(1)$ and so on.

# References

[1] Rick Durrett. *Probability: theory and examples*, volume 31 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, fourth edition, 2010.