

# 確率・統計入門

森 立平

2025-12-10

# 目次

前書き	2
第 I 部 確率論	3
第 1 章 はじめに	4
1.1 なぜ確率論と統計学を学ぶか？	4
1.2 本書の構成	4
1.3 その他の参考文献	4
第 2 章 集合論	5
2.1 集合	5
2.2 集合の関係	5
2.3 集合の演算	6
2.4 補集合	6
2.5 集合族	7
第 3 章 確率空間	9
3.1 確率論を数学的に定式化するには	9
3.2 有限集合上の確率空間	10
3.3 可算無限集合上の確率	10
3.4 すべての部分集合を可測にはできない	10
3.5 確率論の公理	11
3.6 確率の性質	12
第 4 章 確率変数と確率分布	15
4.1 確率変数	15
4.2 確率分布	16
4.3 離散確率分布の例	19
4.3.1 様々な離散確率分布	19
4.3.2 二項分布	19
4.3.3 幾何分布	19
4.3.4 超幾何分布	20

4.3.5	ポアソン分布 . . . . .	20
4.3.6	確率空間は何か? . . . . .	20
4.4	連続確率分布の例 . . . . .	21
4.4.1	一様分布 . . . . .	21
4.4.2	正規分布 (ガウス分布) . . . . .	21
4.4.3	指数分布 . . . . .	22
4.5	確率密度関数 . . . . .	22
<b>第 5 章</b>	<b>複数の確率変数</b>	<b>25</b>
5.1	事象の条件付き確率と独立性 . . . . .	25
5.2	確率変数の条件付き確率と独立性 . . . . .	25
5.3	離散型確率変数と確率質量関数 . . . . .	26
5.4	連続型確率変数と確率密度関数 . . . . .	27
5.5	三つ以上の独立確率変数 . . . . .	28
5.6	独立な離散型確率変数の和 . . . . .	29
5.7	独立な連続型確率変数の和 . . . . .	30
<b>第 6 章</b>	<b>期待値、分散、モーメント</b>	<b>32</b>
6.1	期待値 . . . . .	32
6.2	分散 . . . . .	34
6.3	共分散 . . . . .	36
6.4	モーメントとモーメント母関数 . . . . .	37
<b>第 II 部</b>	<b>統計学</b>	<b>41</b>
<b>第 7 章</b>	<b>ベイズ推定</b>	<b>42</b>
7.1	最大事後確率推定 . . . . .	42
7.2	最尤推定 . . . . .	43
7.3	全変動距離 . . . . .	43
7.4	損失関数 . . . . .	45
<b>第 8 章</b>	<b>仮説検定</b>	<b>47</b>
8.1	事前分布を仮定しない推定問題 . . . . .	47
8.2	伝統的な仮説検定 . . . . .	47
8.3	単純仮説 . . . . .	48
8.4	最強力検定 . . . . .	48
8.5	尤度比検定 . . . . .	49
<b>第 III 部</b>	<b>漸近論</b>	<b>51</b>
<b>第 9 章</b>	<b>大数の法則と集中不等式</b>	<b>52</b>

9.1	大数の弱法則 . . . . .	52
9.2	チェルノフ上界 . . . . .	52
9.3	キュムラント母関数の性質 . . . . .	53
9.4	キュムラント母関数の例 . . . . .	55
9.5	ルジャンドル変換 . . . . .	55
9.6	クラメールの定理 . . . . .	56
<b>第 10 章</b>	<b>正規分布と中心極限定理</b>	<b>59</b>
10.1	中心極限定理 . . . . .	59
10.2	特性関数 . . . . .	59
10.3	特性関数の応用 . . . . .	60
10.4	連続性定理 . . . . .	60
10.5	中心極限定理の証明 . . . . .	60
<b>第 11 章</b>	<b>サノフの定理、KL ダイバージェンス</b>	<b>61</b>
	<b>参考文献</b>	<b>62</b>

# 前書き

これは確率論と統計学の入門書である。確率論を数学的に取り扱うには通常は測度論とルベーグ積分を用いる。本書では測度論を学ぶ前の数学専攻の学生を対象に確率論と統計学の基礎を解説する。測度論とルベーグ積分を省略するため、しばしば積分と極限の交換などの等式を証明なしに用いる。後で測度論を学んだ後にぜひ振り返って欲しい。

## 第 I 部

### 確率論

## 第 1 章

# はじめに

### 1.1 なぜ確率論と統計学を学ぶか？

This is a book created from markdown and executable code.

See Knuth (1984) for additional discussion of literate programming.

### 1.2 本書の構成

### 1.3 その他の参考文献

## 第 2 章

# 集合論

### 2.1 集合

集合は対象の集まりである。例えば

$$A = \{1, 2, 3\}, \quad B = \{\text{赤}, \text{青}, \text{黄}, \text{緑}\}$$

のように表す。集合を構成するものを**要素**もしくは**元**という。

- $x \in A$  で要素  $x$  は集合  $A$  に含まれる
- $x \notin A$  で要素  $x$  は集合  $A$  に含まれない

を表す。

$|A|$  で集合  $A$  の要素数を表すことにする。要素数が無限の集合を考えることもできる。例えば

$$A = \{n \in \mathbb{N} \mid n \text{ は偶数}\}, \quad B = \{x \in \mathbb{R} \mid x \text{ は無理数}\}$$

という集合は無限集合の例となる。何も要素を持たない集合 (要素数が零の集合) を**空集合**といい、 $\emptyset$  で表す。

### 2.2 集合の関係

集合  $A$  の要素がすべて集合  $B$  に含まれるとき、 $A$  を  $B$  の**部分集合** という。部分集合に関連する集合関係について記号を以下のように定義する。

$$\begin{aligned} A \subseteq B &\stackrel{\text{def}}{\iff} \forall x (x \in A \implies x \in B) && (A \text{ は } B \text{ の部分集合}) \\ A \supseteq B &\stackrel{\text{def}}{\iff} B \subseteq A && (A \text{ は } B \text{ の上位集合}) \\ A = B &\stackrel{\text{def}}{\iff} (A \subseteq B \wedge A \supseteq B) && (A \text{ と } B \text{ は等しい}) \\ A \neq B &\stackrel{\text{def}}{\iff} \neg(A = B) && (A \text{ と } B \text{ は等しくない}) \end{aligned}$$



また、

$$\begin{aligned} A \subsetneq B &\stackrel{\text{def}}{\iff} (A \subseteq B \wedge A \neq B) && (A \text{ は } B \text{ の真部分集合}) \\ A \supsetneq B &\stackrel{\text{def}}{\iff} B \subsetneq A && (A \text{ は } B \text{ の真上位集合}) \end{aligned}$$

とする。部分集合や上位集合のような集合間の関係を**包含関係**という。集合  $A, B$  について、 $A \subseteq B$  と  $B \subseteq A$  のどちらかが成り立つとき、「 $A$  と  $B$  の間に包含関係が成り立つ」という。

包含関係は以下の三条件を満たす。

- (反射律)  $A \subseteq A$ .
- (反対称律)  $(A \subseteq B \wedge B \subseteq A) \iff A = B$ .
- (推移律)  $(A \subseteq B \wedge B \subseteq C) \implies A \subseteq C$ .

## 2.3 集合の演算

複数の集合から新しい集合を作る演算がある。

$$\begin{aligned} A \cup B &:= \{x \mid x \in A \text{ または } x \in B\} && (\text{和集合}) \\ A \cap B &:= \{x \mid x \in A \text{ かつ } x \in B\} && (\text{積集合}) \\ A \setminus B &:= \{x \mid x \in A \text{ かつ } x \notin B\} && (\text{差集合}) \end{aligned}$$

これらの演算は以下の法則を満たす。

- (交換法則)  $A \cup B = B \cup A, \quad A \cap B = B \cap A$ .
- (結合法則)  $(A \cup B) \cup C = A \cup (B \cup C), \quad (A \cap B) \cap C = A \cap (B \cap C)$ .
- (分配法則)  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C), \quad A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ .
- (冪等法則)  $A \cup A = A, \quad A \cap A = A$ .
- (吸収法則)  $A \cup (A \cap B) = A, \quad A \cap (A \cup B) = A$ .
- $A \cup \emptyset = A, \quad A \cap \emptyset = \emptyset$ .

和集合と積集合は結合法則を満たすことから、括弧を使わずに  $A \cup B \cup C$  と表すことができる。集合  $A$  と  $B$  が  $A \cap B = \emptyset$  を満たすとき、**確率論文脈では、「 $A$  と  $B$  は排反である」**(集合論文脈では「互いに素である」) という。

## 2.4 補集合

全体の集合  $\Omega$  というのが文脈上存在する場合は、

$$A^c := \Omega \setminus A \quad (\text{補集合})$$

と定義する。

補集合に関連して以下の法則が成り立つ。

- $\Omega^c = \emptyset, \quad \emptyset^c = \Omega.$
- $A \cup \Omega = \Omega, \quad A \cap \Omega = A.$
- $A \cup A^c = \Omega, \quad A \cap A^c = \emptyset.$
- (二重補集合の法則)  $(A^c)^c = A.$
- (ド・モルガンの法則)  $(A \cup B)^c = A^c \cap B^c, \quad (A \cap B)^c = A^c \cup B^c$

## 2.5 集合族

集合  $\Omega$  について  $2^\Omega$  を  $\Omega$  のすべての部分集合からなる集合を表す。

$$2^\Omega := \{A \subseteq \Omega\}.$$

これを  $\Omega$  の**冪集合**という。例えば  $\Omega = \{1, 2, 3\}$  のとき、

$$2^\Omega = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{2, 3\}, \{1, 3\}, \{1, 2, 3\}\}$$

である。 $\Omega$  が有限集合のとき、 $|2^\Omega| = 2^{|\Omega|}$  が成り立つ。また、 $2^\Omega$  の部分集合を  $\Omega$  上の**(部分) 集合族**と呼ぶ。集合  $\Lambda$  の各  $\lambda \in \Lambda$  に対して集合  $A_\lambda \subseteq \Omega$  が存在するとき、集合族

$$\{A_\lambda \subseteq \Omega \mid \lambda \in \Lambda\} \subseteq 2^\Omega$$

を**添字集合  $\Lambda$  で添字付けられた  $\Omega$  上の集合族**という。添字集合が有限集合の場合は集合族全体の和集合や積集合は二つの集合の和集合と積集合の定義を繰り返し用いることで定義できる。それらは以下のように表す。

$$\bigcup_{\lambda \in \Lambda} A_\lambda, \quad \bigcap_{\lambda \in \Lambda} A_\lambda.$$

添字集合  $\Lambda$  が無限集合の場合は和集合と積集合を以下で定義する。

$$\begin{aligned} \bigcup_{\lambda \in \Lambda} A_\lambda &:= \{x \in \Omega \mid \exists \lambda \in \Lambda, \quad x \in A_\lambda\} \\ \bigcap_{\lambda \in \Lambda} A_\lambda &:= \{x \in \Omega \mid \forall \lambda \in \Lambda, \quad x \in A_\lambda\} \end{aligned}$$

この定義は  $\Lambda$  が有限集合の場合も正しいものである。この場合もド・モルガンの法則は成り立つ。つまり、

$$\begin{aligned} \left( \bigcup_{\lambda \in \Lambda} A_\lambda \right)^c &= \bigcap_{\lambda \in \Lambda} A_\lambda^c \\ \left( \bigcap_{\lambda \in \Lambda} A_\lambda \right)^c &= \bigcup_{\lambda \in \Lambda} A_\lambda^c \end{aligned}$$

が成り立つ。

証明:  $x \in \Omega$  について、

$$\begin{aligned}x \in \left( \bigcup_{\lambda \in \Lambda} A_\lambda \right)^c &\iff x \notin \left( \bigcup_{\lambda \in \Lambda} A_\lambda \right) \\&\iff \neg (\exists \lambda \in \Lambda, \quad x \in A_\lambda) \\&\iff \forall \lambda \in \Lambda, \quad x \notin A_\lambda \\&\iff \forall \lambda \in \Lambda, \quad x \in A_\lambda^c \\&\iff x \in \bigcap_{\lambda \in \Lambda} A_\lambda^c.\end{aligned}$$

## 第 3 章

# 確率空間

### 3.1 確率論を数学的に定式化するには

確率は身近に現れる (感じられる) ものであるが、それを数学的に定式化することは自明な問題ではない。実際に確率論には複数の数学的定式化が存在する。その中で圧倒的に一般的なのが測度論的確率論と呼ばれる定式化である。測度というのは集合の「面積」のようなものであり、確率を測度と捉えるのが測度論的確率論である。これは多くの人間の直感にも自然なものであろう。

まず、確率を考える集合について考えよう。例えばコインを投げて表もしくは裏が出る確率を考えたいときは

$$\Omega = \{H, T\}$$

という集合になる。また、明日の天気の高確率を考えたいときは

$$\Omega = \{ \text{” 晴”}, \text{” 雨”}, \text{” 雪”} \}$$

という集合になるだろう。この  $\Omega$  の部分集合に確率を与える関数  $P: 2^\Omega \rightarrow \mathbb{R}$  を定義しよう。偏りのないコインの場合は以下ようになる。

$$P(\emptyset) = 0, \quad P(\{H\}) = \frac{1}{2}, \quad P(\{T\}) = \frac{1}{2}, \quad P(\{H, T\}) = 1.$$

また、天気の場合は例えば以下ようになる。

$$\begin{aligned} P(\emptyset) &= 0, & P(\{\text{晴}\}) &= 0.7, & P(\{\text{雨}\}) &= 0.2, & P(\{\text{雪}\}) &= 0.1 \\ P(\{\text{晴}, \text{雨}\}) &= 0.9, & P(\{\text{雨}, \text{雪}\}) &= 0.3, & P(\{\text{雪}, \text{晴}\}) &= 0.8, & P(\{\text{晴}, \text{雨}, \text{雪}\}) &= 1. \end{aligned}$$

このように集合  $\Omega$  の部分集合に確率を与えることを考える。 $\Omega$  の要素一つずつに確率を与えれば十分であるようにも思えるが、 $\Omega$  が連続的な場合には  $\Omega$  の一つの要素の確率は 0 になってしまうことが多い。例えば明日の正午の気温が  $10\pi^\circ\text{C}$  になる確率は 0 であろう。そのため、 $\Omega$  の要素ではなく部分集合に確率を与えることにする。そのために測度論が適しているわけである。

### 3.2 有限集合上の確率空間

確率を考える集合を  $\Omega$  とする。この  $\Omega$  のことを**標本空間**という。また、 $\Omega$  の部分集合のことを**事象**という。そして、事象に確率を与える関数  $P: 2^\Omega \rightarrow \mathbb{R}_{\geq 0}$  を**確率測度**という。確率測度は以下の条件を満たす。

1.  $P(\Omega) = 1$ .
2.  $\forall A, B \subseteq \Omega, \quad A \cap B = \emptyset \implies P(A \cup B) = P(A) + P(B)$ .

一つ目の条件は全体の確率が1であるという意味の条件である。二つ目の条件は排反な二つの事象の和集合の確率はそれぞれの事象の確率の和であるという意味の条件である。この二つ目の条件を**有限加法性**という。例えば

$$P(\{\text{” 晴”}, \text{” 雨”}\}) = P(\{\text{” 晴”}\}) + P(\{\text{” 雨”}\})$$

という等式は「晴れもしくは雨になる確率 = 晴れになる確率 + 雨になる確率」という意味の等式になる。よって有限加法性が自然な条件であることが分かるだろう。また、これらのことから、 $P$  は各要素  $\omega \in \Omega$  に対する確率  $P(\{\omega\})$  から一意に定まることが分かる。この標本空間と確率測度のペア  $(\Omega, P)$  を**確率空間**という。

### 3.3 可算無限集合上の確率

標本空間  $\Omega$  が可算無限集合のときも、同様に確率測度を定義することもできるが、ここではより強い以下の条件を考える。

1.  $P(\Omega) = 1$ .
2.  $\forall (A_n \subseteq \Omega)_{n \geq 0}, \quad \forall i \neq j, A_i \cap A_j = \emptyset \implies P\left(\bigcup_{n \geq 0} A_n\right) = \sum_{n \geq 0} P(A_n)$ .

ここで二つ目の条件における  $n \geq 0$  という添字において  $n$  は非負整数全体をわたる。今後もこの記法を用いる。この二つ目の条件を**完全加法性**もしくは **$\sigma$ -加法性**という。右辺の無限和において  $P(A_n) \geq 0$  なので、無限和  $\sum_{n \geq 0} P(A_n)$  が存在するときこれは絶対収束する。よって、この無限和は事象列の並び順に依存しないことが分かる。完全加法性ではなく有限加法性だけを使って確率論を構築する試みもあるが、確率測度の連続性などの性質が失なわれるため標準的な確率論では完全加法性を課す。

### 3.4 すべての部分集合を可測にはできない

より一般に  $\Omega$  が非可算無限集合の場合を考えよう。例えば  $\Omega = [0, 1)$  の場合が考えられる。このとき、すべての  $A \subseteq \Omega$  に確率を与えることができるのだろうか？ そうすると、ごく自然な性質を持つような確率測度が**存在しない**ことが、選択公理を認めると示される。

**定理 3.1** (非可測集合の存在).  $\Omega = [0, 1)$  とする。また、集合の平行移動を

$$A + x := \{a + x - \lfloor a + x \rfloor \mid a \in A\}$$

と定義する。このとき、

$$\forall x \in \Omega, A \subseteq \Omega, \quad P(A+x) = P(A) \quad (\text{平行移動不変性})$$

を満たす確率測度  $P: 2^{[0,1)} \rightarrow \mathbb{R}_{\geq 0}$  は存在しない。

証明.  $\Omega$  上の同値関係を  $x \sim y \stackrel{\text{def}}{\iff} x-y \in \mathbb{Q}$  と定義する。選択公理より、この同値関係の同値類から一つずつ要素を含む集合  $V$  が存在する (Vitali 集合)。任意の  $x \in [0,1)$  について、ある  $z \in V$  が唯一存在して  $x \sim z$  である。よって、任意の  $x \in [0,1)$  について、ある  $z \in V$  と  $q \in \mathbb{Q} \cap [0,1)$  が唯一存在して  $x = z + q - \lfloor z+q \rfloor$  であることから

$$[0,1) = \bigcup_{q \in \mathbb{Q} \cap [0,1)} (V+q)$$

であり、右辺は互いに排反である。よって条件を満たす  $P$  が存在すると仮定すると、

$$\begin{aligned} 1 = P([0,1)) &= P\left(\bigcup_{q \in \mathbb{Q} \cap [0,1)} (V+q)\right) \\ &= \sum_{q \in \mathbb{Q} \cap [0,1)} P(V+q) \quad (\text{完全加法性}) \\ &= \sum_{q \in \mathbb{Q} \cap [0,1)} P(V) \quad (\text{平行移動不変性}) \end{aligned}$$

ここで  $P(V)$  をどのように定めても、それを無限回足して 1 にすることはできない。よって  $P$  は存在しない。  $\square$

### 3.5 確率論の公理

定理 3.1 より、

1. 選択公理
2. 確率測度の完全加法性
3.  $[0,1)$  上の平行移動不変な確率測度
4.  $\Omega$  のすべての部分集合に確率を与える

のどれかを諦めないといけない。標準的な確率論では 4 を諦める。以下に確率空間の厳密な定義を述べる。

**定義 3.1** (完全加法族).  $\Omega$  を集合とする。  $\Omega$  上の集合族  $\mathcal{F} \subseteq 2^\Omega$  が以下を満たすとする。

1.  $\Omega \in \mathcal{F}$ .
2.  $\forall A \in \mathcal{F}, \quad A^c = \Omega \setminus A \in \mathcal{F}$ .
3.  $\forall (A_n \in \mathcal{F})_{n \geq 0}, \quad \bigcup_{n \geq 0} A_n \in \mathcal{F}$ .

このとき、 $\mathcal{F}$  を  $\Omega$  上の完全加法族もしくは  $\sigma$ -加法族という。

**定義 3.2** (確率空間).  $\Omega$  を集合とし、 $\mathcal{F} \subseteq 2^\Omega$  を  $\Omega$  上の完全加法族とする。

また、 $P: \mathcal{F} \rightarrow \mathbb{R}_{\geq 0}$  が以下を満たすとする。

1.  $P(\Omega) = 1$ .

2.  $\forall (A_n \in \mathcal{F})_{n \geq 0}, \quad \forall i \neq j, A_i \cap A_j = \emptyset \implies P\left(\bigcup_{n \geq 0} A_n\right) = \sum_{n \geq 0} P(A_n)$ .

このとき、 $(\Omega, \mathcal{F}, P)$  を確率空間という。また、 $\Omega$  を標本空間、 $\mathcal{F}$  を事象集合、 $P$  を確率測度という。

このように事象集合  $\mathcal{F}$  の元についてのみ確率が与えられる。

### ノート

ペア  $(\Omega, \mathcal{F})$  を可測空間という。確率空間  $(\Omega, \mathcal{F}, P)$  から  $P(\Omega) = 1$  の条件を除いたものが一般の測度空間である。

定理 3.1 の証明では選択公理を用いたが、実際に選択公理を認めないと定理 3.1 が成立しないことが分かっている。以降、 $\Omega$  が非可算無限集合である場合も含めて、 $\mathcal{F} = 2^\Omega$  だと思ふことにする。これは正しくない場合もあるのだが、選択公理を使わない限り矛盾は導かれないので、問題になることはほとんどない。よって、以降は  $(\Omega, P)$  を確率空間とする。

$\Omega = [0, 1]$  で任意の  $0 \leq a < b \leq 1$  について  $P([a, b]) = b - a$  であるような確率空間はとても基本的なものである。今後この性質を満たす確率空間が存在すると仮定して話を進める。

## 3.6 確率の性質

**補題 3.1** (確率のいくつかの性質). 確率空間  $(\Omega, P)$  と任意の  $A, B \subseteq \Omega$  について以下が成り立つ。

1.  $P(A^c) = 1 - P(A)$ .
2.  $B \subseteq A \implies P(B) \leq P(A)$
3.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .
4.  $P(A \cup B) \leq P(A) + P(B)$  (ブールの不等式、union bound).

証明.

- 1.

$$1 = P(\Omega) = P(A \cup A^c) = P(A) + P(A^c).$$

- 2.

$$P(A) = P(B \cup (A \setminus B)) = P(B) + P(A \setminus B) \geq P(B).$$

- 3.

$$\begin{aligned} P(A \cup B) &= P(A \cup (B \setminus A)) = P(A) + P(B \setminus A) \\ P(B) &= P((B \setminus A) \cup (A \cap B)) = P(B \setminus A) + P(A \cap B) \end{aligned}$$

より、 $P(B \setminus A)$  を消去することで得られる。

4. 3 より自明

□

**補題 3.2** (ユニオンバウンド). 確率空間  $(\Omega, P)$  と  $(A_n \subseteq \Omega)_{n \geq 0}$  について、

$$P\left(\bigcup_{n \geq 0} A_n\right) \leq \sum_{n \geq 0} P(A_n).$$

証明.  $B_0 := A_0$ ,  $B_n := A_n \setminus \bigcup_{k=0}^{n-1} A_k$  とおくと、

$$\begin{aligned} P\left(\bigcup_{n \geq 0} A_n\right) &= P\left(\bigcup_{n \geq 0} B_n\right) \\ &= \sum_{n \geq 0} P(B_n) \\ &\leq \sum_{n \geq 0} P(A_n). \end{aligned}$$

□

**定理 3.2** (確率測度の連続性). 確率空間  $(\Omega, P)$  と事象列  $A_0 \subseteq A_1 \subseteq \dots \subseteq \Omega$  について

$$P\left(\bigcup_{n \geq 0} A_n\right) = \lim_{n \rightarrow \infty} P(A_n).$$

また、事象列  $\Omega \supseteq A_0 \supseteq A_1 \supseteq \dots$  について

$$P\left(\bigcap_{n \geq 0} A_n\right) = \lim_{n \rightarrow \infty} P(A_n).$$

証明. 事象列  $A_0 \subseteq A_1 \subseteq \dots \subseteq \Omega$  について考える。  $B_0 := A_0$ ,  $n \geq 1$  について  $B_n := A_n \setminus A_{n-1}$  とおく。このとき、 $i \neq j$  について  $B_i \cap B_j = \emptyset$ 。また  $k \geq 0$  について、 $\bigcup_{n=0}^k B_n = \bigcup_{n=0}^k A_n = A_k$  が成り立つ。

$$\begin{aligned} P\left(\bigcup_{n \geq 0} A_n\right) &= P\left(\bigcup_{n \geq 0} B_n\right) = \sum_{n \geq 0} P(B_n) \\ &= \lim_{k \rightarrow \infty} \sum_{n=0}^k P(B_n) = \lim_{k \rightarrow \infty} P\left(\bigcup_{n=0}^k B_n\right) = \lim_{k \rightarrow \infty} P(A_k). \end{aligned}$$

事象列  $\Omega \supseteq A_0 \supseteq A_1 \supseteq \dots$  について考える。このとき、 $A'_n := A_n^c$  とおくと、 $A'_0 \subseteq A'_1 \subseteq \dots \subseteq \Omega$  を満たす。よって、

$$P\left(\bigcup_{n \geq 0} A'_n\right) = \lim_{n \rightarrow \infty} P(A'_n)$$



である。ド・モルガンの公式より、

$$\begin{aligned} P\left(\left(\bigcap_{n \geq 0} A_n\right)^c\right) &= \lim_{n \rightarrow \infty} P(A_n^c) \\ \Leftrightarrow 1 - P\left(\bigcap_{n \geq 0} A_n\right) &= \lim_{n \rightarrow \infty} (1 - P(A_n)) \\ \Leftrightarrow P\left(\bigcap_{n \geq 0} A_n\right) &= \lim_{n \rightarrow \infty} P(A_n). \end{aligned}$$

□

## 第 4 章

# 確率変数と確率分布

### 4.1 確率変数

確率空間の上の様々な部分集合の確率を調べたい。そのためには確率変数を導入すると便利である。

**定義 4.1** (確率変数). 確率空間  $(\Omega, P)$  について、関数  $X: \Omega \rightarrow \mathbb{R}$  を**確率変数**という。

また、任意の  $A \subseteq \mathbb{R}$  について

$$\Pr(X \in A) := P(\{\omega \in \Omega \mid X(\omega) \in A\})$$

と定義する。さらに、任意の  $a \in \mathbb{R}$  について

$$\Pr(X \geq a) := P(X \in \{x \in \mathbb{R} \mid x \geq a\}) = P(\{\omega \in \Omega \mid X(\omega) \geq a\})$$

と定義する。同様に  $\Pr(X > a)$ ,  $\Pr(X \leq a)$ ,  $\Pr(X < a)$ ,  $\Pr(X = a)$  など定義される。

#### **i** ノート

本当の確率論では確率変数は「確率空間から位相空間への写像で開集合の逆像が可測集合になるような関数」として定義される。また、文献によっては確率変数は「確率空間から可測空間への写像で可測集合の逆像が可測集合になるような関数」として定義されることもある。開集合族が生成する完全加法族をボレル集合族と呼ぶが、この方法で位相空間から可測空間を作ることができる。そうすると前者の定義は「確率空間から位相空間への写像でボレル集合の逆像が可測集合になるような関数」と等価なので、後者の定義に含まれる。後者の定義の中の可測空間が位相空間由来のものに限定したのが前者の定義であると言える。

任意の関数  $f: \mathbb{R} \rightarrow \mathbb{R}$  について  $f \circ X$  は確率変数である。このとき、 $\Pr(f \circ X \in A)$  と書く代わりに  $\Pr(f(X) \in A)$  と書く。

### i ノート

本当の確率論の言葉で述べると任意のボレル関数 (開集合 (ボレル集合としても等価) の逆像がボレル集合となる関数)  $f$  について  $f \circ X$  は確率変数である。

共通の確率空間  $(\Omega, P)$  上の確率変数  $X, Y$  についてその和  $X + Y$  や積  $XY$  も確率変数である。

**例 4.1** (コイン投げ). コインを 2 回独立に投げる場合の確率空間  $(\Omega, P)$  を以下で定義する。

- $\Omega = \{HH, HT, TH, TT\}$ .
- $P(A) = \frac{|A|}{4}$ .

確率変数  $X_1$  と  $X_2$  を

$$\begin{array}{llll} X_1(HH) = 1 & X_1(HT) = 1 & X_1(TH) = 0 & X_1(TT) = 0 \\ X_2(HH) = 1 & X_2(HT) = 0 & X_2(TH) = 1 & X_2(TT) = 0 \end{array}$$

とすると、 $k \in \{1, 2\}$  について、 $X_k$  は  $k$  回目のコイン投げが表だったときに 1、裏だったときに 0 になる。また、 $X = X_1 + X_2$  は表が出た回数を表す確率変数となる。表が一回以上出る確率は以下のように表せる。

$$\Pr(X \geq 1) = P(\{\omega \in \Omega \mid X(\omega) \geq 1\}) = P(\{HH, HT, TH\}) = \frac{3}{4}.$$

## 4.2 確率分布

**定義 4.2** (累積分布関数). 確率変数  $X$  について**累積分布関数**  $F_X: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  は以下で定義される。

$$F_X(x) := \Pr(X \leq x) \quad \forall x \in \mathbb{R}.$$

累積分布関数は定義より単調非減少関数であることが分かる。累積分布関数は連続とは限らないが右連続であることは以下のように確認できる。

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr\left(X \leq x + \frac{1}{n}\right) &= \lim_{n \rightarrow \infty} P\left(\left\{\omega \in \Omega \mid X(\omega) \leq x + \frac{1}{n}\right\}\right) \\ &= P\left(\bigcap_{n=0}^{\infty} \left\{\omega \in \Omega \mid X(\omega) \leq x + \frac{1}{n}\right\}\right) \\ &= P(\{\omega \in \Omega \mid X(\omega) \leq x\}) \\ &= \Pr(X \leq x). \end{aligned}$$

累積分布関数  $F_X: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  は以下の性質を持つ。

1.  $F_X$  は単調非減少。
2.  $F_X$  は右連続。

$$3. \lim_{x \rightarrow \infty} F_X(x) = 1.$$

$$4. \lim_{x \rightarrow -\infty} F_X(x) = 0.$$

確率変数  $X$  の像が高々可算集合のとき、 $X$  を**離散型確率変数**という。離散型確率変数でない確率変数を**連続型確率変数**という。

$X$  が離散型確率変数のとき、 $A$  を  $X$  の像の部分集合とすると

$$\begin{aligned} \Pr(X \in A) &= P(\{\omega \in \Omega \mid X(\omega) \in A\}) \\ &= P\left(\bigcup_{x \in A} \{\omega \in \Omega \mid X(\omega) = x\}\right) \\ &= \sum_{x \in A} P(\{\omega \in \Omega \mid X(\omega) = x\}) \\ &= \sum_{x \in A} \Pr(X = x) \end{aligned}$$

であるので、

$$f_X(x) := \Pr(X = x) \quad \forall x \in \mathbb{R}$$

という関数を用いて

$$\Pr(X \in A) = \sum_{x \in A} f_X(x)$$

と表せる。この  $f_X(x)$  を  $X$  の**確率質量関数**という。

$X$  が連続型確率変数のとき、

$$\Pr(X \in A) = \int_A f_X(x) dx \quad \forall A \subseteq \mathbb{R}$$

を満たす  $f_X: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  が存在するとき、この  $f_X(x)$  を  $X$  の**確率密度関数**という。しかし、確率密度関数は常に存在するとは限らない。

**例 4.2.**  $\Omega = [0, 1)$  とし、 $P: 2^\Omega \rightarrow \mathbb{R}_{\geq 0}$  を  $P([a, b)) = b - a$  を満たす確率測度とする。確率空間  $(\Omega, P)$  上の確率変数  $X: [0, 1) \rightarrow \mathbb{R}$  を

$$X(\omega) = \begin{cases} 2\omega & \text{if } \omega < \frac{1}{2} \\ \frac{1}{2} & \text{otherwise.} \end{cases}$$

と定義する。

この時  $X$  の像は  $[0, 1)$  であるので  $X$  は連続型確率変数である。

また、 $\Pr(X = \frac{1}{2}) = P([\frac{1}{2}, 1) \cup \{\frac{1}{4}\}) = \frac{1}{2}$  である。よって  $X$  の確率密度関数は存在しない。

連続確率変数  $X$  が確率密度関数を持つとき、

$$F_X(x) = \int_{-\infty}^x f_X(z) dz$$

という関係が成り立つ。よって累積分布関数が微分可能なときは

$$f_X(x) = \frac{dF_X}{dx}$$

とおいても構わない。確率密度関数はその積分値にのみ意味を持つ。そのため確率密度関数は一意には定まらない。

例 4.2 の確率変数  $X$  の累積分布関数は以下のようなになる。



図 4.1:  $X$  の累積分布関数

## 4.3 離散確率分布の例

### 4.3.1 様々な離散確率分布

- 二項分布  $\text{Binom}(n, p)$ : 表が出る確率が  $p$  のコインを  $n$  回独立に投げたとき、表が出る回数の分布

$$\Pr(X = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

- 幾何分布: 表が出るまでに投げるコインの回数の分布

$$\Pr(X = k) = (1-p)^{k-1} p.$$

- 超幾何分布: 袋の中に  $N$  個のボールがあって、そのうち  $K$  個が当たりとし、 $n$  個引いたときの当たりの個数の分布

$$\Pr(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}.$$

- ポアソン分布:

$$\Pr(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}.$$

### 4.3.2 二項分布

二項分布の確率質量関数は

$$f_X(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

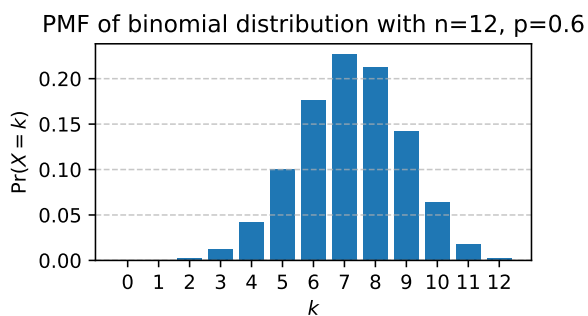


図 4.2: 二項分布の確率質量関数



図 4.3: 二項分布の累積分布関数

### 4.3.3 幾何分布

$$f_X(k) = (1-p)^{k-1} p, \quad \forall k \geq 1$$



図 4.4: 幾何分布の確率質量関数



図 4.5: 幾何分布の累積分布関数

#### 4.3.4 超幾何分布

$$f_X(k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}, \quad \forall k \in \{0, \dots, \min\{K, n\}\}$$

PMF of hypergeometric distribution with K=12, N=40, n = 16



図 4.6: 超幾何分布の確率質量関数

CDF of hypergeometric distribution with K=12, N=40, n = 16



図 4.7: 超幾何分布の累積分布関数

#### 4.3.5 ポアソン分布

$$f_X(k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad \forall k \geq 0$$

#### 4.3.6 確率空間は何か？

確率変数が従う確率分布のみに注目した議論をする場合、確率空間を陽に考えないことがしばしばある。

- 二項分布:  $n$  回のコイン投げ、もしくは無限回のコイン投げ。
- 幾何分布: 無限回のコイン投げ。



図 4.8: ポアソン分布の確率質量関数



図 4.9: ポアソン分布の累積分布関数

- 超幾何分布:  $K$  個の当たりと  $N - K$  個のはずれのボールの  $\binom{N}{K}$  通りの並び順。
- ポアソン分布: ?

## 4.4 連続確率分布の例

- 一様分布
- 正規分布
- 指数分布

### 4.4.1 一様分布

$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise} \end{cases},$$

$$F_X(x) = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } x \in [a, b] \\ 1 & \text{otherwise.} \end{cases}$$

### 4.4.2 正規分布 (ガウス分布)

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

$$F_X(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(z-\mu)^2}{2\sigma^2}} dz$$

特に  $\mu = 0, \sigma = 1$  のとき標準正規分布という。





図 4.10: 一様分布の確率密度関数

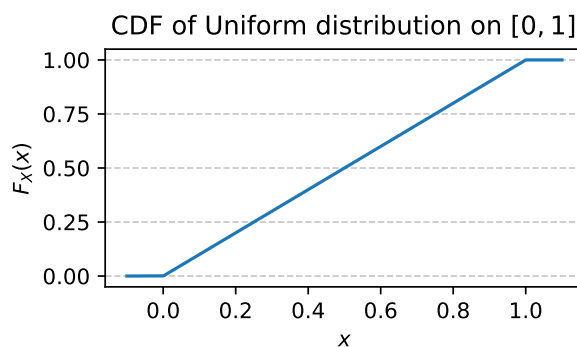


図 4.11: 一様分布の累積分布関数

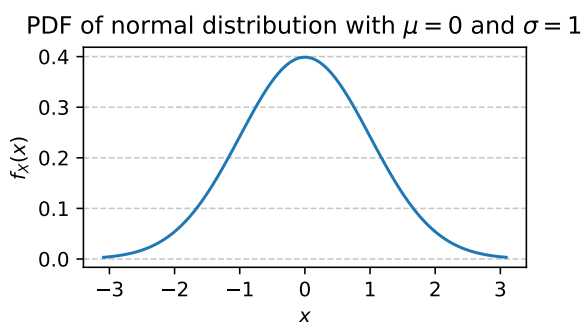


図 4.12: 標準正規分布の確率密度関数



図 4.13: 標準正規分布の累積分布関数

#### 4.4.3 指数分布

$$f_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ \lambda e^{-\lambda x} & \text{otherwise} \end{cases},$$

$$F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 - e^{-\lambda x} & \text{otherwise} \end{cases}$$

### 4.5 確率密度関数

補題 4.1. 確率変数  $X$  について、

$$\begin{aligned} f_{X+a}(x) &= f_X(x-a) & \forall a \in \mathbb{R} \\ f_{aX}(x) &= \frac{1}{|a|} f_X(x/a) & \forall a \in \mathbb{R}_{\neq 0} \end{aligned}$$

はそれぞれ  $X+a$  と  $aX$  の確率密度関数になる。

証明. 関数  $f_Z: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  が

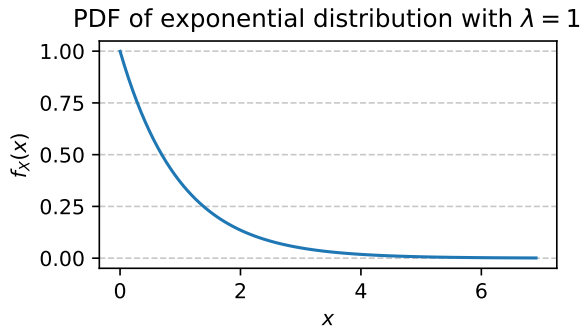


図 4.14: 指数分布の確率密度関数

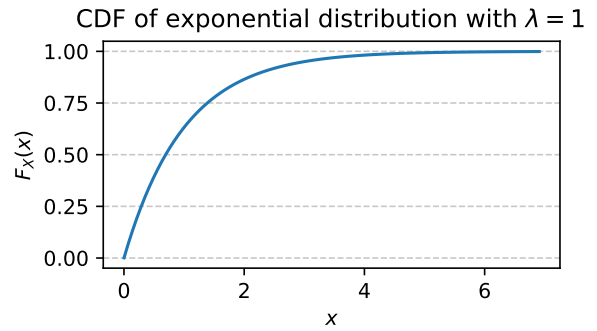


図 4.15: 指数分布の累積分布関数

$$\Pr(Z \leq x) = \int_{-\infty}^x f_Z(z) dz$$

を満たすとき、 $f_Z$  は確率変数  $Z$  の確率密度関数となる。

任意の  $a \in \mathbb{R}$  について、

$$\begin{aligned} \Pr(X + a \leq x) &= \Pr(X \leq x - a) \\ &= \int_{-\infty}^{x-a} f_X(z) dz \\ &= \int_{-\infty}^x f_X(z' - a) dz' \quad (z' = z + a). \end{aligned}$$

任意の  $a > 0$  について、

$$\begin{aligned} \Pr(aX \leq x) &= \Pr(X \leq x/a) \\ &= \int_{-\infty}^{x/a} f_X(z) dz \\ &= \int_{-\infty}^x \frac{1}{a} f_X(z'/a) dz' \quad (z' = az). \end{aligned}$$

任意の  $a < 0$  について、

$$\begin{aligned} \Pr(aX \leq x) &= \Pr(X \geq x/a) \\ &= \int_{x/a}^{\infty} f_X(z) dz \\ &= \int_x^{-\infty} \frac{1}{a} f_X(z'/a) dz' \quad (z' = az) \\ &= - \int_{-\infty}^x \frac{1}{a} f_X(z'/a) dz'. \end{aligned}$$

□

**補題 4.2.**  $J \subseteq \mathbb{R}$  を有界とは限らない区間とし、 $g: J \rightarrow \mathbb{R}$  を  $J$  の内点で微分可能で  $g'(x) > 0$  とする。確率変数  $X$  が  $\Pr(X \in J) = 1$  を満たし確率密度関数を持つとき、

$$f_{g(X)}(x) = \begin{cases} \frac{1}{g'(g^{-1}(x))} f_X(g^{-1}(x)) & \text{if } x \in \text{Image}(g) \\ 0 & \text{otherwise} \end{cases}$$

は  $g(X)$  の確率密度関数になる。

証明.  $\Pr(g(X) \in \text{Image}(g)) = 1$  なので、 $x \notin \text{Image}(g)$  について  $f_{g(X)}(x) = 0$  とおいてよい。

任意の  $x \in \text{Image}(g)$  について

$$\begin{aligned} \Pr(g(X) \leq x) &= \Pr(X \leq g^{-1}(x)) \\ &= \int_{-\infty}^{g^{-1}(x)} f_X(z) dz \\ &= \int_{\inf J}^{g^{-1}(x)} f_X(z) dz \\ &= \int_{\inf \text{Image}(g)}^x f_X(g^{-1}(z')) \frac{1}{g'(g^{-1}(z'))} dz' \quad (z' = g(z)) \\ &= \int_{-\infty}^x f_{g(X)}(z') dz'. \end{aligned}$$

□

## 第 5 章

# 複数の確率変数

### 5.1 事象の条件付き確率と独立性

**定義 5.1** (条件付き確率). 確率空間  $(\Omega, P)$  の事象  $A, B \subseteq \Omega$  について  $P(B) > 0$  のとき、 $B$  における  $A$  の条件付き確率は以下で定義される。

$$P(A | B) := \frac{P(A \cap B)}{P(B)}.$$

二つの事象  $A, B \subseteq \Omega$  を考える文脈では  $P(A \cap B)$  を同時確率、 $P(A)$ ,  $P(B)$  を周辺確率という。

**定義 5.2** (事象の独立性). 確率空間  $(\Omega, P)$  の事象  $A, B \subseteq \Omega$  について

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

を満たすとき、事象  $A$  と  $B$  は独立であるという。

**例 5.1** (二回のコイン投げ). 標本空間を  $\Omega = \{HH, HT, TH, TT\}$  とし、確率測度を  $P(A) = \frac{|A|}{4} \quad \forall A \subseteq \Omega$  とする。このとき、 $A = \{HH, HT\}$ ,  $B = \{HH, TH\}$  とおくと、

$$P(A \cap B) = \frac{1}{4}, \quad P(A) = P(B) = \frac{1}{2}$$

より  $P(A \cap B) = P(A)P(B)$  を満たすことが分かる。よって事象  $A, B$  は独立である。

事象  $A$  と  $B$  が独立であり、 $P(B) > 0$  であるとき、 $P(A | B) = P(A)$  である。

### 5.2 確率変数の条件付き確率と独立性

事象は確率変数を通じて表すことが多い。そのため確率変数を用いた条件付き確率も定義する。

**定義 5.3** (確率変数). 確率空間  $(\Omega, P)$  上の確率変数  $X_1, X_2$  について同時確率を

$$\Pr(X_1 \in A, X_2 \in B) := P(\{\omega \in \Omega \mid X_1(\omega) \in A\} \cap \{\omega \in \Omega \mid X_2(\omega) \in B\}) \quad \forall A, B \subseteq \mathbb{R}$$

と定義する。確率変数が三つ以上の場合も同様に定義する。また  $\Pr(X_2 \in B) > 0$  のとき、条件付き確率は以下で定義する。

$$\begin{aligned}\Pr(X_1 \in A \mid X_2 \in B) &:= P(\{\omega \in \Omega \mid X_1(\omega) \in A\} \mid \{\omega \in \Omega \mid X_2(\omega) \in B\}) \\ &= \frac{\Pr(X_1 \in A, X_2 \in B)}{\Pr(X_2 \in B)} \quad \forall A, B \subseteq \mathbb{R}.\end{aligned}$$

任意の  $A, B \subseteq \mathbb{R}$  について、

$$\Pr(X_1 \in A, X_2 \in B) = \Pr(X_1 \in A) \Pr(X_2 \in B)$$

を満たすとき、**確率変数  $X_1$  と  $X_2$  は独立である**という。

**補題 5.1.** 確率空間  $(\Omega, P)$  上の確率変数  $X_1, X_2$  が独立であるとする。このとき任意の関数  $f_1, f_2: \mathbb{R} \rightarrow \mathbb{R}$  について、 $f_1(X_1), f_2(X_2)$  は独立である。

証明.

$$\begin{aligned}\Pr(f_1(X_1) \in A, f_2(X_2) \in B) &= \Pr(X_1 \in f_1^{-1}(A), X_2 \in f_2^{-1}(B)) \\ &= \Pr(X_1 \in f_1^{-1}(A)) \Pr(X_2 \in f_2^{-1}(B)) \\ &= \Pr(f_1(X_1) \in A) \Pr(f_2(X_2) \in B).\end{aligned}$$

□

二つ以上の確率変数の累積分布関数を

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) := \Pr(X_1 \leq x_1, \dots, X_n \leq x_n)$$

と定義する。

### 5.3 離散型確率変数と確率質量関数

**定義 5.4.** 確率空間  $(\Omega, P)$  上の離散型確率変数  $X_1, X_2$  について、**同時確率質量関数**を

$$f_{X_1, X_2}(x_1, x_2) := \Pr(X_1 = x_1, X_2 = x_2)$$

と定義する。

同時確率質量関数からそれぞれの確率変数の確率質量関数が得られる。

$$\begin{aligned}f_{X_1}(x_1) &= \sum_{x_2} f_{X_1, X_2}(x_1, x_2) \\ f_{X_2}(x_2) &= \sum_{x_1} f_{X_1, X_2}(x_1, x_2)\end{aligned}$$

それぞれの確率変数の確率質量関数を**周辺質量関数**と呼ぶ。同時確率質量関数から周辺質量関数を計算する操作のことを**周辺化**という。

**定義 5.5.** 確率空間  $(\Omega, P)$  上の離散型確率変数  $X_1, X_2$  について、条件付き確率質量関数を

$$f_{X_1|X_2}(x_1 | x_2) := \Pr(X_1 = x_1 | X_2 = x_2)$$

と定義する。

条件付き確率の定義より

$$f_{X_1, X_2}(x_1, x_2) = f_{X_1|X_2}(x_1 | x_2) f_{X_2}(x_2)$$

が成り立つ。

**補題 5.2.** 確率空間  $(\Omega, P)$  上の離散型確率変数  $X_1, X_2$  について、 $X_1$  と  $X_2$  が独立  $\iff$

$$f_{X_1, X_2}(x_1, x_2) = f_{X_1}(x_1) f_{X_2}(x_2) \quad \forall x_1 \in \text{Image}(X_1), x_2 \in \text{Image}(X_2).$$

証明.  $\implies$  は自明。 $\impliedby$  を示す。

$$\begin{aligned} \Pr(X_1 \in A, X_2 \in B) &= \sum_{x_1 \in A} \sum_{x_2 \in B} \Pr(X_1 = x_1, X_2 = x_2) \\ &= \sum_{x_1 \in A} \sum_{x_2 \in B} \Pr(X_1 = x_1) \Pr(X_2 = x_2) \\ &= \Pr(X_1 \in A) \Pr(X_2 \in B) \end{aligned}$$

□

**例 5.2** (二回のコイン投げ). 標本空間を  $\Omega = \{\text{HH}, \text{HT}, \text{TH}, \text{TT}\}$  とし、確率測度を  $P(A) = \frac{|A|}{4} \quad \forall A \subseteq \Omega$  とする。

$$\begin{aligned} X_1(\text{HH}) &= X_1(\text{HT}) = 1, & X_1(\text{TH}) &= X_1(\text{TT}) = 0 \\ X_2(\text{HH}) &= X_2(\text{TH}) = 1, & X_2(\text{HT}) &= X_2(\text{TT}) = 0 \end{aligned}$$

と定義する。このとき、

$$\begin{aligned} f_{X_1, X_2}(x_1, x_2) &= \frac{1}{4} \quad \forall x_1, x_2 \in \{0, 1\} \\ f_{X_1}(x_1) &= f_{X_2}(x_2) = \frac{1}{2} \quad \forall x_1, x_2 \in \{0, 1\} \end{aligned}$$

より  $f_{X_1, X_2}(x_1, x_2) = f_{X_1}(x_1) f_{X_2}(x_2)$  を満たすことが分かる。よって確率変数  $X_1, X_2$  は独立である。

## 5.4 連続型確率変数と確率密度関数

確率空間  $(\Omega, P)$  上の連続型確率変数  $X_1, X_2$  について、同時確率密度関数を

$$\Pr(X_1 \in A, X_2 \in B) = \int_A \left( \int_B f_{X_1, X_2}(x_1, x_2) dx_2 \right) dx_1$$

を満たすものと定義する。 $X_1$  と  $X_2$  が確率密度関数を持つ場合でも  $X_1$  と  $X_2$  の同時確率密度関数が存在するとは限らない。例えば  $X_1 = X_2$  の場合がその例である。逆に  $X_1$  と  $X_2$  が同時確率密度関数を持つとき、それぞれの確率密度関数は

$$\begin{aligned} f_{X_1}(x_1) &= \int_{-\infty}^{\infty} f_{X_1, X_2}(x_1, x_2) dx_2 \\ f_{X_2}(x_2) &= \int_{-\infty}^{\infty} f_{X_1, X_2}(x_1, x_2) dx_1 \end{aligned}$$

で得られる。この操作を確率密度関数の**周辺化**という。同時確率密度関数を持つ確率変数  $X_1, X_2$  が独立であるとき、

$$f_{X_1, X_2}(x_1, x_2) = f_{X_1}(x_1) f_{X_2}(x_2)$$

が成り立つ。

## 5.5 三つ以上の独立確率変数

三つ以上の確率変数についても同時確率質量関数、同時確率密度関数を同様に定義する。独立性についても同様に定義する。

**定義 5.6.** 確率空間  $(\Omega, P)$  上の確率変数  $X_1, X_2, \dots, X_n$  が**独立**  $\stackrel{\text{def}}{\iff}$

$$\Pr(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{k=1}^n \Pr(X_k \in A_k) \quad \forall A_1, \dots, A_n \subseteq \mathbb{R}.$$

**定義 5.7.** 確率空間  $(\Omega, P)$  上の確率変数  $X_1, X_2, \dots, X_n$  が**互いに独立**  $\stackrel{\text{def}}{\iff}$  任意の  $1 \leq i < j \leq n$  について、 $X_i$  と  $X_j$  が独立。

確率変数  $X_1, \dots, X_n$  が独立であるとき、それらは互いに独立であることは以下から分かる。

$$\begin{aligned} \Pr(X_1 \in A, X_2 \in B) &= \Pr(X_1 \in A, X_2 \in B, X_3 \in \mathbb{R}, \dots, X_n \in \mathbb{R}) \\ &= \Pr(X_1 \in A) \Pr(X_2 \in B) \prod_{k=3}^n \Pr(X_k \in \mathbb{R}) \\ &= \Pr(X_1 \in A) \Pr(X_2 \in B). \end{aligned}$$

一方で確率変数  $X_1, \dots, X_n$  が互いに独立であっても、それらが独立であるとは限らない。例えば、離散型確率変数  $X_1, \dots, X_n$  を

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \begin{cases} \frac{1}{2^{n-1}} & \text{if } \sum_{k=1}^n x_k \text{ is even} \\ 0 & \text{otherwise} \end{cases} \quad \forall x_1, \dots, x_n \in \{0, 1\}$$

と定義する。このとき、確率変数  $X_n$  を周辺化すると

$$\begin{aligned} f_{X_1, \dots, X_{n-1}}(x_1, \dots, x_{n-1}) &= f_{X_1, \dots, X_n}(x_1, \dots, x_{n-1}, 0) + f_{X_1, \dots, X_n}(x_1, \dots, x_{n-1}, 1) \\ &= \frac{1}{2^{n-1}} \end{aligned}$$

となる。つまり、 $X_1, \dots, X_{n-1}$  は  $\{0, 1\}^{n-1}$  上の一様分布に従う。この確率分布は  $X_1, \dots, X_n$  について対称なので、どの確率変数を周辺化しても一様分布に従う。一様分布は独立なので、 $n \geq 3$  のとき、どの二つの確率変数も独立である。

一方で、 $n \geq 2$  のとき、これらの確率変数の周辺確率は一様である。つまり、

$$f_{X_k}(0) = f_{X_k}(1) = \frac{1}{2} \quad \forall k = 1, 2, \dots, n.$$

しかし、

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{k=1}^n f_{X_k}(x_k) \quad \forall x_1, \dots, x_n \in \{0, 1\}$$

は成り立たないので独立ではない。

## 5.6 独立な離散型確率変数の和

離散型確率変数  $X_1$  と  $X_2$  が独立であるとする。このとき、 $X_1 + X_2$  の確率質量関数は

$$f_{X_1+X_2}(x) = \sum_z f_{X_1}(z) f_{X_2}(x-z)$$

で与えられる。これを確率質量関数の**畳み込み**という。

**例 5.3.** 二項分布の  $N = 1$  の場合をベルヌーイ分布  $\text{Ber}(p)$  と呼ぶ。つまり、 $X \sim \text{Ber}(p) \stackrel{\text{def}}{\iff}$

$$\Pr(X = 0) = 1 - p, \quad \Pr(X = 1) = p$$

である。確率変数  $X_1, X_2$  が独立で  $\text{Ber}(p)$  に従うとする。このとき、 $X_1 + X_2$  は  $\text{Binom}(2, p)$  に従う。

$$\begin{aligned} \Pr(X_1 + X_2 = 0) &= \Pr(X_1 = 0, X_2 = 0) \\ &= \Pr(X_1 = 0) \Pr(X_2 = 0) = (1 - p)^2 \\ \Pr(X_1 + X_2 = 1) &= \Pr(X_1 = 0, X_2 = 1) + \Pr(X_1 = 1, X_2 = 0) \\ &= \Pr(X_1 = 0) \Pr(X_2 = 1) + \Pr(X_1 = 1) \Pr(X_2 = 0) \\ &= 2p(1 - p) \\ \Pr(X_1 + X_2 = 2) &= \Pr(X_1 = 1, X_2 = 1) \\ &= \Pr(X_1 = 1) \Pr(X_2 = 1) = p^2 \end{aligned}$$

と計算できる。よって

$$\Pr(X_1 + X_2 = k) = \binom{2}{k} p^k (1 - p)^{2-k}$$

が成り立ち  $X_1 + X_2 \sim \text{Binom}(2, p)$  であることが分かる。一般に、独立確率変数  $X_1, X_2$  について  $X_1 \sim$



$\text{Binom}(n, p)$ ,  $X_2 \sim \text{Binom}(m, p)$  のとき、 $X_1 + X_2 \sim \text{Binom}(n + m, p)$  である。

$$\begin{aligned}
 \Pr(X_1 + X_2 = k) &= \sum_{\ell \geq 0} \Pr(X_1 = \ell) \Pr(X_2 = k - \ell) \\
 &= \sum_{\ell \geq 0} \binom{n}{\ell} p^\ell (1-p)^{n-\ell} \binom{m}{k-\ell} p^{k-\ell} (1-p)^{m-k+\ell} \\
 &= \left( \sum_{\ell \geq 0} \binom{n}{\ell} \binom{m}{k-\ell} \right) p^k (1-p)^{n+m-k} \\
 &= \binom{n+m}{k} p^k (1-p)^{n+m-k}.
 \end{aligned}$$

よって、 $X_1, \dots, X_n$  が独立で  $\text{Ber}(p)$  に従うとき、 $X_1 + \dots + X_n$  は  $\text{Binom}(n, p)$  に従う。

**例 5.4.** 確率変数  $X_1, X_2$  が独立で  $X_1 \sim \text{Poisson}(\lambda_1)$ ,  $X_2 \sim \text{Poisson}(\lambda_2)$  とする。このとき、

$$\begin{aligned}
 f_{X_1+X_2}(x) &= \sum_{z=0}^x f_{X_1}(z) f_{X_2}(x-z) \\
 &= \sum_{z=0}^x \frac{\lambda_1^z}{z!} e^{-\lambda_1} \frac{\lambda_2^{x-z}}{(x-z)!} e^{-\lambda_2} \\
 &= \frac{1}{x!} e^{-(\lambda_1+\lambda_2)} \sum_{z=0}^x \binom{x}{z} \lambda_1^z \lambda_2^{x-z} \\
 &= \frac{1}{x!} e^{-(\lambda_1+\lambda_2)} (\lambda_1 + \lambda_2)^x
 \end{aligned}$$

が成り立つ。よって  $X_1 + X_2 \sim \text{Poisson}(\lambda_1 + \lambda_2)$  である。独立なポアソン分布の和はポアソン分布に従う。

これらの例のように、確率分布の族 (集合) が畳み込みに閉じているとき、確率分布の族は**再生性**を持つという。二項分布の場合はパラメータ  $p$  を固定したときに再生性を持つ。

## 5.7 独立な連続型確率変数の和

連続確率変数  $X_1$  と  $X_2$  が独立で密度関数を持つとき、

$$\begin{aligned}
 \Pr(X_1 + X_2 \in A) &= \int \int_{y+z \in A} f_{X_1}(y) f_{X_2}(z) dy dz \\
 &= \int \int_{x \in A} f_{X_1}(y) f_{X_2}(x-y) dy dx \quad (x = y + z) \\
 &= \int_A \left( \int_{-\infty}^{\infty} f_{X_1}(y) f_{X_2}(x-y) dy \right) dx \quad (x = y + z)
 \end{aligned}$$

と表せるので、 $X_1 + X_2$  は確率密度関数を持ち

$$f_{X_1+X_2}(x) = \int_{-\infty}^{\infty} f_{X_1}(z) f_{X_2}(x-z) dz$$

とすることができる。確率変数  $X_1, X_2$  が独立で  $X_1 \sim N(\mu_1, \sigma_1^2), X_2 \sim N(\mu_2, \sigma_2^2)$  とする。このとき、

$$\begin{aligned}
f_{X_1+X_2}(x) &= \int_{-\infty}^{\infty} f_{X_1}(z) f_{X_2}(x-z) dz \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(z-\mu_1)^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(x-z-\mu_2)^2}{2\sigma_2^2}} dz \\
&= \frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2}} \int_{-\infty}^{\infty} e^{-\frac{\sigma_2^2(z-\mu_1)^2 + \sigma_1^2(x-z-\mu_2)^2}{2\sigma_1^2\sigma_2^2}} dz \\
&= \frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2}} \int_{-\infty}^{\infty} e^{-\frac{(\sigma_1^2+\sigma_2^2)\left(z-\frac{\sigma_2^2\mu_1+\sigma_1^2(x-\mu_2)}{\sigma_1^2+\sigma_2^2}\right)^2 + \sigma_2^2\mu_1^2 + \sigma_1^2(x-\mu_2)^2 - \frac{(\sigma_2^2\mu_1+\sigma_1^2(x-\mu_2))^2}{\sigma_1^2+\sigma_2^2}}{2\sigma_1^2\sigma_2^2}} dz \\
&= \frac{1}{\sqrt{2\pi(\sigma_1^2+\sigma_2^2)}} e^{-\frac{\sigma_2^2\mu_1^2 + \sigma_1^2(x-\mu_2)^2 - \frac{(\sigma_2^2\mu_1+\sigma_1^2(x-\mu_2))^2}{\sigma_1^2+\sigma_2^2}}{2\sigma_1^2\sigma_2^2}} \\
&= \frac{1}{\sqrt{2\pi(\sigma_1^2+\sigma_2^2)}} e^{-\frac{(x-(\mu_1+\mu_2))^2}{2(\sigma_1^2+\sigma_2^2)}}
\end{aligned}$$

が成り立つので  $X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$  であることが分かる。よって正規分布は再生性を持つ。

## 第 6 章

# 期待値、分散、モーメント

### 6.1 期待値

**定義 6.1** (期待値). 離散型確率変数  $X$  の期待値は

$$\begin{aligned}\mathbb{E}[X] &:= \sum_{x \in \text{Image}(X)} x \Pr(X = x) \\ &= \sum_{x \in \text{Image}(X)} x f_X(x)\end{aligned}$$

と定義される。ここで、右辺の和が絶対収束しない場合は (適当な順番で和を取って収束したとしても) 期待値は定義されない。

連続型確率変数  $X$  が確率密度関数  $f_X$  を持つとき、その期待値は

$$\mathbb{E}[X] := \int_{-\infty}^{\infty} x f_X(x) dx$$

と定義される。ただし、広義積分で上記の積分が存在する場合でも

$$\int_{-\infty}^{\infty} |x| f_X(x) dx$$

が存在しない場合は期待値は定義されない。

連続型確率変数の期待値に関する様々な証明はルベグ積分の知識を必要とするので本書では扱わない。以下、証明はすべて離散確率変数の場合に限って与える。

**補題 6.1.** 確率変数  $X, Y$  について

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

である。また、 $X$  と  $Y$  が独立のとき、

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]$$

である。

証明.

$$\begin{aligned} \mathbb{E}[X + Y] &= \sum_z z f_{X+Y}(z) \\ &= \sum_z z \sum_x f_{X,Y}(x, z-x) \\ &= \sum_{x,y} (x+y) f_{X,Y}(x, y) \quad (y = z-x) \\ &= \sum_{x,y} x f_{X,Y}(x, y) + \sum_{x,y} y f_{X,Y}(x, y) \\ &= \sum_x x f_X(x) + \sum_y y f_Y(y) \\ &= \mathbb{E}[X] + \mathbb{E}[Y]. \end{aligned}$$

$X$  と  $Y$  が独立のとき、

$$\begin{aligned} \mathbb{E}[XY] &= \sum_z z f_{XY}(z) \\ &= \sum_z z \sum_{x \neq 0} f_{X,Y}(x, z/x) \\ &= \sum_z z \sum_{x \neq 0} f_X(x) f_Y(z/x) \\ &= \sum_{x \neq 0, y} xy f_X(x) f_Y(y) \quad (y = z/x) \\ &= \sum_{x,y} xy f_X(x) f_Y(y) \\ &= \left( \sum_x x f_X(x) \right) \left( \sum_y y f_Y(y) \right) \\ &= \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

□

**補題 6.2** (Law of the unconscious statistician (LOTUS)). 任意の関数  $g: \mathbb{R} \rightarrow \mathbb{R}$  について、

1.  $X$  が離散型確率変数のとき、

$$\mathbb{E}[g(X)] = \sum_x g(x) f_X(x).$$

2.  $X$  が連続型確率変数で確率密度関数を持つとき、

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

証明.  $X$  を離散型確率変数とする。

$$\begin{aligned}
\mathbb{E}[g(X)] &= \sum_x x \Pr(g(X) = x) \\
&= \sum_x x P(\{\omega \in \Omega \mid g(X(\omega)) = x\}) \\
&= \sum_x x P\left(\bigcup_{y \in \text{Image}(X): g(y)=x} \{\omega \in \Omega \mid X(\omega) = y\}\right) \\
&= \sum_x \sum_{y \in \text{Image}(X): g(y)=x} x P(\{\omega \in \Omega \mid X(\omega) = y\}) \\
&= \sum_{y \in \text{Image}(X)} g(y) f_X(y).
\end{aligned}$$

□

**命題 6.1** (期待値の性質). 任意の確率変数  $X$  と  $a \in \mathbb{R}$  について

$$\begin{aligned}
\mathbb{E}[X + a] &= \mathbb{E}[X] + a \\
\mathbb{E}[aX] &= a\mathbb{E}[X].
\end{aligned}$$

**定理 6.1** (マルコフの不等式). 任意の**非負**確率変数  $X$  と  $a > 0$  について

$$\Pr(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

証明.

$$\begin{aligned}
\mathbb{E}[X] &= \sum_{x \in \text{Image}(X)} f_X(x)x \\
&= \sum_{\text{Image}(X): x \geq a} f_X(x)x + \sum_{\text{Image}(X): x < a} f_X(x)x \\
&\geq \sum_{\text{Image}(X): x \geq a} f_X(x)x \quad (\Pr(X \geq 0) = 1) \\
&\geq \sum_{\text{Image}(X): x \geq a} f_X(x)a \\
&= \Pr(X \geq a)a.
\end{aligned}$$

□

## 6.2 分散

**定義 6.2** (分散). 確率変数  $X$  が期待値を持つとき、その**分散**を

$$\text{Var}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2]$$

と定義する。また、分散の平方根を**標準偏差**という。

確率変数  $X$  が分散を持つとき、

$$\begin{aligned}\mathrm{Var}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

である。分散は定義より非負の値を取る。

**命題 6.2** (分散の性質). 任意の確率変数  $X$  と  $a \in \mathbb{R}$  について

$$\begin{aligned}\mathrm{Var}[X + a] &= \mathrm{Var}[X] \\ \mathrm{Var}[aX] &= a^2 \mathrm{Var}[X].\end{aligned}$$

分散は直感的には期待値からのはずれ具合を表す値である。

**定理 6.2** (チェビシェフの不等式). 任意の確率変数  $X$  と  $a > 0$  について

$$\Pr(|X - \mathbb{E}[X]| \geq a) \leq \frac{\mathrm{Var}[X]}{a^2}.$$

証明.

$$\begin{aligned}\Pr(|X - \mathbb{E}[X]| \geq a) &= \Pr((X - \mathbb{E}[X])^2 \geq a^2) \\ &\leq \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{a^2} = \frac{\mathrm{Var}[X]}{a^2}.\end{aligned}$$

□

**補題 6.3** (互いに独立な確率変数の和). 確率変数  $X_1, X_2, \dots, X_n$  が互いに独立のとき

$$\mathrm{Var}[X_1 + \dots + X_n] = \mathrm{Var}[X_1] + \dots + \mathrm{Var}[X_n].$$

証明.

$$\begin{aligned}\mathrm{Var}[X_1 + \dots + X_n] &= \mathbb{E}[(X_1 + \dots + X_n - \mathbb{E}[X_1 + \dots + X_n])^2] \\ &= \mathbb{E}[(X_1 - \mathbb{E}[X_1]) + \dots + (X_n - \mathbb{E}[X_n])^2] \\ &= \mathbb{E}\left[\sum_i (X_i - \mathbb{E}[X_i])^2 + 2 \sum_{i < j} (X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])\right] \\ &= \sum_i \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] + 2 \sum_{i < j} \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] \\ &= \sum_i \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] + 2 \sum_{i < j} \mathbb{E}[X_i - \mathbb{E}[X_i]] \mathbb{E}[X_j - \mathbb{E}[X_j]] \\ &= \sum_i \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] = \sum_i \mathrm{Var}[X_i].\end{aligned}$$

□

**例 6.1.** 互いに独立な確率変数  $X_1, \dots, X_n$  のそれぞれが確率変数  $X$  と同分布であるとし、

$$Y := \frac{1}{n}(X_1 + \dots + X_n)$$

と定義する。このとき、

$$\begin{aligned}\mathbb{E}[Y] &= \mathbb{E}[X] \\ \text{Var}[Y] &= \frac{1}{n} \text{Var}[X]\end{aligned}$$

である。互いに独立な確率変数の平均を取ると期待値は変わらず、分散は小さくなる。

**例 6.2.** 独立確率変数  $X_1, \dots, X_n \sim \text{Ber}(1/2)$  について

$$Y_S := \sum_{i \in S} X_i \pmod{2} \quad \forall S \subseteq \{1, 2, \dots, n\}$$

と定義すると、これらは互いに独立である。また、 $S \neq \emptyset$  について  $Y_S \sim \text{Ber}(1/2)$  である。任意の関数  $g: \{0, 1\} \rightarrow \mathbb{R}$  について、

$$Y := \frac{1}{2^n - 1} \sum_{S \subseteq \{1, \dots, n\}: S \neq \emptyset} g(Y_S)$$

と定義すると、 $X \sim \text{Ber}(1/2)$  について、

$$\begin{aligned}\mathbb{E}[Y] &= \mathbb{E}[g(X)] \\ \text{Var}[Y] &= \frac{1}{2^n - 1} \text{Var}[g(X)].\end{aligned}$$

## 6.3 共分散

**定義 6.3** (共分散). 確率変数  $X, Y$  が期待値を持つとき、その**共分散**を

$$\begin{aligned}\text{Cov}[X, Y] &:= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y]\end{aligned}$$

と定義する。共分散がゼロである確率変数のペアを**無相関**であるという。

定義より、 $\text{Cov}[X, X] = \text{Var}[X]$  であることが分かる。

**命題 6.3.** 独立確率変数  $X, Y$  は無相関である。

逆に無相関であっても独立とは限らない。

例 6.3. 確率変数  $X$  を

$$f_X(0) = f_X(+1) = f_X(-1) = \frac{1}{3}$$

を満たすものとし、 $Y = X^2$  とする。このとき、

$$\begin{aligned}\text{Cov}[X, Y] &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \\ &= \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2] \\ &= \mathbb{E}[X] - \mathbb{E}[X]\mathbb{E}[X^2] \\ &= \mathbb{E}[X](1 - \mathbb{E}[X^2]) \\ &= 0\end{aligned}$$

なので、 $X$  と  $Y$  は無相関である。一方で

$$\begin{aligned}f_Y(0) &= \frac{1}{3} & f_Y(1) &= \frac{2}{3} \\ f_{X,Y}(0,0) &= \frac{1}{3} & f_{X,Y}(1,1) &= \frac{1}{3} & f_{X,Y}(-1,1) &= \frac{1}{3}\end{aligned}$$

なので  $X$  と  $Y$  は独立ではない。

共分散は正の値も負の値も取り得る。大雑把に言うと、

- $X$  と  $Y$  の共分散が正  $\iff X$  が大きいとき  $Y$  も大きい
- $X$  と  $Y$  の共分散が負  $\iff X$  が大きいとき  $Y$  は小さい

という意味になる。

補題 6.4. 任意の確率変数  $X_1, \dots, X_n$  について

$$\text{Var}\left[\sum_i X_i\right] = \sum_i \text{Var}[X_i] + 2 \sum_{i < j} \text{Cov}[X_i, X_j].$$

証明. 補題 6.3 の証明参照。 □

## 6.4 モーメントとモーメント母関数

定義 6.4 (モーメント). 確率変数  $X$  と正の整数  $n \geq 1$  について、

$$\mu_n(X) := \mathbb{E}[X^n]$$

を  $X$  の  $n$  次モーメントという。

定義 6.5 (モーメント母関数 (積率母関数)). 確率変数  $X$  について、

$$M_X(t) := \mathbb{E}[e^{tX}] \quad t \in \mathbb{R}$$



を  $X$  のモーメント母関数という。すべての  $t \in \mathbb{R}$  で  $M_X(t)$  が存在しない場合もある。また、

$$K_X(t) := \log M_X(t)$$

を  $X$  のキュムラント母関数という。

今後は以下の補題を認めることにする。証明にはルベーグ積分の知識が必要である。

**定理 6.3.** 確率変数  $X$  について、ある  $\epsilon > 0$  が存在し、モーメント母関数  $M_X(t)$  が  $t \in (-\epsilon, \epsilon)$  で存在するとき、

$$\begin{aligned} M_X(t) &= \sum_{n \geq 0} \frac{\mathbb{E}[X^n]}{n!} t^n \quad \forall t \in (-\epsilon, \epsilon) \\ \mu_n(X) &= \left. \frac{d^n M_X(t)}{dt^n} \right|_{t=0}. \end{aligned}$$

証明. 前半の証明を与える。離散型確率変数  $X$  について、

$$\begin{aligned} M_X(t) &= \mathbb{E}[e^{tX}] \\ &= \sum_x e^{tx} f_X(x) \\ &= \sum_x \left( \sum_{n \geq 0} \frac{(tx)^n}{n!} \right) f_X(x) \end{aligned}$$

である。ここで、任意の  $t \in (-\epsilon, \epsilon)$  について

$$\begin{aligned} &\sum_x \sum_{n \geq 0} \left| \frac{(tx)^n}{n!} f_X(x) \right| \\ &= \sum_x \sum_{n \geq 0} \frac{|tx|^n}{n!} f_X(x) \\ &= \sum_x e^{|tx|} f_X(x) \\ &\leq \sum_x (e^{tx} + e^{-tx}) f_X(x) \\ &= \mathbb{E}[e^{tX}] + \mathbb{E}[e^{-tX}] \\ &= M_X(t) + M_X(-t) < \infty \end{aligned}$$

よって、無限和

$$\sum_x \sum_{n \geq 0} \frac{(tx)^n}{n!} f_X(x)$$

は任意の  $t \in (-\epsilon, \epsilon)$  について**絶対収束する**。そのため、和の順序を変えても収束値は変化しない。よって、

$$\begin{aligned}
 M_X(t) &= \mathbb{E}[e^{tX}] \\
 &= \sum_x \left( \sum_{n \geq 0} \frac{(tx)^n}{n!} \right) f_X(x) \\
 &= \sum_{n \geq 0} \sum_x \frac{(tx)^n}{n!} f_X(x) \\
 &= \sum_{n \geq 0} \frac{\sum_x x^n f_X(x)}{n!} t^n \\
 &= \sum_{n \geq 0} \frac{\mathbb{E}[X^n]}{n!} t^n \quad \forall t \in (-\epsilon, \epsilon).
 \end{aligned}$$

□

**系 6.1.** 確率変数  $X$  について、ある  $\epsilon > 0$  が存在し、モーメント母関数  $M_X(t)$  が  $t \in (-\epsilon, \epsilon)$  で**存在する**とき、

$$\begin{aligned}
 \left. \frac{dK_X(t)}{dt} \right|_{t=0} &= \mathbb{E}[X] \\
 \left. \frac{d^2 K_X(t)}{dt^2} \right|_{t=0} &= \text{Var}[X].
 \end{aligned}$$

証明.

$$\begin{aligned}
 \left. \frac{dK_X(t)}{dt} \right|_{t=0} &= \left. \frac{M'_X(t)}{M_X(t)} \right|_{t=0} = \mathbb{E}[X] \\
 \left. \frac{d^2 K_X(t)}{dt^2} \right|_{t=0} &= \left. \frac{M''_X(t)M_X(t) - M'_X(t)^2}{M_X(t)^2} \right|_{t=0} = M''_X(0) - M'_X(0)^2 = \text{Var}[X].
 \end{aligned}$$

□

また、重要度は低くなるが、

$$\begin{aligned}
 \left. \frac{d^3 K_X(t)}{dt^3} \right|_{t=0} &= \mathbb{E}[(X - \mathbb{E}[X])^3] \\
 \left. \frac{d^4 K_X(t)}{dt^4} \right|_{t=0} &= \mathbb{E}[(X - \mathbb{E}[X])^4] - 3\text{Var}[X]^2
 \end{aligned}$$

が成り立つ。一般に

$$\kappa_n(X) := \left. \frac{d^n K_X(t)}{dt^n} \right|_{t=0}$$

を  $X$  の  $n$  次キュムラントと呼ぶ。

**定理 6.4.** 確率変数  $X$  と  $Y$  のモーメント母関数が  $0$  を含む開区間  $(-\epsilon, \epsilon)$  で存在し、それらが等しいとき、 $X$  の分布と  $Y$  の分布は等しい。

証明.  $\text{Image}(X)$  と  $\text{Image}(Y)$  が有限の場合に限って証明を与える (この場合はモーメント母関数は  $\mathbb{R}$  全体で存在するのだが)。

$$\{x_0, \dots, x_{N-1}\} := \text{Image}(X) \cup \text{Image}(Y)$$

とする。

$$M_X(t) = \sum_{k=0}^{N-1} f_X(x_k) e^{tx_k}$$

$$M_Y(t) = \sum_{k=0}^{N-1} f_Y(x_k) e^{tx_k}$$

なので、

$$0 = M_X(t) - M_Y(t) = \sum_{k=0}^{N-1} (f_X(x_k) - f_Y(x_k)) e^{tx_k} \quad \forall t \in (-\epsilon, \epsilon)$$

各  $k \in \{0, 1, \dots, N-1\}$  について、 $t_k := \epsilon \frac{k}{N}$  とおくと、

$$\sum_{k=0}^{N-1} (f_X(x_k) - f_Y(x_k)) e^{t_\ell x_k} = 0 \quad \forall \ell \in \{0, 1, \dots, N-1\} \quad (6.1)$$

である。ここで、 $N \times N$  実行列  $V$  を

$$V_{\ell k} = e^{t_\ell x_k} = e^{\frac{\epsilon x_k}{N} \ell} \quad \forall k, \ell \in \{0, 1, \dots, N-1\}$$

とおく。この行列  $V$  は Vandermonde 行列の転置であり正則なので、

$$\sum_{k=0}^{N-1} V_{\ell k} g_k = 0 \quad \forall \ell \in \{0, 1, \dots, N-1\}$$

$$\Rightarrow g_k = 0 \quad \forall k \in \{0, 1, \dots, N-1\}$$

よって、

$$f_X(x_k) = f_Y(x_k) \quad \forall k \in \{0, 1, \dots, N-1\}$$

である。 □

定理 6.4 より、モーメント母関数には確率変数の分布のすべての情報が含まれていると言える。ただし、モーメント母関数は原点まわりで存在しないこともあるので、分布の情報をすべて含む関数としては**特性関数**

$$\varphi_X(t) := \mathbb{E}[e^{itX}] \quad \forall t \in \mathbb{R}$$

の方が優秀である。特性関数は常に存在する。一方でモーメント母関数は確率の集中を示す文脈では中心的な役割を果たす。

## 第 II 部

# 統計学

## 第 7 章

# ベイズ推定

統計的推論とは現実の推定問題を確率論に基づきモデル化し、誤り確率を最小化するように推論する方法論である。統計的推論には大きく分けて二種類の流派がある。

- ベイズ主義: 推論する対象の分布 (事前分布) を仮定する。
- 頻度主義: 推論する対象の分布 (事前分布) を仮定しない。

例えば統計的推論は以下のような問題に適用されている。

### 7.1 最大事後確率推定

データが取り得る値の集合を  $\mathcal{X}$  とし、分布のパラメータの取り得る値の集合を  $\Theta$  とする。簡単のため、 $\mathcal{X}$  と  $\Theta$  は高々可算集合とする。データ  $x \in \mathcal{X}$  からパラメータ  $\theta \in \Theta$  を推定する問題を考える。このとき、 $x$  と  $\theta$  が何かしらの確率分布に従っていると仮定する。パラメータ  $\theta$  に対する  $x$  の確率質量関数を  $p(x | \theta)$  と表す。また、パラメータ  $\theta$  の確率質量関数を  $\pi(\theta)$  と表す。つまり、パラメータ  $\theta \in \Theta$  とデータ  $x \in \mathcal{X}$  が選ばれる確率は

$$\pi(\theta)p(x | \theta)$$

である。また、

$$p(x) = \sum_{\theta \in \Theta} \pi(\theta)p(x | \theta), \quad p(\theta | x) := \frac{\pi(\theta)p(x | \theta)}{p(x)}$$

と定義する。ベイズ推定の文脈では

- $\pi(\theta)$ : 事前確率
- $p(x | \theta)$ : 尤度
- $p(\theta | x)$ : 事後確率

と呼ぶ。

得られたデータ  $x \in \mathcal{X}$  からパラメータ  $\theta \in \Theta$  を推定する関数  $\hat{\theta}: \mathcal{X} \rightarrow \Theta$  を**推定量 (estimator)** もしくは推定関数という。推定量  $\hat{\theta}$  の誤り確率を

$$P_{\text{err}}(\hat{\theta}) := \sum_{\theta \in \Theta} \pi(\theta) \sum_{x \in \mathcal{X}} p(x | \theta) \mathbb{I}\{\hat{\theta}(x) \neq \theta\}$$

と定義する。このとき、

$$\begin{aligned} P_{\text{err}}(\hat{\theta}) &= \sum_{\theta \in \Theta} \pi(\theta) \sum_{x \in \mathcal{X}} p(x | \theta) \mathbb{I}\{\hat{\theta}(x) \neq \theta\} \\ &= \sum_{\theta \in \Theta} \pi(\theta) \sum_{x \in \mathcal{X}} p(x | \theta) (1 - \mathbb{I}\{\hat{\theta}(x) = \theta\}) \\ &= 1 - \sum_{\theta \in \Theta} \pi(\theta) \sum_{x \in \mathcal{X}} p(x | \theta) \mathbb{I}\{\hat{\theta}(x) = \theta\} \\ &= 1 - \sum_{x \in \mathcal{X}} \sum_{\theta \in \Theta} \pi(\theta) p(x | \theta) \mathbb{I}\{\hat{\theta}(x) = \theta\} \\ &= 1 - \sum_{x \in \mathcal{X}} \pi(\hat{\theta}(x)) p(x | \hat{\theta}(x)) \\ &\geq 1 - \sum_{x \in \mathcal{X}} \max_{\theta \in \Theta} \pi(\theta) p(x | \theta) \end{aligned}$$

と下から抑えることができ、

$$\begin{aligned} \hat{\theta}_{\text{MAP}}(x) &:= \arg \max_{\theta \in \Theta} p(\theta | x) \\ &= \arg \max_{\theta \in \Theta} p(\theta | x) p(x) \\ &= \arg \max_{\theta \in \Theta} \pi(\theta) p(x | \theta) \end{aligned}$$

という推定量により等号が達成される。この推定量を最大事後確率 (maximum a posteriori; MAP) 推定量と呼ぶ。

## 7.2 最尤推定

MAP 推定は誤り確率を最小化する推定方法であるが、事前確率  $\pi(\theta)$  を仮定しないと用いることができない。一方で尤度を最大化する推定量

$$\hat{\theta}_{\text{ML}}(x) := \arg \max_{\theta \in \Theta} p(x | \theta)$$

を最尤推定 (maximum a priori; ML) 量という。 $\Theta$  が有限集合で、事前確率が一様分布  $\pi(\theta) = \frac{1}{|\Theta|}$  のとき、最尤推定は最大事後確率推定と一致する。

## 7.3 全変動距離

特にパラメータが二値である場合を考える。この章では  $\Theta = \{0, 1\}$  とする。また、

$$p^{(0)}(x) := p(x | 0) \qquad p^{(1)}(x) := p(x | 1)$$

とする。このとき、

$$\begin{aligned} P_{\text{err}}(\hat{\theta}_{\text{MAP}}) &= 1 - \sum_{x \in \mathcal{X}} \max_{\theta \in \{0,1\}} \pi(\theta) p(x | \theta) \\ &= 1 - \sum_{x \in \mathcal{X}} \max \{ \pi(0)p^{(0)}(x), \pi(1)p^{(1)}(x) \} \end{aligned}$$

である。ここで、

$$\begin{aligned} \max\{a, b\} - \min\{a, b\} &= |a - b| \\ \max\{a, b\} + \min\{a, b\} &= a + b \end{aligned} \quad \forall a, b \in \mathbb{R}$$

であるので、

$$\max\{a, b\} = \frac{1}{2}(a + b + |a - b|) \quad \forall a, b \in \mathbb{R}.$$

よって、

$$\max \{ \pi(0)p^{(0)}(x), \pi(1)p^{(1)}(x) \} = \frac{1}{2} (p(x) + |\pi(0)p^{(0)}(x) - \pi(1)p^{(1)}(x)|)$$

よって、

$$\begin{aligned} P_{\text{err}}(\hat{\theta}_{\text{MAP}}) &= 1 - \sum_{x \in \mathcal{X}} \max \{ \pi(0)p^{(0)}(x), \pi(1)p^{(1)}(x) \} \\ &= 1 - \sum_{x \in \mathcal{X}} \frac{1}{2} (p(x) + |\pi(0)p^{(0)}(x) - \pi(1)p^{(1)}(x)|) \\ &= 1 - \frac{1}{2} \left( 1 + \sum_{x \in \mathcal{X}} |\pi(0)p^{(0)}(x) - \pi(1)p^{(1)}(x)| \right) \\ &= \frac{1}{2} \left( 1 - \sum_{x \in \mathcal{X}} |\pi(0)p^{(0)}(x) - \pi(1)p^{(1)}(x)| \right) \end{aligned}$$

**定義 7.1** (全変動距離). 高々可算集合  $\mathcal{X}$  上の関数  $f$  について、

$$\|f\|_1 := \sum_{x \in \mathcal{X}} |f(x)|$$

と定義する。また、 $\mathcal{X}$  上の確率質量関数  $p^{(0)}, p^{(1)}$  について、

$$d_{\text{TV}}(p^{(0)}, p^{(1)}) := \frac{1}{2} \|p^{(0)} - p^{(1)}\|_1$$

を  $p^{(0)}$  と  $p^{(1)}$  の**全変動距離**という。

これらの記法を用いると、

$$P_{\text{err}}(\hat{\theta}_{\text{MAP}}) = \frac{1}{2} (1 - \|\pi(0)p^{(0)} - \pi(1)p^{(1)}\|_1)$$

と表せる。また、 $\pi(0) = \pi(1) = 1/2$  のとき、

$$P_{\text{err}}(\hat{\theta}_{\text{MAP}}) = \frac{1}{2} (1 - d_{\text{TV}}(p^{(0)}, p^{(1)}))$$

である。

## 7.4 損失関数

パラメータが取り得る値の集合が実数の部分集合  $\Theta \subseteq \mathbb{R}$  であると仮定する。真のパラメータとその推定値の間の「誤差」を表す関数  $L: \Theta \times \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  を**損失関数**と呼ぶ。また、**期待損失**  $R: \Theta \times (\mathcal{X} \rightarrow \mathbb{R}) \rightarrow \mathbb{R}_{\geq 0}$  を

$$\begin{aligned} R(\theta, \hat{\theta}) &:= \mathbb{E} [L(\theta, \hat{\theta}(X)) \mid \theta] \\ &= \sum_{x \in \mathcal{X}} L(\theta, \hat{\theta}(x)) p(x \mid \theta) \end{aligned}$$

と定義する。また、**ベイズリスク**  $\rho: \mathcal{P}(\Theta) \times (\mathcal{X} \rightarrow \mathbb{R}) \rightarrow \mathbb{R}_{\geq 0}$  を

$$\begin{aligned} \rho(\pi, \hat{\theta}) &:= \mathbb{E} [L(\theta, \hat{\theta}(X))] \\ &= \sum_{\theta \in \Theta} \sum_{x \in \mathcal{X}} L(\theta, \hat{\theta}(x)) p(x \mid \theta) \pi(\theta) \end{aligned}$$

と定義する。 $\Theta \subseteq \mathbb{R}$  が非可算無限集合の場合は  $\pi(\theta)$  を確率密度関数とし、

$$\begin{aligned} \rho(\pi, \hat{\theta}) &:= \mathbb{E} [L(\theta, \hat{\theta}(X))] \\ &= \int \sum_{x \in \mathcal{X}} L(\theta, \hat{\theta}(x)) p(x \mid \theta) \pi(\theta) d\theta \end{aligned}$$

と定義する。このとき、 $p(x \mid \theta)$  は条件付き確率というより、 $\theta \in \Theta$  というパラメータを持った確率質量関数 ( $\theta \in \Theta$  から定まる確率質量関数) と理解すれば十分である。

例えば  $\Theta$  が高々可算集合で、

$$L(\theta, \theta') = \mathbb{I}\{\theta \neq \theta'\} \quad \forall \theta, \theta' \in \Theta$$

と定義すると、ベイズリスク  $\rho(\pi, \hat{\theta})$  は推定量  $\hat{\theta}$  の誤り確率  $P_{\text{err}}(\hat{\theta})$  である。

その他の重要な損失関数の例として二乗誤差がある。

$$L(\theta, \theta') = (\theta - \theta')^2.$$

損失関数  $L$  を定めたときに、ベイズリスクを最小化する推定量

$$\begin{aligned} \hat{\theta} &= \arg \min_{\hat{\theta}} \rho(\pi, \hat{\theta}) \\ \Leftrightarrow \hat{\theta}(x) &= \arg \min_{\theta' \in \Theta} \sum_{\theta \in \Theta} L(\theta, \theta') p(x \mid \theta) \pi(\theta) \quad \forall x \in \mathcal{X} \\ \Leftrightarrow \hat{\theta}(x) &= \arg \min_{\theta' \in \Theta} \sum_{\theta \in \Theta} L(\theta, \theta') p(\theta \mid x) p(x) \quad \forall x \in \mathcal{X} \\ \Leftrightarrow \hat{\theta}(x) &= \arg \min_{\theta' \in \Theta} \sum_{\theta \in \Theta} L(\theta, \theta') p(\theta \mid x) \quad \forall x \in \mathcal{X} \end{aligned}$$

ここで、

$$\mathbb{E} [L(\theta, \theta') \mid x] = \sum_{\theta \in \Theta} L(\theta, \theta') p(\theta \mid x) \quad \forall x \in \mathcal{X}$$



を損失関数の事後平均という。各  $x \in \mathcal{X}$  について、 $\theta' = \hat{\theta}(x)$  が損失関数の事後平均を最小化するとき、ベイズリスクを最小化する。

損失関数が  $L(\theta, \theta')$  が各  $\theta \in \Theta$  を固定したときに  $\theta'$  について凸関数であるとき、事後平均  $\mathbb{E}[L(\theta, \theta') \mid x]$  も  $\theta'$  について凸関数となる (凸関数の非負倍は凸関数であり、凸関数の和は凸関数なので)。さらに、 $L(\theta, \theta')$  が  $\theta'$  について微分可能なとき、

$$\frac{\partial \mathbb{E}[L(\theta, \theta') \mid x]}{\partial \theta'} = \sum_{\theta \in \Theta} \frac{\partial L(\theta, \theta')}{\partial \theta'} p(\theta \mid x) = 0 \quad \forall x \in \mathcal{X}$$

を満たす  $\theta'$  を  $\hat{\theta}(x)$  として選択するのが最適である。二乗誤差  $L(\theta, \theta') = (\theta - \theta')^2$  のとき、この条件は

$$0 = \sum_{\theta} 2(\theta' - \theta)p(\theta \mid x) = 2 \left( \theta' - \sum_{\theta} \theta p(\theta \mid x) \right)$$

となり、

$$\hat{\theta}(x) = \sum_{\theta} \theta p(\theta \mid x)$$

とするのが最適であることが分かる。この右辺の値をパラメータ  $\theta$  の事後平均という。

## 第 8 章

# 仮説検定

### 8.1 事前分布を仮定しない推定問題

ベイズ推定では推定したいパラメータ  $\theta \in \Theta$  の事前分布  $\pi(\theta)$  を既知として仮定した。しかし、現実の問題ではこの事前分布を適切に仮定する方法がない場合もある。例えば

- 開発中の薬に効果があるかないか
- ある患者が病気かどうか
- サイコロに偏りがあるかどうか

といった問題について事前分布をどのように仮定するのが適切か不明瞭である。そのような状況でデータ  $x \in \mathcal{X}$  からパラメータ  $\theta \in \Theta$  を推定する問題を考える。尤度  $p(x | \theta)$  は既知とする。

### 8.2 伝統的な仮説検定

パラメータの集合を二つの部分集合  $\Theta_0$  と  $\Theta_1$  に分割する。つまり、 $\Theta_0 \cup \Theta_1 = \Theta$ ,  $\Theta_0 \cap \Theta_1 = \emptyset$  である。そしてパラメータが  $\Theta_0$  に属するか  $\Theta_1$  に属するかを知りたいとする。このとき二つの命題

$$\begin{aligned}H_0: \theta \in \Theta_0 \\ H_1: \theta \in \Theta_1\end{aligned}$$

を仮説という。その二つの仮説のうちの**通常成り立っていると考える方**を  $H_0$  とし (対応するパラメータ集合は  $\Theta_0$ ) **帰無仮説**と呼ぶ。また、そうでない方を  $H_1$  とし (対応するパラメータ集合は  $\Theta_1$ ) **対立仮説**と呼ぶ。帰無仮説の例として

- 開発中の薬に効果はない
- ある患者が病気ではない
- サイコロに偏りはない

などがある。それらに対応する対立仮説はそれぞれ

- 開発中の薬に効果がある

- ある患者が病気である
- サイコロに偏りがある

となる。仮説検定の考え方では帰無仮説を棄却するかしないかを定める。帰無仮説を棄却した場合、対立仮説を正しいと考え、帰無仮説を棄却しなかった場合は何も言えないと結論づける。

### 8.3 単純仮説

仮説  $\Theta_0$  と  $\Theta_1$  がそれぞれ一元集合であるとき、 $H_0$  と  $H_1$  を**単純仮説**という。 $\Theta_0 = \{\theta_0\}$ ,  $\Theta_1 = \{\theta_1\}$  として、 $p_0(x) := p(x | \theta_0)$  と  $p_1(x) := p(x | \theta_1)$  とする。データから仮説を推定する関数  $E: \mathcal{X} \rightarrow [0, 1]$  を**検定関数**という。各  $x \in \mathcal{X}$  について、 $E(x)$  は**帰無仮説を棄却する確率**とする。この検定関数について二種類の誤り確率を

$$\alpha_E := \mathbb{E}[E(X) | \theta_0] = \sum_x E(x)p(x | \theta_0)$$

$$\beta_E := 1 - \mathbb{E}[E(X) | \theta_1] = 1 - \sum_x E(x)p(x | \theta_1)$$

と定義する。このとき、 $\alpha_E$  は帰無仮説が正しいときに帰無仮説を棄却する確率であり、**第一種誤り確率**もしくは**有意水準**という。また、 $\beta_E$  は対立仮説が正しいときに帰無仮説を棄却しない確率であり、**第二種誤り確率**という。第一種誤り確率だけを小さくしたければ  $E(x) = 0$  とすればよいし、第二種誤り確率だけを小さくしたければ  $E(x) = 1$  とすればよい。

### 8.4 最強力検定

**定義 8.1.** 検定関数  $E: \mathcal{X} \rightarrow [0, 1]$  が有意水準  $\alpha \in [0, 1]$  の**最強力検定**  $\stackrel{\text{def}}{\iff} \alpha_E = \alpha$  であり、任意の  $F: \mathcal{X} \rightarrow [0, 1]$  について、 $\alpha_F \leq \alpha$  ならば  $\beta_F \geq \beta$  が成り立つ。

実現可能な誤り確率  $(\alpha, \beta)$  の集合

$$C = \{(\alpha_E, \beta_E) | E: \mathcal{X} \rightarrow [0, 1]\}$$

について考える。自明に  $(1, 0)$  と  $(0, 1)$  は実現可能である。

この集合  $C$  は凸集合である。

$$\alpha_{pE+(1-p)F} = p\alpha_E + (1-p)\alpha_F$$

$$\beta_{pE+(1-p)F} = p\beta_E + (1-p)\beta_F$$

であることから、

$$(\alpha_{pE+(1-p)F}, \beta_{pE+(1-p)F}) = p(\alpha_E, \beta_E) + (1-p)(\alpha_F, \beta_F)$$

が確認できる。また、検定結果を反転させた検定関数  $1 - E(x)$  を考えると、

$$(\alpha_{1-E}, \beta_{1-E}) = (1, 1) - (\alpha_E, \beta_E)$$

である。

まとめると、実現可能な  $(\alpha, \beta)$  の集合  $C$  は

1.  $(1, 0), (0, 1) \in C$
2.  $C$  は 凸集合
3.  $(\alpha, \beta) \in C \iff (1 - \alpha, 1 - \beta) \in C$

を満たす。



図 8.1: 実現可能な  $(\alpha, \beta)$  の集合の例

各  $\alpha \in [0, 1]$  について、 $\beta$  を最小化するのが最強力検定であるので、この図の下のカークが最強力検定で実現される  $(\alpha, \beta)$  となる。

## 8.5 尤度比検定

ベイズ推定の枠組みでは MAP 推定量が誤り確率を最小化する推定量であった。この MAP 推定量は

$$E_{\text{MAP}}(x) = \begin{cases} 0 & \text{if } \frac{p(x|\theta_0)}{p(x|\theta_1)} \geq \frac{\pi(\theta_1)}{\pi(\theta_0)} \\ 1 & \text{otherwise} \end{cases}$$

と表すことができる。

一般的に  $\eta > 0, \kappa \in [0, 1]$  について

$$E(x) = \begin{cases} 0 & \text{if } \frac{p(x|\theta_0)}{p(x|\theta_1)} > \eta \\ 1 & \text{if } \frac{p(x|\theta_0)}{p(x|\theta_1)} < \eta \\ \kappa & \text{otherwise} \end{cases}$$

という形の検定関数を尤度比検定という。

**補題 8.1** (ネイマン・ピアソンの補題). 任意の尤度比検定は最強力検定である (逆も成り立つ)。

証明. 尤度比検定  $E$  における尤度比の閾値を  $\eta > 0$  とする。任意の  $F: \mathcal{X} \rightarrow [0, 1]$  について

$$\begin{aligned} & (F(x) - E(x))(p_0(x) - \eta p_1(x)) \geq 0 \quad \forall x \in \mathcal{X} \\ \Leftrightarrow & F(x)p_0(x) - E(x)p_0(x) - \eta F(x)p_1(x) + \eta E(x)p_1(x) \geq 0 \quad \forall x \in \mathcal{X} \\ \Rightarrow & \alpha_F - \alpha_E - \eta(1 - \beta_F) + \eta(1 - \beta_E) \geq 0 \\ \Leftrightarrow & (\alpha_F - \alpha_E) + \eta(\beta_F - \beta_E) \geq 0. \end{aligned}$$

よって  $\alpha_F \leq \alpha_E$  ならば  $\beta_F \geq \beta_E$  である。 □

データ  $\mathcal{X}$  が連続な場合は確率質量関数の代わりに確率密度関数を考える。

**例 8.1.** コインを持っており、表が出る確率は  $p_0$  か  $p_1$  のどちらかである。コインを独立に  $N$  回投げて表が出る確率が  $p_0$  か  $p_1$  かを推定したい。 $\mathcal{X} = \{0, 1\}^N$  とし、0 は裏、1 は表に対応するものとする。このとき尤度は

$$p(\mathbf{x} | p_k) = \prod_{i=1}^N p_k^{x_i} (1 - p_k)^{1-x_i} \quad \text{for } k \in \{0, 1\}$$

である。このとき尤度比は

$$\begin{aligned} \frac{p(\mathbf{x} | p_0)}{p(\mathbf{x} | p_1)} &= \frac{\prod_i p_0^{x_i} (1 - p_0)^{1-x_i}}{\prod_i p_1^{x_i} (1 - p_1)^{1-x_i}} \\ &= \left(\frac{p_0}{p_1}\right)^{\sum_i x_i} \left(\frac{1-p_0}{1-p_1}\right)^{N-\sum_i x_i} \end{aligned}$$

である。よって尤度比は表が出た回数  $T(\mathbf{x}) = \sum_i x_i$  から定まる。なので  $p_1 > p_0$  とすると、表が出た回数が多いときに  $E(\mathbf{x}) = 1$  とすることになる。

## 第 III 部

## 漸近論

## 第 9 章

# 大数の法則と集中不等式

### 9.1 大数の弱法則

表が出る確率が  $1/2$  のコインを 100 回独立に投げたときに表が出る回数は大体 50 回くらいになるだろう。それを一般的な形で述べたものが大数の法則である。

**定理 9.1** (大数の弱法則 (分散有限を仮定)). 確率変数  $X$  が分散を持つとする。確率変数  $X_1, \dots, X_N$  が独立同分布で  $X$  と同じ分布に従うとする。このとき、任意の  $\epsilon > 0$  について

$$\lim_{N \rightarrow \infty} \Pr \left( \left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X] \right| \geq \epsilon \right) = 0$$

が成り立つ。

証明.

$$\begin{aligned} \Pr \left( \left| \frac{1}{n} \sum_{i=1}^N X_i - \mathbb{E}[X] \right| \geq \epsilon \right) &= \Pr \left( \left( \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X] \right)^2 \geq \epsilon^2 \right) \\ &= \Pr \left( \left( \sum_{i=1}^N X_i - N\mathbb{E}[X] \right)^2 \geq \epsilon^2 N^2 \right) \\ &\leq \frac{N\text{Var}[X]}{\epsilon^2 N^2} = \frac{\text{Var}[X]}{\epsilon^2 N} \rightarrow 0. \end{aligned}$$

□

### 9.2 チェルノフ上界

上記の大数の弱法則の証明では確率が 0 に収束するスピードは  $O(1/N)$  であった。より詳しく確率が 0 にくくスピードを解析しよう。

**補題 9.1** (チェルノフ上界). 任意の確率変数  $X$  と  $a \in \mathbb{R}$  について、

$$\begin{aligned}\Pr(X \geq a) &\leq \frac{M_X(t)}{e^{at}} = e^{K_X(t)-at} \quad \forall t \geq 0 \\ \Pr(X \leq a) &\leq \frac{M_X(t)}{e^{at}} = e^{K_X(t)-at} \quad \forall t \leq 0.\end{aligned}$$

証明.  $t = 0$  のときは不等式の右辺は 1 となるので、不等式は自明に成り立つ。任意の  $t > 0$  について、

$$\begin{aligned}\Pr(X \geq a) &= \Pr(e^{tX} \geq e^{ta}) \\ &\leq \frac{M_X(t)}{e^{at}} \quad (\text{マルコフの不等式}).\end{aligned}$$

が成り立つ。もう一つの不等式も同様に示すことができる。  $\square$

マルコフの不等式は非負の確率変数にしか適用できないが、チェルノフ上界は任意の確率変数に適用できる。

**補題 9.2** (確率変数の和に対するチェルノフ上界). 任意の確率変数  $X$  と  $a \in \mathbb{R}$  について、

$$\begin{aligned}\Pr\left(\frac{1}{N} \sum_{i=1}^N X_i \geq a\right) &\leq e^{-(at - K_X(t))N} \quad \forall t \geq 0 \\ \Pr\left(\frac{1}{N} \sum_{i=1}^N X_i \leq a\right) &\leq e^{-(at - K_X(t))N} \quad \forall t \leq 0.\end{aligned}$$

証明.

$$\begin{aligned}\Pr\left(\frac{1}{N} \sum_{i=1}^N X_i \geq a\right) &= \Pr\left(\sum_{i=1}^N X_i \geq aN\right) \\ &\leq e^{K_{\sum_i X_i}(t) - atN} \quad (\text{チェルノフ上界}) \\ &= e^{(K_X(t) - at)N} \quad \left(K_{\sum_i X_i}(t) = \sum_i K_{X_i}(t) = NK_X(t)\right).\end{aligned}$$

$\square$

このようにチェルノフ上界を使うと  $N$  について指数関数の上界が得られる。係数  $at - K_X(t)$  が正であれば、確率は指数関数的に小さいことになる。ここで、最適な  $t$  を選ぶことで、この係数  $at - K_X(t)$  を最大化することを考える。

### 9.3 キュムラント母関数の性質

キュムラント母関数

$$K_X(t) = \log M_X(t) = \log \mathbb{E}[e^{tX}]$$

の性質を改めて考えよう。発散する場合は  $+\infty$  に値を取るとみなして  $K_X: \mathbb{R} \rightarrow (-\infty, +\infty]$  と考えることにする。



まず、 $K_X(0) = 0$  である。

ある  $t > 0$  について、 $M_X(t) < +\infty$  と仮定すると、任意の  $s \in (0, t)$  について

$$\begin{aligned} M_X(s) &= \mathbb{E}[e^{sX}] \\ &= \mathbb{E}[e^{sX} \mathbb{1}_{\{X \geq 0\}}] + \mathbb{E}[e^{sX} \mathbb{1}_{\{X < 0\}}] \\ &\leq \mathbb{E}[e^{tX} \mathbb{1}_{\{X \geq 0\}}] + 1 \\ &\leq M_X(t) + 1 < +\infty \end{aligned}$$

である。同様に、ある  $t < 0$  について、 $M_X(t) < +\infty$  と仮定すると、任意の  $s \in (t, 0)$  について  $M_X(s) < +\infty$  である。よって、 $M_X(t)$  や  $K_X(t)$  が有限となる範囲は  $0$  を含む区間となる。ここでいう区間とは一般的に空集合、もしくは  $a < b$  について、

$$[a, a] \quad (a, b) \quad [a, b) \quad (a, b] \quad [a, b] \quad (a, +\infty) \quad [a, +\infty) \quad (-\infty, b) \quad (-\infty, b] \quad (-\infty, +\infty)$$

のいずれかの形の集合を指す。この区間を

$$\text{dom}(K_X) := \{t \in \mathbb{R} \mid K_X(t) < +\infty\}$$

と表す。区間は 1 次元の凸集合と一言で理解できる。

証明はしないが、キュムラント母関数  $K_X(t)$  は  $\text{dom}(K_X)$  の内点で何回でも微分可能であり、無限和や積分を取る前に微分しても構わない。

$$\begin{aligned} \frac{dK_X(t)}{dt} &= \frac{\mathbb{E}[Xe^{tX}]}{\mathbb{E}[e^{tX}]} \\ \frac{d^2K_X(t)}{dt^2} &= \frac{\mathbb{E}[X^2e^{tX}]\mathbb{E}[e^{tX}] - \mathbb{E}[Xe^{tX}]^2}{\mathbb{E}[e^{tX}]^2} \\ &= \frac{\mathbb{E}[X^2e^{tX}]}{\mathbb{E}[e^{tX}]} - \left( \frac{\mathbb{E}[Xe^{tX}]}{\mathbb{E}[e^{tX}]} \right)^2 \\ &= \frac{\mathbb{E}\left[\left(X - \frac{\mathbb{E}[Xe^{tX}]}{\mathbb{E}[e^{tX}]}\right)^2 e^{tX}\right]}{\mathbb{E}[e^{tX}]} \geq 0 \end{aligned}$$

ここで確率変数  $Z_t$  を導入し、確率質量関数

$$f_{Z_t}(x) = \frac{f_X(x)e^{tx}}{\sum_x f_X(x)e^{tx}}$$

を持つものとする、

$$\begin{aligned} \frac{dK_X(t)}{dt} &= \mathbb{E}[Z_t] \\ \frac{d^2K_X(t)}{dt^2} &= \text{Var}[Z_t] \end{aligned}$$

であることが分かる。また、 $X$  が決定的 ( $\Pr(X = \mathbb{E}[X]) = 1$ ) でない限り、 $K_X$  は  $\text{dom}(K_X)$  で狭義凸である。

よって  $X$  のキュムラント母関数  $K_X(t)$  が原点付近で存在すると仮定すると、 $K_X(t)$  は

- 0 を含む区間で定義され、
- 原点を通り、
- 凸関数で、
- 原点の傾きは  $\mathbb{E}[X]$

であることが分かる。

### i ノート

関数  $f: \mathbb{R} \rightarrow (-\infty, +\infty]$  が下半連続であるとは、任意の  $\alpha \in \mathbb{R}$  について  $\{t \in \mathbb{R} \mid f(t) \leq \alpha\}$  が閉集合であることをいう。キュムラント母関数は下半連続の凸関数である。

また、キュムラント母関数の簡単な性質として以下が成り立つ。任意の独立確率変数  $X$  と  $Y$  と  $a \in \mathbb{R}$  について

$$\begin{aligned} K_{X+a}(t) &= \log \mathbb{E}[e^{t(X+a)}] = \log(\mathbb{E}[e^{tX}] \cdot e^{ta}) = K_X(t) + at \\ K_{aX}(t) &= \log \mathbb{E}[e^{t(aX)}] = K_X(at) \\ K_{X+Y}(t) &= \log \mathbb{E}[e^{t(X+Y)}] = \log(\mathbb{E}[e^{tX}] \mathbb{E}[e^{tY}]) = K_X(t) + K_Y(t). \end{aligned}$$

## 9.4 キュムラント母関数の例

## 9.5 ルジャンドル変換

チェルノフ上界に現れる係数  $at - K_X(t)$  の最大化はルジャンドル変換を用いて  $K_X^*(a)$  と表せる。ルジャンドル変換を定義する際は  $+\infty$  という値を許して  $(-\infty, +\infty]$  を値域として考えると都合がよい。このように  $+\infty$  を値として許した場合にも凸性を通常関数と同じように定義する。一般に区間上で定義された凸関数を実数全体に拡張し、元の定義域の外で  $+\infty$  を取ることにするとやはり凸関数になる。また、 $f: \mathbb{R} \rightarrow (-\infty, +\infty]$  の有効領域を

$$\text{dom}(f) := \{x \in \mathbb{R} \mid f(x) < +\infty\}$$

と定義する。凸関数の有効領域は区間になる。

**定義 9.1** (ルジャンドル変換).

恒等的に  $+\infty$  ではない凸関数  $f: \mathbb{R} \rightarrow (-\infty, +\infty]$  のルジャンドル変換  $f^*: \mathbb{R} \rightarrow (-\infty, +\infty]$  を以下で定義する。

$$f^*(a) := \sup_{t \in \mathbb{R}} \{at - f(t)\}.$$

ルジャンドル変換において  $f$  は凸関数なので、 $at - f(t)$  は凹関数 (上に凸の関数) になる。簡単のため  $f$  が  $\text{dom}(f)$  の内点で微分可能であると仮定しよう (キュムラント母関数はこの仮定は満たす)。このとき、

- $at - f(t)$  の微分が 0 になる点、つまり  $f'(t_a) = a$  を満たす  $t_a \in \text{dom}(f)^\circ$  が存在するとき、 $t_a$  で  $at - f(t)$  は最大化される。

- そのような  $t_a \in \text{dom}(f)^\circ$  が存在しないときは、 $\text{dom}(f)$  の端への極限で  $\sup$  が達成される。

凸関数は  $\text{dom}(f)$  の上では接線の集合で表すことができる。接線  $ax + b$  は傾き  $a$  と切片  $b$  のペアで表すことができる。傾き  $a$  を持つ  $f$  の接線は

$$a(x - t_a) + f(t_a) = ax - (at_a - f(t_a)) = ax - f^*(a)$$

この傾き  $a$  から接線  $ax + b$  の切片の  $-1$  倍である  $-b$  への関数が  $f^*$  である。

例えば  $f$  が  $0$  で微分可能であるとき

$$f^*(f'(0)) = 0$$

である。

また、ルジャンドル変換  $f^*$  は凸関数である。

**補題 9.3.** 凸関数  $f: \mathbb{R} \rightarrow (-\infty, +\infty]$  のルジャンドル変換  $f^*: \mathbb{R} \rightarrow (-\infty, +\infty]$  は凸関数である。

証明. 任意の  $a_1, a_2 \in \mathbb{R}$  と  $\lambda \in [0, 1]$  について

$$\begin{aligned} f^*(\lambda a_1 + (1 - \lambda)a_2) &= \sup_{t \in \mathbb{R}} \{(\lambda a_1 + (1 - \lambda)a_2)t - f(t)\} \\ &= \sup_{t \in \mathbb{R}} \{\lambda(a_1 t - f(t)) + (1 - \lambda)(a_2 t - f(t))\} \\ &\leq \sup_{t \in \mathbb{R}} \{\lambda(a_1 t - f(t))\} + \sup_{t \in \mathbb{R}} \{(1 - \lambda)(a_2 t - f(t))\} \\ &\leq \lambda f^*(a_1) + (1 - \lambda)f^*(a_2). \end{aligned}$$

□

ルジャンドル変換を凸でない関数  $f$  についても同様に定義した場合でも、ルジャンドル変換  $f^*$  は同様に凸関数となる。

#### ノート

またルジャンドル変換  $f^*$  は下半連続である。ルジャンドル変換は凸で下半連続な関数を凸で下半連続な関数に写す。また、凸で下半連続な関数  $f$  について  $f = f^{**}$  である。

## 9.6 クラメールの定理

**補題 9.4** (イェンセンの不等式). 任意の凸関数  $f: \mathbb{R} \rightarrow \mathbb{R}$  と確率変数  $X$  について

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]).$$

$X$  の像が有限のときの証明. 確率変数  $X$  は  $k = 1, 2, \dots, m$  について確率  $p_k$  で値  $a_k$  をとると仮定する.  $m$  についての帰納法で示す.  $m = 1$  のときは明らかに成り立つ.  $X$  の像のサイズが  $m$  未満のときにイェンセンの不等式が成り立つと仮定すると、

$$\begin{aligned}\mathbb{E}[f(X)] &= \sum_{k=1}^m p_k f(a_k) = \left( \sum_{k=1}^{m-1} p_k \right) \sum_{k=1}^{m-1} \frac{p_k}{\sum_{\ell=1}^{m-1} p_\ell} f(a_k) + p_m f(a_m) \\ &\geq \left( \sum_{k=1}^{m-1} p_k \right) f \left( \sum_{k=1}^{m-1} \frac{p_k}{\sum_{\ell=1}^{m-1} p_\ell} a_k \right) + p_m f(a_m) \quad (\text{帰納法の仮定}) \\ &\geq f \left( \left( \sum_{k=1}^{m-1} p_k \right) \sum_{k=1}^{m-1} \frac{p_k}{\sum_{\ell=1}^{m-1} p_\ell} a_k + p_m a_m \right) \quad (f \text{ の凸性}) \\ &= f \left( \sum_{k=1}^m p_k a_k \right) = f(\mathbb{E}[X]).\end{aligned}$$

□

補題 9.2 の右辺を  $t$  について最小化することを考えよう。

補題 9.5. レート関数  $I_X: \mathbb{R} \rightarrow [0, +\infty]$  を

$$I_X(a) := K_X^*(a) = \sup_{t \in \mathbb{R}} \{at - K_X(t)\}$$

と定義する。このとき、

1.  $\text{dom}(K_X) = \{0\}$  のとき、 $I_X(a) = 0$ .
2. ある  $\epsilon > 0$  が存在して  $K_X(\epsilon) < +\infty$  のとき、 $\mathbb{E}[X] < +\infty$  であり、

$$\sup_{t \geq 0} \{at - K_X(t)\} = \begin{cases} I_X(a) & \text{if } a > \mathbb{E}[X] \\ 0 & \text{otherwise.} \end{cases}$$

3. ある  $\epsilon > 0$  が存在して  $K_X(-\epsilon) < +\infty$  のとき、 $\mathbb{E}[X] > -\infty$  であり、

$$\sup_{t \leq 0} \{at - K_X(t)\} = \begin{cases} I_X(a) & \text{if } a < \mathbb{E}[X] \\ 0 & \text{otherwise.} \end{cases}$$

証明. 1 は自明. 2 を示す. ある  $\epsilon > 0$  について、 $K_X(\epsilon) < +\infty$  と仮定する. 一般に、

$$e^{tx} \geq tx + 1 \quad \forall t, x \in \mathbb{R}$$

より、 $M_X(t) \geq t\mathbb{E}[X] + 1$  である。よって、

$$\mathbb{E}[X] \leq \frac{M_X(\epsilon) - 1}{\epsilon} < +\infty$$

である。また、イェンセンの不等式より、

$$K_X(t) = \log \mathbb{E}[e^{tX}] \geq \mathbb{E}[\log e^{tX}] = t\mathbb{E}[X] \quad \forall t \in \mathbb{R}$$

である。よって、

$$I_X(\mathbb{E}[X]) = \sup_{t \in \mathbb{R}} \{t\mathbb{E}[X] - K_X(t)\} = 0$$

である。任意の  $a > \mathbb{E}[X]$  と  $t < 0$  について、

$$ta - K_X(t) \leq t\mathbb{E}[X] - K_X(t) \leq 0$$

であるので、任意の  $a > \mathbb{E}[X]$  について

$$I_X(a) := \sup_{t \in \mathbb{R}} \{at - K_X(t)\} = \sup_{t \geq 0} \{at - K_X(t)\}$$

また、 $\sup_{t \geq 0} \{at - K_X(t)\}$  は  $a$  について単調なので、任意の  $a < \mathbb{E}[X]$  について

$$\sup_{t \geq 0} \{at - K_X(t)\} = 0$$

である。

3 は 2 と同様に示せる。

□

よってチェルノフ上界を最適化することで確率の指数的な上界を得る。

**定理 9.2** (最適化されたチェルノフ上界).

$$\begin{aligned} \Pr\left(\frac{1}{N} \sum_{i=1}^N X \geq a\right) &\leq e^{-I_X(a)N} & \forall a > \mathbb{E}[X] \\ \Pr\left(\frac{1}{N} \sum_{i=1}^N X \leq a\right) &\leq e^{-I_X(a)N} & \forall a < \mathbb{E}[X]. \end{aligned}$$

この指数は漸近的に最適である。証明は少し難しいので紹介しない。

**定理 9.3** (クラメールの定理).

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{1}{N} \log \Pr\left(\frac{1}{N} \sum_{i=1}^N X \geq a\right) &= -I_X(a) & \forall a > \mathbb{E}[X] \\ \lim_{N \rightarrow \infty} \frac{1}{N} \log \Pr\left(\frac{1}{N} \sum_{i=1}^N X \leq a\right) &= -I_X(a) & \forall a < \mathbb{E}[X]. \end{aligned}$$

## 第 10 章

# 正規分布と中心極限定理

### 10.1 中心極限定理

大数の法則により  $\frac{1}{N} \sum_{k=1}^N X_k$  は期待値周辺に集中することが分かる。また、クラメールの定理によりその集中のスピードも精密に評価することができる。一方で期待値周辺の挙動を

$$\frac{1}{\sqrt{N}} \left( \sum_{k=1}^N X_k - N\mathbb{E}[X] \right)$$

より詳しく見る。

**定理 10.1** (中心極限定理). 分散を持つ確率変数  $X$  とその  $i$

$$\frac{1}{\sqrt{N\text{Var}[X]}} \left( \sum_{k=1}^N X_k - N\mathbb{E}[X] \right)$$

### 10.2 特性関数

**定義 10.1.** 確率変数  $X$  について、特性関数  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  を以下で定義する。

$$\varphi_X(t) := \mathbb{E}[e^{itX}].$$

**定理 10.2.** 確率変数  $X, Y$  について

$$\varphi_X = \varphi_Y \iff F_X = F_Y.$$

確率変数の像が有限の場合の証明.  $\text{Image}(X)$  と  $\text{Image}(Y)$  が有限の場合に限って証明を与える。

$$\{x_0, \dots, x_{N-1}\} := \text{Image}(X) \cup \text{Image}(Y)$$

とする。

$$\varphi_X(t) = \sum_{k=0}^{N-1} f_X(x_k) e^{itx_k}, \quad \varphi_Y(t) = \sum_{k=0}^{N-1} f_Y(x_k) e^{itx_k}$$

なので、

$$0 = \varphi_X(t) - \varphi_Y(t) = \sum_{k=0}^{N-1} (f_X(x_k) - f_Y(x_k)) e^{itx_k} \quad \forall t \in \mathbb{R}.$$

よって、

$$\sum_{k=0}^{N-1} (f_X(x_k) - f_Y(x_k)) e^{i\ell x_k} = 0 \quad \forall \ell \in \{0, 1, \dots, N-1\} \quad (10.1)$$

である。ここで、 $N \times N$  実行列  $V$  を

$$V_{\ell k} = e^{i\ell x_k} \quad \forall k, \ell \in \{0, 1, \dots, N-1\}$$

とおく。この行列  $V$  は Vandermonde 行列の転置であり正則なので、

$$\begin{aligned} \sum_{k=0}^{N-1} V_{\ell k} g_k &= 0 \quad \forall \ell \in \{0, 1, \dots, N-1\} \\ \Rightarrow g_k &= 0 \quad \forall k \in \{0, 1, \dots, N-1\} \end{aligned}$$

よって、

$$f_X(x_k) = f_Y(x_k) \quad \forall k \in \{0, 1, \dots, N-1\}$$

である。 □

### 10.3 特性関数の応用

### 10.4 連続性定理

### 10.5 中心極限定理の証明

## 第 11 章

# サノフの定理、KL ダイバージェンス



## 参考文献

Knuth, Donald E. 1984. “Literate Programming.” *Comput. J.* 27 (2): 97–111. <https://doi.org/10.1093/comjnl/27.2.97>.