

# 経験分布関数の理論と応用

塚原 英敦

成城大学経済学部

([tsukahar@seijo.ac.jp](mailto:tsukahar@seijo.ac.jp))

2012年9月@統数研 夏期大学院

# 目次

1. はじめに — 経験分布関数とその基本的性質
2. 経験過程の理論 — 古典的アプローチ
3. 統計学への応用 — 古典的な例中心
4. 経験過程の現代的アプローチ — 動機付け
5. おわりに — その他の結果

## 1. はじめに

$X_1, \dots, X_n : (\Omega, \mathcal{F}, P)$  上の i.i.d. 実確率変数列

$F(x) := P(X_1 \leq x)$  : 共通の分布関数

### 経験分布関数

$$\mathbb{F}_n(x, \omega) = \mathbb{F}_n(x) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \leq x\}}$$

- $x \in \mathbb{R}$  を固定すると

$$n\mathbb{F}_n(x) \sim \text{Bin}(n, F(x))$$

よって,  $E(\mathbb{F}_n(x)) = F(x)$  (不偏性)

- また, 大数の強法則により, 各  $x \in \mathbb{R}$  に対して

$$\mathbb{F}_n(x) \xrightarrow{\text{a.s.}} F(x)$$

が成り立つ (強一緻性)

- さらに , Lindeberg-Lévy の中心極限定理から ,  
各  $x \in \mathbb{R}$  に対して

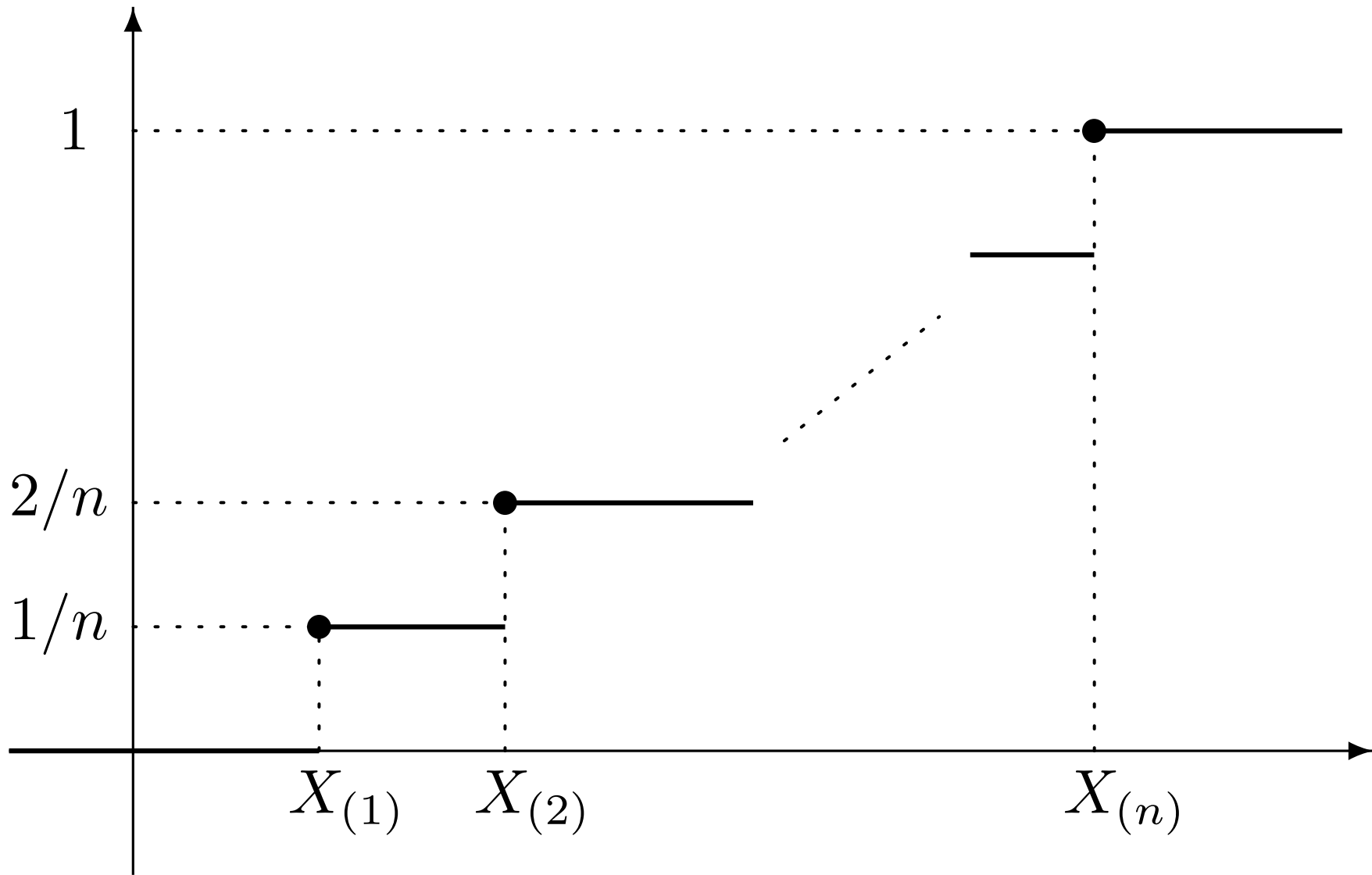
$$\sqrt{n} (\mathbb{F}_n(x) - F(x)) \xrightarrow{\mathcal{L}} N(0, F(x)(1 - F(x)))$$

を得る (漸近正規性)

- ▶▶  $\omega \in \Omega$  を固定して ,  $x$  についての関数とみると ,

$$X_{(1)} \leq X_{(2)} \leq \cdots \leq X_{(n)}$$

を順序統計量として ,



つまり, “ 経験分布関数  $\xleftrightarrow{1-1}$  順序統計量 ” である

$\mathcal{F}_{ac} :=$  絶対連続な分布関数全体

$F \in \mathcal{F}_{ac}$  ならば,

$$\mathbb{P} \left[ (X_1, \dots, X_n) = (x_{(i_1)}, \dots, x_{(i_n)}) \mid \right. \\ \left. (X_{(1)}, \dots, X_{(n)}) = (x_{(1)}, \dots, x_{(n)}) \right] = \frac{1}{n!}$$

$\implies (X_{(1)}, \dots, X_{(n)}) : \mathcal{F}_{ac}$  に対する **十分統計量**

さらに  $T = (X_{(1)}, \dots, X_{(n)})$  が  $\mathcal{F}_{ac}$  について**完備**であることを示すことができる。つまり、

$$\forall F \in \mathcal{F}_{ac}, E_F(g(T)) = 0$$

$$\Rightarrow g = 0, \text{ a.e. } \{\mathcal{L}_F(T) : F \in \mathcal{F}_{ac}\}$$

(例えば、鍋谷 (1978) , 定理 2.5.1, 2.5.2 参照)

▶▶ Lehmann-Scheffé の定理により、

$\mathbb{F}_n(x) : F(x)$  の**一様最小分散不偏 (UMVU) 推定量**



## ノンパラメトリック (NP) MLE

$$\text{NP 尤度} : L(F) = \prod_{i=1}^n [F(x_i) - F(x_{i-})]$$

- $z_1 < \cdots < z_m : x_1, \dots, x_n$  のうちの異なる値
- $n_j := \#\{x_i : x_i = z_j\}$
- $p_j := F(z_j) - F(z_{j-}), \quad \hat{p}_j := n_j/n$

ここで,  $\forall j, p_j > 0$ , かつ  $\exists j, p_j \neq \hat{p}_j$  とすると,

$$\begin{aligned} \log \frac{L(F)}{L(\mathbb{F}_n)} &= \log \prod_{i=1}^n \frac{p_j^{n_j}}{\hat{p}_j^{n_j}} = \sum_{j=1}^m n_j \log \frac{p_j}{\hat{p}_j} \\ &= n \sum_{j=1}^m \hat{p}_j \log \frac{p_j}{\hat{p}_j} < n \sum_{j=1}^m \hat{p}_j \left( \frac{p_j}{\hat{p}_j} - 1 \right) \leq 0 \end{aligned}$$

$\implies \mathbb{F}_n$  は  $F$  の NPMLE

## 2. 経験過程の理論（古典的アプローチ）

$F$  の分位関数 :  $F^{-1}(u) := \inf\{x : F(x) \geq u\}$

（1次元の場合の）基本定理

任意の分布関数  $F$  に対して，

$$1. F(x) \geq t \quad \Leftrightarrow \quad F^{-1}(t) \leq x$$

$$2. \xi \sim U(0, 1) \text{ のとき, } X := F^{-1}(\xi) \sim F$$

## 一様分布の場合への帰着

$\xi_1, \dots, \xi_n$  : i.i.d.  $U(0, 1)$  確率変数列

$$\mathbb{G}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\xi_i \leq t\}}$$

$$\mathbb{G}_n(F(x)) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\xi_i \leq F(x)\}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{F^{-1}(\xi_i) \leq x\}}$$

$$\implies \mathbb{F}_n \stackrel{\mathcal{L}}{=} \mathbb{G}_n(F)$$

## 確率積分変換

$X \sim F$  ならば,  $P(F(X) \leq t) \leq t, t \in [0, 1]$ .

等号成立は  $t \in \overline{\text{range}(F)}$  の場合に限る. 故に

$$F \text{ が連続} \implies F(X) \sim U(0, 1)$$

$F$  の  $\forall$  連続点  $x$  に対して,  $F_n(x) \rightarrow F(x)$



$F^{-1}$  の  $\forall$  連続点  $x$  に対して,  $F_n^{-1}(x) \rightarrow F^{-1}(x)$

## Glivenko-Cantelli の定理 (一様大数の法則)

$$\|\mathbb{G}_n - I\| := \sup_{0 < t < 1} |\mathbb{G}_n(t) - t| \xrightarrow{\text{a.s.}} 0, \quad n \rightarrow \infty$$

[証]  $0 < k \leq M$  に対して, SLLN より

$$\mathbb{G}_n(k/M) - k/M \xrightarrow{\text{a.s.}} 0, \quad n \rightarrow \infty$$

$\forall \varepsilon > 0, \exists M$  s.t.  $M^{-1} < \varepsilon$ .

よって,  $(k-1)/M \leq t \leq k/M$  ならば,

$$\mathbb{G}_n(t) - t \leq \mathbb{G}_n\left(\frac{k}{M}\right) - \frac{k-1}{M} = \mathbb{G}_n\left(\frac{k}{M}\right) - \frac{k}{M} + \frac{1}{M}$$

$$\mathbb{G}_n(t) - t \geq \mathbb{G}_n\left(\frac{k-1}{M}\right) - \frac{k}{M} = \mathbb{G}_n\left(\frac{k-1}{M}\right) - \frac{k-1}{M} + \frac{1}{M}$$

よって,

$$\sup_{0 < t < 1} |\mathbb{G}_n(t) - t| \leq \max_{0 \leq k \leq M} \left| \mathbb{G}_n\left(\frac{k}{M}\right) - \frac{k}{M} \right| + \frac{1}{M}$$

$$\xrightarrow{\text{a.s.}} 0 + \frac{1}{M} < \varepsilon \quad \blacksquare$$

▶▶  $\sqrt{n} (\mathbb{F}_n(x) - F(x)) \xrightarrow{\mathcal{L}} N(0, F(x)(1 - F(x)))$   
の一様化 (関数版) は?

**一様経験過程** :  $\mathbb{U}_n(t) := \sqrt{n}(\mathbb{G}_n(t) - t)$  を用いると

$$\begin{aligned}\sqrt{n} (\mathbb{F}_n(x) - F(x)) &\stackrel{\mathcal{L}}{=} \sqrt{n} (\mathbb{G}_n(F(x)) - I(F(x))) \\ &= \mathbb{U}_n(F(x))\end{aligned}$$

となり, 一様分布の場合に帰着される.



## 有限次元分布

- 確率過程  $X := \{X(t)\}_{t \in [0,1]}$  の有限次元分布とは ,  
 $(X(t_1), \dots, X(t_k))$  の  $\mathbb{R}^k$  上の分布のことである  
( $t_1 < \dots < t_k, k \in \mathbb{N}$ )
- $X \stackrel{\mathcal{L}}{=} Y$  :  $X$  と  $Y$  の全ての有限次元分布が一致
- $X_n \xrightarrow{\text{f.d.}} X$  とは ,  $\forall k \in \mathbb{N}, \forall t_1 < \dots < t_k,$   
$$(X_n(t_1), \dots, X_n(t_k)) \xrightarrow{\mathcal{L}} (X(t_1), \dots, X(t_k))$$

経験過程の場合，

$$\begin{pmatrix} \mathbb{U}_n(t_1) \\ \vdots \\ \mathbb{U}_n(t_k) \end{pmatrix} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} \mathbf{1}_{\{\xi_i \leq t_1\}} - t_1 \\ \vdots \\ \mathbf{1}_{\{\xi_i \leq t_k\}} - t_k \end{pmatrix}$$

$E(\mathbf{1}_{\{\xi_i \leq t_j\}}) = t_j$  と  $\text{Cov}(\mathbf{1}_{\{\xi_i \leq t_j\}}, \mathbf{1}_{\{\xi_i \leq t_l\}}) = t_j \wedge t_l - t_j t_l$  だから，多次元 CLT より，

$$\mathbb{U}_n \xrightarrow{\text{f.d.}} \mathbb{U}$$

ただし， $\mathbb{U}$  は **ブラウン橋** (Brownian bridge)

## ブラウン橋 (Brownian bridge)

ブラウン橋  $\mathbb{U} = \{\mathbb{U}(t)\}_{t \in [0,1]}$  は

$$E(\mathbb{U}(t)) = 0, \quad \text{Cov}(\mathbb{U}(s), \mathbb{U}(t)) = s \wedge t - st$$

を満たす, 標本路が連続なガウス過程である.

ブラウン運動  $S$  から,

$$\mathbb{U}(t) = S(t) - tS(1)$$

としても構成できる.

Doob(1949) :

“We shall assume, until a contradiction frustrates our devotion to heuristic reasoning, that *in calculating asymptotic  $\mathbb{U}_n(t)$  process distributions when  $n \rightarrow \infty$  we may simply replace the  $\mathbb{U}_n(t)$  processes by the  $\mathbb{U}(t)$  process*. It is clear that this cannot be done in all possible situations, but let the reader who has never used this sort of reasoning exhibit the first counter example.”

Kolmogorov-Smirnov 統計量  $\|U_n\|$  の漸近分布が  $\|U\|$  の分布と一致するが,  $U_n \xrightarrow{\text{f.d.}} U$  だけからは導かれない.



標本路の属する関数空間上の, 写像  $\omega \mapsto U_n(\cdot, \omega)$  により誘導される確率測度列の収束



弱収束理論 (weak convergence theory) が必要

各  $\omega$  に対して,  $\mathbb{U}_n(\cdot, \omega)$  は  $[0, 1]$  上の実数値関数:

$$\mathbb{U}_n(\cdot, \omega) \in \mathbb{R}^{[0,1]}$$

$\mathbb{R}^{[0,1]}$  は大きすぎる.  $\mathbb{U}_n(\cdot, \omega)$  は連続ではないが,

$$\mathbb{U}_n(\cdot, \omega) \in D[0, 1]$$

$D[0, 1]$ :  $[0, 1]$  上の右連続で左極限をもつ実関数全体

一方,  $\mathbb{U}(\cdot, \omega) \in C[0, 1] := [0, 1]$  上の実連続関数全体

## 距離空間上の確率測度の弱収束

$(S, d)$  : 距離空間 ,  $\mathcal{B}_d(S)$  : Borel 集合体

$P, P_n$  ( $n \in \mathbb{N}$ ) :  $(S, \mathcal{B}_d(S))$  上の確率測度

定義 :  $P_n$  が  $P$  に弱収束するとは ,

$$\int_S f(s) P_n(ds) \rightarrow \int_S f(s) P(ds), \quad f \in C_b(S)$$

が成り立つことをいい ,  $P_n \xrightarrow{w} P$  と書く .

$X$  と  $X_n$  :  $S$  値確率変数

i.e.,  $\forall B \in \mathcal{B}_d(S), X^{-1}(B) \in \mathcal{F}$  &  $X_n^{-1}(B) \in \mathcal{F}$

定義 :  $X_n$  が  $X$  に**法則収束 (分布収束)** するとは,  $X_n$  の分布  $P_{X_n}$  が  $X$  の分布  $P_X$  に弱収束することをいう. このとき,  $X_n \xrightarrow{\mathcal{L}} X$  と書く.

注 :  $P_X(B) = \mathbb{P}(X \in B), B \in \mathcal{B}_d(S)$

$P_{X_n}(B) = \mathbb{P}(X_n \in B), B \in \mathcal{B}_d(S)$



## Portmanteau 定理

次の 5 つの条件は同値である :

$$(i) P_n \xrightarrow{w} P$$

$$(ii) \forall f \in C_{ub}(S) \text{ に対して, } \int f dP_n \rightarrow \int f dP$$

$$(iii) \forall \text{閉集合 } F \text{ に対して, } \limsup_n P_n(F) \leq P(F)$$

$$(iv) \forall \text{開集合 } G \text{ に対して, } \liminf_n P_n(G) \geq P(G)$$

$$(v) \forall P\text{-連続集合 } A \text{ (i.e., } P(\partial A) = 0 \text{)} \text{ に対して,}$$

$$P_n(A) \rightarrow P(A)$$

## 連続写像定理

$(S, d), (S', d')$  : 距離空間

$h: S \rightarrow S'$  : Borel 可測

$D_h$  :  $h$  の不連続点全体から成る集合

$P_n, P$  :  $(S, d)$  上の確率測度

定理

$$P_n \xrightarrow{w} P, P(D_h) = 0 \quad \Rightarrow \quad P_n h^{-1} \xrightarrow{w} P h^{-1}$$

$P_n$  と  $P$  がそれぞれ  $S$  値確率変数  $X_n$  と  $X$  の分布であると考えると, 上の定理は

$$X_n \xrightarrow{\mathcal{L}} X \text{ かつ } P(X \in D_h) = 0 \\ \implies h(X_n) \xrightarrow{\mathcal{L}} h(X)$$

となることを意味する .

▶▶ この定理における関数  $h$  が  $n$  に依存する場合への拡張は応用上有用 (Topsøe)

## 弱収束の十分条件

- $(S, d)$  上の確率測度の集合  $\{P_\alpha\}$  が**一様に緊密 (uniformly tight)** であるとは,

$$\forall \varepsilon > 0, \exists K_\varepsilon \text{ コンパクト s.t. } \inf_{\alpha} P_\alpha(K_\varepsilon) \geq 1 - \varepsilon$$

- $(S, d)$  上の確率測度の集合  $\{P_\alpha\}$  が**相対コンパクト (relatively compact)** であるとは,  $\{P_\alpha\}$  に含まれる任意の列  $(P_n)$  が弱収束する部分列をもつことをいう.

$$\left. \begin{array}{l} \{P_n\} \text{ が相対コンパクト} \\ \text{極限 } P \text{ を一意に識別する条件} \end{array} \right\} \implies P_n \xrightarrow{w} P$$

▶▶  $C[0, 1]$  や  $D[0, 1]$  の場合，極限  $P$  を一意に識別する条件として有限次元分布の収束が有効

相対コンパクト性は直接示し難いため，チェックしやすい（必要）十分条件が欲しい



## Prohorov の定理

- (i)  $\{P_\alpha\}$  が一様に緊密ならば,  $\{P_\alpha\}$  は相対コンパクトである.
- (ii)  $(S, d)$  が完備かつ可分のとき,  $\{P_\alpha\}$  が相対コンパクトならば,  $\{P_\alpha\}$  は一様に緊密である.

★  $\left. \begin{array}{l} \{P_n\} \text{ が一様に緊密} \\ \text{極限 } P \text{ を一意に識別する条件} \end{array} \right\} \implies P_n \xrightarrow{w} P$

## Skorohod-Dudley-Wichura の表現定理

$(P_n) \xrightarrow{w} P$  ならば, 適当な確率空間上に確率変数  $Y, Y_1, Y_2, \dots$  を構成して,  $Y_n \sim P_n, Y \sim P, Y_n \xrightarrow{\text{a.s.}} Y$  が成り立つようにできる.

- $(S, d)$  が完備かつ可分するとき : Skorohod (1956)
- $(S, d)$  が可分するとき : Dudley (1968)
- $(S, d)$  は一般の距離空間で,  $P$  の台が可分なとき :  
Wichura (1970)

## 距離空間 $C$ と $D$

$$\text{sup 距離} : \|x - y\| := \sup_{0 \leq t \leq 1} |x(t) - y(t)|$$

▶▶  $(C, \|\cdot\|)$  は完備かつ可分

3 つの  $\sigma$  集合体

$$\mathcal{C} := \sigma(\text{有限次元集合}), \quad \mathcal{C}_{\|\cdot\|} := \sigma(\|\cdot\| \text{-開集合})$$

$$\mathcal{C}_{\|\cdot\|}^B := \sigma(\|\cdot\| \text{-開球})$$

は一致する .



▶▶  $(D, \|\cdot\|)$  は可分でない

$\mathcal{D} := \sigma(\text{有限次元集合}), \quad \mathcal{D}_{\|\cdot\|} := \sigma(\|\cdot\| \text{-開集合})$

$\mathcal{D}_{\|\cdot\|}^B := \sigma(\|\cdot\| \text{-開球})$

について,  $\mathcal{D} = \mathcal{D}_{\|\cdot\|}^B \subsetneq \mathcal{D}_{\|\cdot\|}$

$[\phi_t = \mathbf{1}_{[t,1]}$  で定義される  $\phi: [0, 1] \rightarrow D$  は可測 ( $\mathcal{D}$ )

$A = \bigcup_{t \in H} B_{\|\cdot\|}(\phi_t, 1/2)$  は  $\|\cdot\|$  開集合,  $\phi^{-1}A = H$

$H$  を非可測集合ととれば  $\phi_t$  は可測 ( $\mathcal{D}_{\|\cdot\|}$ ) でない]

$\omega \mapsto \mathbb{U}_n(\cdot, \omega)$  は可測 ( $\mathcal{F} / \mathcal{D}_{\|\cdot\|}$ ) でない可能性

$\implies \mathbb{U}_n$  の分布  $P\mathbb{U}_n^{-1}$  が定義できない.

★ 対処法 :

1. Skorohod の  $J_1$  位相  $\Rightarrow D$  は完備可分
2. 距離は  $\|\cdot\|$  のまま,  $(D, \mathcal{D}_{\|\cdot\|}^B)$  を考える  
(Dudley, Wichura ‘ $\overset{w}{\rightarrow}$ ’ 理論)
3. 可測性の要請を捨てる  
(Hoffmann-Jørgensen & Dudley 理論)

## 文献案内

- 2 を採用して，1 次元経験過程を主に扱った大書：  
Shorack & Wellner('86)
- 1 はセミマルチンゲールや確率積分の弱収束：  
Ethier & Kurtz('86), Jacod & Shiryaev('87)
- 3 が経験過程を扱う現代の流儀：  
Dudley('99), van der Vaart & Wellner('96)

## $C$ における一様緊密性

連続係数 :  $w_x(\delta) = \sup_{|s-t| \leq \delta} |x(s) - x(t)|$

$(C, \mathcal{C})$  上の確率測度列  $(P_n)$  が一様に緊密  $\Leftrightarrow$

(i)  $\forall \eta > 0, \exists a \text{ \& } n_0 \text{ s.t.}$

$$P_n(x : |x(0)| \geq a) \leq \eta, \quad n \geq n_0$$

(ii)  $\forall \varepsilon, \eta > 0, \exists \delta \text{ \& } n_0 \text{ s.t.}$

$$P_n(x : w_x(\delta) \geq \varepsilon) \leq \eta, \quad n \geq n_0$$

## $D$ における一様緊密性

極限の台が  $D$  である場合，その条件はやや複雑  
( Billingsley('99) 参照 )

### 経験過程の場合に便利な条件

$(D, \mathcal{D})$  上の確率測度列  $(P_n)$  に対して， $C$  での一様緊密性の条件 (i), (ii) が成り立っていれば， $(P_n)$  は一様緊密であり，収束部分列の極限  $P$  は  $P(C) = 1$  を満たす．

この条件を用いることによって,

### Donsker の定理

$(D, \mathcal{D}, \|\cdot\|)$  において,

$$U_n \xrightarrow{\mathcal{L}} U, \quad n \rightarrow \infty$$

統計への応用にはこのままでは不便

以下, S&W ('86), Shorack(2000) に従う

## いくつかの経験過程

$\xi_{n1}, \dots, \xi_{nn} : \text{i.i.d. } U(0, 1)$  確率変数の 3 角列

$$\mathbb{G}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\xi_{ni} \leq t\}}$$

$$\mathbb{U}_n(t) := \sqrt{n}(\mathbb{G}_n(t) - t) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{1}_{\{\xi_{ni} \leq t\}} - t]$$

前に計算したように,

$$\mathbb{E}(\mathbb{U}_n(t)) = 0, \quad \text{Cov}(\mathbb{U}_n(s), \mathbb{U}_n(t)) = s \wedge t - st$$

經驗分位関数 (empirical quantile) :

$$\mathbb{G}_n^{-1}(t) := \inf\{u \in [0, 1] : \mathbb{G}_n(u) \geq t\}$$

經驗分位過程 (empirical quantile process) :

$$\mathbb{V}_n(t) := \sqrt{n}(\mathbb{G}_n^{-1}(t) - t)$$

恒等式

$$\mathbb{U}_n(t) = -\mathbb{V}_n(\mathbb{G}_n) + \sqrt{n}(\mathbb{G}_n^{-1} \circ \mathbb{G}_n - I)$$

$$\mathbb{V}_n(t) = -\mathbb{U}_n(\mathbb{G}_n^{-1}) + \sqrt{n}(\mathbb{G}_n \circ \mathbb{G}_n^{-1} - I)$$



図を描くとわかること：

$$\|G_n \circ G_n^{-1} - I\| = \frac{1}{n}$$

$$\|G_n^{-1} \circ G_n - I\| = \max_{1 \leq i \leq n+1} (\xi_{n:i} - \xi_{n:i-1})$$

さらに，

$$\|G_n - I\| = \|G_n^{-1} - I\|$$

Glivenko-Cantelli 定理より

$$\|G_n - I\| \xrightarrow{\text{a.s.}} 0, \quad \|G_n^{-1} - I\| \xrightarrow{\text{a.s.}} 0$$

## 加重経験過程

$c_{n1}, \dots, c_{nn}$  : (基準化された) 定数の列

$$\bar{c}_n := \frac{1}{n} \sum_{i=1}^n c_{ni} = 0, \quad \sigma_{c,n}^2 := \frac{1}{n} \sum_{i=1}^n (c_{ni} - \bar{c}_n)^2 = 1$$

- **UAN 条件** :  $\max_{1 \leq i \leq n} \frac{|c_{ni}|}{\sqrt{n}} \rightarrow 0, \quad n \rightarrow \infty$

を通常仮定する .

## 加重経験過程 (weighted empirical process)

$$\mathbb{W}_n(t) := \frac{1}{\sqrt{n}} \sum_{i=1}^n c_{ni} [\mathbf{1}_{\{\xi_{ni} \leq t\}} - t]$$

$$\mathbb{E}(\mathbb{W}_n(t)) = 0$$

$$\text{Cov}(\mathbb{W}_n(s), \mathbb{W}_n(t)) = \sigma_{c,n}^2 (s \wedge t - st) = s \wedge t - st$$

$$\text{Cov}(\mathbb{U}_n(s), \mathbb{W}_n(t)) = \bar{c}_n (s \wedge t - st) = 0$$

$\Rightarrow \mathbb{W}_n \xrightarrow{\text{f.d.}} \mathbb{W}$ ,  $\mathbb{W}$  は  $\mathbb{U}$  と独立なブラウン橋

$(R_{n1}, \dots, R_{nn}) : \xi_{n1}, \dots, \xi_{nn}$  の順位

$(D_{n1}, \dots, D_{nn}) : \text{反順位}$

$$\xi_{nD_{ni}} = \xi_{n:i}, \quad \xi_{ni} = \xi_{n:R_{ni}}$$

有限抽出過程 (empirical finite sampling process) :

$$\mathbb{R}_n(t) := \frac{1}{\sqrt{n}} \sum_{i=1}^{[(n+1)t]} c_{nD_{ni}}$$

## 恒等式

$$W_n = R_n(\tilde{G}_n), \quad R_n = W_n(\tilde{G}_n^{-1})$$

$[\tilde{G}_n, \tilde{G}_n^{-1}]$  は  $G_n, G_n^{-1}$  の線形化バージョン]

順序統計量  $\xi_{n:1}, \dots, \xi_{n:n}$  と順位  $(R_{n1}, \dots, R_{nn})$  は  
独立



$R_n$  と  $V_n$  は独立

## 単純線形順位統計量

$$\begin{aligned} T_n &:= \frac{1}{\sqrt{n}} \sum_{i=1}^n c_{ni} K \left( \frac{R_{ni}}{n+1} \right) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n K \left( \frac{i}{n+1} \right) c_{ni} D_{ni} \\ &= \int_0^1 K \, d\mathbb{R}_n = - \int_0^1 \mathbb{R}_n \, dK \end{aligned}$$

[最後の = は  $K$  が有界変動かつ左連続なら OK]

## 特別構成 (special construction) 定理 :

単一の確率空間  $(\Omega, \mathcal{F}, P)$  上に、行独立な  $U(0, 1)$  3角列  $\xi_{n1}, \dots, \xi_{nn}$  と独立な2つのブラウン橋  $U = -V$  と  $W$  が次の条件を満たすように構成できる :

$$\|U_n - U\| \xrightarrow{\text{a.s.}} 0, \quad \|V_n - V\| \xrightarrow{\text{a.s.}} 0$$

$$\|W_n - W\| \xrightarrow{\text{a.s.}} 0, \quad \|R_n - R\| \xrightarrow{\text{a.s.}} 0$$

ただし  $\{c_{ni}\}$  は UAN 条件を満たし  $\bar{c}_n = 0, \sigma_{c,n}^2 = 1$

## Pyke-Shorack の定理 :

$q : (0, 1)$  上の正の関数 ,  $(0, \frac{1}{2}]$  で増加 ,  $[\frac{1}{2}, 1)$  で減少

かつ  $\int_0^1 [q(t)]^{-2} dt < \infty$  [例:  $[t(1-t)]^{1/2-\delta}$ ]

このとき , 上の特別構成に対して ,

$$\|(\mathbb{U}_n - \mathbb{U})/q\| \xrightarrow{\mathbb{P}} 0, \quad \|(\mathbb{V}_n - \mathbb{V})/q\| \xrightarrow{\mathbb{P}} 0$$

$$\|(\mathbb{W}_n - \mathbb{W})/q\| \xrightarrow{\mathbb{P}} 0, \quad \|(\mathbb{R}_n - \mathbb{R})/q\| \xrightarrow{\mathbb{P}} 0$$



## 適合度検定

特別構成の  $\xi_{n1}, \dots, \xi_{nn}$  を用いて,  $X_{ni} := F^{-1}(\xi_{ni})$

$\Rightarrow X_{n1}, \dots, X_{nn}$  i.i.d.  $F$

$$\mathbb{F}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_{ni} \leq x\}} \quad : \text{経験分布関数}$$

$\mathbb{E}_n(x) := \sqrt{n}(\mathbb{F}_n(x) - F(x)) = \mathbb{U}_n(F(x))$  とおくと

$$\|\mathbb{E}_n(x) - \mathbb{U}(F)\| \leq \|\mathbb{U}_n - \mathbb{U}\| \xrightarrow{\text{a.s.}} 0$$

## Kolmogorov-Smirnov 統計量 ( $F$ : 連続)

$$\text{片側} : \sqrt{n} \sup_x (\mathbb{F}_n(x) - F(x))^+ = \|\mathbb{U}_n^+\| \xrightarrow{\text{a.s.}} \|\mathbb{U}^+\|$$

$$\text{両側} : \sqrt{n} \sup_x |\mathbb{F}_n(x) - F(x)| = \|\mathbb{U}_n\| \xrightarrow{\text{a.s.}} \|\mathbb{U}\|$$

▶▶  $\mathbb{R}_+$  上の BM  $S$  について,  $c \geq 0, d > 0$

$$(i) \ P(\exists t \geq 0, S(t) \geq ct + d) = \exp(-2cd)$$

$$(ii) \ P(\exists t \geq 0, |S(t)| \geq ct + d) \\ = 2 \sum_{k=1}^{\infty} (-1)^{k+1} \exp(-2k^2 cd)$$

$$(i) \mathbf{P}(\|\mathbb{U}^+\| > b) = \exp(-2b^2)$$

$$(ii) \mathbf{P}(\|\mathbb{U}\| > b) = 2 \sum_{k=1}^{\infty} (-1)^{k+1} \exp(-2k^2b^2)$$

$$\because \mathbf{P}(\|\mathbb{U}^+\| > b) = \mathbf{P}(\exists t \in (0, 1), \mathbb{U}(t) > b)$$

$$= \mathbf{P}(\exists t \in (0, 1), (1-t)\mathbb{S}\left(\frac{t}{1-t}\right) > b)$$

$$= \mathbf{P}(\exists r > 0, \frac{1}{1+r}\mathbb{S}(r) > b)$$

$$= \mathbf{P}(\exists r > 0, \mathbb{S}(r) > b + rb) = \exp(-2b^2)$$

## Cramer-von Mises 統計量 ( $F$ : 連続)

$$\begin{aligned} & n \int [\mathbb{F}_n(x) - F(x)]^2 dF(x) \\ &= \int_0^1 \mathbb{U}_n^2(t) dt \xrightarrow{\text{a.s.}} \int_0^1 \mathbb{U}^2(t) dt \end{aligned}$$

$\mathbb{U}$  の Karhunen-Loève 展開 :

$$Y(t) := \sum_{k=1}^{\infty} \phi_k(t) \frac{1}{\pi k} Z_k$$

$Z_1, Z_2, \dots : \text{i.i.d. } N(0, 1)$

$\phi_k(t) = \sqrt{2} \sin(\pi kt), k = 1, 2, \dots : L^2$  の正規直交系

Parseval の恒等式から

$$\int_0^1 Y^2(t) dt := \sum_{k=1}^{\infty} \frac{1}{\pi^2 k^2} Z_k^2$$

$U \stackrel{\mathcal{L}}{=} Y$  を示せば  $\int_0^1 U^2(t) dt$  の分布が近似計算できる。

## Anderson-Darling 統計量 ( $F$ : 連続)

$$\begin{aligned} A_n &:= n \int \frac{[\mathbb{F}_n(x) - F(x)]^2}{F(x)(1 - F(x))} dF(x) \\ &= \int_0^1 \frac{U_n^2(t)}{t(1-t)} dt \xrightarrow{\text{a.s.}} \int_0^1 \frac{U^2(t)}{t(1-t)} dt \end{aligned}$$

$$\int_0^1 \frac{U^2(t)}{t(1-t)} dt \stackrel{\mathcal{L}}{=} \sum_{k=1}^{\infty} \frac{1}{k(k+1)^2} Z_k^2$$

が知られている。

## $L$ 統計量の漸近正規性

$X_{n1}, \dots, X_{nn}$  i.i.d.  $F$  (特別構成から)

$X_{n:1} < \dots < X_{n:n}$  : 順序統計量

$$T_n := \frac{1}{n} \sum_{i=1}^n c_{ni} h(X_{n:i})$$

の形の統計量を  $L$  統計量 という .

(linear combination of functions of order statistics)

$$T_n = \int_0^1 h(\mathbb{F}_n^{-1}(u)) J_n(u) \, du = \int_{[0,1]} h(\mathbb{F}_n^{-1}(u)) \, d\Psi_n(u)$$

$$J_n(u) := \sum_{i=1}^n c_{ni} \mathbf{1}_{(\frac{i-1}{n}, \frac{i}{n}]}(u) + c_{n1} \mathbf{1}_{\{0\}}(u),$$

$$\Psi_n(u) := \int_{1/2}^u J_n(v) \, dv$$

中心化定数 :  $g := h \circ F^{-1}$  とおいて ,

$$\mu_n := \int_0^1 g(u) J_n(u) \, du = \int_{[0,1]} g(u) \, d\Psi_n(u)$$



仮定：

- $B(u) := Mu^{-b_1}(1-u)^{-b_2}$ ,  $|J_n| \leq B$ ,  $|J| \leq B$
- $h = h_1 - h_2$ ,  $h_1$  &  $h_2$  は左連続増加関数で,

$$|h_i(F^{-1})| \leq H, \quad H(u) := Mu^{-d_1}(1-u)^{-d_2}$$

- $|g|$ -ほとんどすべての  $u$  に対して,  $J$  は  $u$  で連続,  
かつ  $J_n \rightarrow J$  (局所一様)

$b_i + d_i < \frac{1}{2}$ ,  $i = 1, 2$  を仮定する

$$\begin{aligned}
T_n - \mu_n &= \int_{[0,1]} g \circ \mathbb{G}_n^{-1} d\Psi_n - \int_{[0,1]} g d\Psi \\
&= \int_{[0,1]} g d[\Psi_n \circ \mathbb{G}_n - \Psi_n] \\
&= - \int_{[0,1]} [\Psi_n \circ \mathbb{G}_n - \Psi_n] dg \\
&= - \int_{[0,1]} [\Psi_n \circ \mathbb{G}_n - \Psi_n - (\mathbb{G}_n - I)J] dg \\
&\quad - \int_{[0,1]} [\mathbb{G}_n - I] J dg \\
&=: -\gamma_n - S_n
\end{aligned}$$

$$\sqrt{n}\gamma_n = \int_0^1 \mathbb{U}_n(u) A_n(u) dg(u)$$

$$\sqrt{n}|\gamma_n| \leq \|\mathbb{U}_n/q\| \int_0^1 |A_n(u)|q(u) d|g|(u)$$

ただし ,

$$A_n(u) := \frac{1}{\mathbb{G}_n(u) - u} \int_u^{\mathbb{G}_n(u)} J_n(v) dv - J(u)$$

Glivenko-Cantelli 定理より ,  $A_n \rightarrow 0$ ,  $|g|$ -a.e.

仮定より ,  $|A_n(u)| \leq B(\mathbb{G}_n(u)) \vee B(u) + B(u)$

▶▶  $\mathbb{G}_n$  と  $\mathbb{G}_n^{-1}$  に対する確率的線形限界

$\forall \varepsilon > 0, \exists \lambda \ \& \ A_{n,\varepsilon}, P(A_{n,\varepsilon}) \geq 1 - \varepsilon$  s.t.  $A_{n,\varepsilon}$  上で

$$\mathbb{G}_n(t) \leq t/\lambda, t \in [0, 1], \quad \mathbb{G}_n(t) \geq \lambda t, t \in [\xi_{n:1}, 1]$$

$$\mathbb{G}_n(1-t) \leq 1 - \lambda(1-t), t \in [0, \xi_{n:n}),$$

$$\mathbb{G}_n(1-t) \geq 1 - (1-t)/\lambda, t \in [0, 1]$$

を用いると ,  $\sqrt{n}|\gamma_n| \xrightarrow{P} 0$  と言える .

$$\sqrt{n}S_n = \int_0^1 \mathbb{U}_n J \, dg = \int_0^1 \frac{\mathbb{U}_n - \mathbb{U}}{q} J q \, dg + \int_0^1 \mathbb{U} J \, dg$$

$$\left| \int_0^1 \frac{\mathbb{U}_n - \mathbb{U}}{q} J q \, dg \right| \leq \left\| \frac{\mathbb{U}_n - \mathbb{U}}{q} \right\| \int_0^1 B q \, dg \xrightarrow{\text{P}} 0$$

よって,

$$\sqrt{n}(T_n - \mu_n) \xrightarrow{\mathcal{L}} - \int_{[0,1]} \mathbb{U}(u) J(u) \, dg(u) \sim N(0, \sigma^2)$$

$$\sigma^2 := \int_0^1 \int_0^1 (u \wedge v - uv) J(u) J(v) \, dg(u) dg(v)$$

## 2 標本問題

$$X_1, \dots, X_m \quad \text{i.i.d. } F_1$$

$$Y_1, \dots, Y_n \quad \text{i.i.d. } F_2 \quad N := m + n$$

帰無仮説  $H_0 : F_1 = F_2$  の検定

特別構成：

$$X_{m1} = F_1^{-1}(\xi_{m1}^{(1)}), \dots, X_{mm} = F_1^{-1}(\xi_{mm}^{(1)})$$

$$Y_{n1} = F_2^{-1}(\xi_{n1}^{(2)}), \dots, Y_{nn} = F_2^{-1}(\xi_{nn}^{(2)})$$

$\mathbb{F}_{1m} : X_{m1}, \dots, X_{mm}$  の経験分布関数

$\mathbb{F}_{2m} : Y_{n1}, \dots, Y_{nn}$  の経験分布関数

$H_0$  の下での

$$\sqrt{\frac{mn}{N}} (\mathbb{F}_{1m} - \mathbb{F}_{2n})$$

分布を求めるには, 共通の df を  $F$ ,  $\lambda_N := \frac{m}{N}$  として,

$$\sqrt{\frac{mn}{N}} (\mathbb{F}_{1m} - \mathbb{F}_{2n}) = \sqrt{1 - \lambda_N} U_{1m}(F) - \sqrt{\lambda_N} U_{2n}(F)$$

## 2 標本 Kolmogorov-Smirnov 検定

$F$  が連続とすると,  $H_0$  の下で ( $\lambda_N \rightarrow \lambda > 0$ )

$$\begin{aligned} \sqrt{\frac{mn}{N}} \|\mathbb{F}_{1m} - \mathbb{F}_{2n}\| &= \left\| \sqrt{1 - \lambda_N} U_{1m} - \sqrt{\lambda_N} U_{2n} \right\| \\ &\xrightarrow{\text{a.s.}} \left\| \sqrt{1 - \lambda} U_1 - \sqrt{\lambda} U_2 \right\| \end{aligned}$$

$U_1, U_2$  は独立なブラウン橋

$\implies \sqrt{1 - \lambda} U_1 - \sqrt{\lambda} U_2$  もブラウン橋



## Chernoff-Savage 問題

$$H_N := \lambda_N F_1 + (1 - \lambda_N) F_2$$

$$\mathbb{H}_N := \lambda_N \mathbb{F}_{1m} + (1 - \lambda_N) \mathbb{F}_{2n}$$

: プールした標本の経験分布関数

- 検定統計量 :  $T_N := \frac{1}{m} \sum_{i=1}^N c_{Ni} Z_{Ni}$ ,  $c_{Ni}$  : 定数 ,

$$Z_{Ni} = \begin{cases} 1 & \text{プールした標本中 } i \text{ 番目に大きいもの } \in \{X_{mi}\} \\ 0 & \text{その他} \end{cases}$$

スコア関数  $J_N$  を

$$J_N(t) := c_{Ni}, \quad \frac{i-1}{N} < t \leq \frac{i}{N}, \quad i = 1, \dots, N$$

で定義すると

$$T_N = \int J_N(\mathbb{H}_N) dF_{1m}$$

中心化定数 :  $\mu_N := \int J(H_N) dF_1$

## 仮定

1.  $\exists \lambda_0 \in (0, \frac{1}{2}), \lambda_0 \leq \lambda_N \leq 1 - \lambda_0$
2.  $\exists J, \frac{1}{\sqrt{m}} \sum_{i=1}^{N-1} |c_{Ni} - J(i/N)| \rightarrow 0 \quad (N \rightarrow \infty)$
3.  $N^{-1/2} c_{NN} \rightarrow 0 \quad (N \rightarrow \infty)$
4.  $J$  は連続な導関数  $J'$  をもち,  $\exists \delta > 0,$   
 $|J| \leq [I(1-I)]^{-1/2+\delta}, \quad |J'| \leq [I(1-I)]^{-3/2+\delta},$

$$\begin{aligned}
& \sqrt{m}(T_N - \mu_N) \\
&= \int \sqrt{m}[J(\mathbb{H}_N) - J(H_N)] d\mathbb{F}_{1m} + \int J(H_N) d\sqrt{m}(\mathbb{F}_{1m} - F_1) \\
&= \int \frac{J_N(\mathbb{H}_N) - J(H_N)}{\mathbb{H}_N - H_N} \sqrt{m}(\mathbb{H}_N - H_N) d\mathbb{F}_{1m} \\
&\qquad\qquad\qquad + \int J(H_N) d\mathbb{U}_{1m} \\
&\stackrel{a}{=} \int J'(H_N)[\lambda_N \mathbb{U}_{1m}(F_1) + \sqrt{\lambda(1-\lambda)} \mathbb{U}_{2n}(F_2)] d\mathbb{F}_{1m} \\
&\qquad\qquad\qquad + \int J(H_N) d\mathbb{U}_{1m}
\end{aligned}$$

$$\begin{aligned}
& \sqrt{m}(T_N - \mu_N) \\
& \stackrel{a}{=} \lambda_N \int J'(H_N) \mathbb{U}_{1m}(F_1) dF_1 \\
& \quad + \sqrt{\lambda(1-\lambda)} \int J'(H_N) \mathbb{U}_{2n}(F_2) dF_1 \\
& \quad - \int \mathbb{U}_{1m}(F_1) J'(H_N) d[\lambda_N F_1 + (1-\lambda_N) F_2] \\
& = \sqrt{1-\lambda_N} \left[ \sqrt{\lambda_N} \int J'(H_N) \mathbb{U}_{2n}(F_2) dF_1 \right. \\
& \quad \left. - \sqrt{1-\lambda_N} \int J'(H_N) \mathbb{U}_{1m}(F_1) dF_2 \right]
\end{aligned}$$

$$\begin{aligned} \therefore \sqrt{m}(T_N - \mu_N) \\ \stackrel{a}{=} \sqrt{1 - \lambda_N} \left[ \sqrt{\lambda_N} \int J'(H_N) \mathbb{U}_2(F_2) dF_1 \right. \\ \left. - \sqrt{1 - \lambda_N} \int J'(H_N) \mathbb{U}_1(F_1) dF_2 \right] \end{aligned}$$

この確率変数は平均 0 の正規分布に従う．分散も  $\mathbb{U}_1$  と  $\mathbb{U}_2$  の独立性から容易に計算できる．

**例** :  $c_{Ni} = i/N$  (Wilcoxon)

$c_{Ni} = \Phi^{-1}(i/(N + 1))$  (van der Waerden)

漸近的最適性 1 : たたみ込み定理 (Beran(1977))

$F$  の推定量の列  $(\hat{F}_n)$  が**正則** (regular) であり, その  
極限過程が  $C$  上の  $\mathbb{Z}$  であるとき,

$$\mathbb{Z} \stackrel{\mathcal{L}}{=} U(F) + W$$

が成り立つ. ここで, ブラウン橋  $U$  と  $W$  は独立である.

## 漸近的最適性 2 : 漸近ミニマックス性

(Dvoretzky et al.(1956), Millar(1979))

$w$  を適当なクラスの損失関数 ,  $D_n$  を確率化決定関数 ,  $(\hat{F}_n)$  を  $F$  の任意の推定量列として ,

$$\lim_{n \rightarrow \infty} \frac{\sup_F \mathbf{E}_F [w(\|\sqrt{n}(\mathbb{F}_n - F)\|)]}{\sup_{b \in D_n} \sup_F \mathbf{E}_F \left[ \int w(\|\sqrt{n}(\hat{F}_n - F)\|) b(d\hat{F}_n, \mathbf{X}) \right]} = 1$$



DKW 最大不等式 : Dvoretzky, Kiefer and Wolfowitz

$$P(\sup_x |\mathbb{F}_n(x) - F(x)| > z) \leq C e^{-2nz^2}, \quad z > 0$$

ここで  $C > 0$  は  $F$  に依存しない定数である ( $C = 2$  が最良).

この不等式の多次元への拡張は Kiefer によってなされた (Kiefer(1961) 参照). この不等式は非常に強力であり, これを用いて例えば Glivenko-Cantelli の定理を強めた結果を証明できる.

## 重複対数の法則 (laws of the iterated logarithm)

- Smirnov :

$$\limsup_{n \rightarrow \infty} \frac{\|\mathbb{U}_n\|}{\sqrt{2 \log \log n}} = \frac{1}{2}$$

- Chung :  $\lambda_n \nearrow$  に対して

$$P(\|\mathbb{U}_n\| \leq \lambda_n, \text{ i.o.}) = \begin{cases} 0 & \sum_{n=1}^{\infty} (\lambda_n^2/n) \exp(-2\lambda_n^2) < \infty \\ 1 & \sum_{n=1}^{\infty} (\lambda_n^2/n) \exp(-2\lambda_n^2) = \infty \end{cases}$$

## Strassen 型の関数重複対数の法則

$D$  に標本路をもつ確率過程

$$\frac{U_n}{\sqrt{2 \log \log n}}$$

は  $D$  上  $\| \cdot \|$  に関して P-a.s. 相対コンパクトであり ,  
その集積点全体は

$$\mathcal{H} = \left\{ h : \text{絶対連続} , h(0) = h(1) = 0 , \int_0^1 [h'(t)]^2 dt \leq 1 \right\}$$

( Finkelstein , James )

Hungarian 構成 (強不変原理) の一例 :

$U(0, 1)$  確率変数列  $(\xi_n)$  とブラウン橋の列  $(\mathbb{B}_n)$  をあ  
る (共通の) 確率空間上に

$$\limsup_{n \rightarrow \infty} \frac{\sqrt{n}}{(\log n)^2} \|\mathbb{U}_n - \mathbb{B}_n\| < \infty, \quad \text{a.s.}$$

を満たすように構成できる . ここで  $\mathbb{U}_n$  は  $\xi_1, \dots, \xi_n$   
に基づく経験過程である .

(Csörgő and Révész, Komlós, Major and Tusnády)

## Bahadur-Kiefer 表現

$F(x)$  が  $x = F^{-1}(u)$  において 2 回連続微分可能で  $f(F^{-1}(u)) > 0$  ( $f = F'$ ) を満たすとき ,

$$\mathbb{F}_n^{-1}(u) = F^{-1}(u) + \frac{u - \mathbb{F}_n(F^{-1}(u))}{f(F^{-1}(u))} + R_n(u),$$

ただし  $R_n(u) = O(n^{-3/4}(\log n)^{3/4})$  , a.s. ( $n \rightarrow \infty$ )  
が成り立つ .

Kiefer のより精緻な結果 :

$$\limsup_{n \rightarrow \infty} \pm \frac{n^{3/4} R_n(u)}{(\log \log n)^{3/4}} = \frac{2^{5/4} [u(1-u)]^{1/4}}{3^{3/4}}, \quad \text{a.s.}$$

$R_n^* \triangleq \sup_{0 < u < 1} f(F^{-1}(u)) |R_n(u)|$  に対して ,

$$\limsup_{n \rightarrow \infty} \frac{n^{3/4} R_n^*}{(\log n)^{1/2} (\log \log n)^{1/4}} = 2^{-1/4}, \quad \text{a.s.}$$

## 最後に (much more to learn)

- $U$  統計量と  $U$  過程 (de la Peña & Giné)
- パターン認識・分類
- ブートストラップ法
- 関数デルタ法
- 生存解析

● ● ●