

経験的ベイズ手法

(複合型決定問題を含む)

統計教理研 鈴木 義一郎

§1. 序

この報告の目的は、表題に掲げた主題に關する H. Robbins の一連の論文 [1], [2], [3], [4] (主に [1], [3]) の簡単なスケッチと、これらを綜括的に議論した筆者 [5] の要旨の紹介である。

論文 [1] の重要な意義は、通例妥当と考えられてきたミニマックスな法則に対して有力な競争相手の決定法則を構成し、従来数学的、理論的面に於てのみ重視されてきた“ベイズ的考え”を統計的、實際的面に取り入れ、いわゆる“経験的ベイズ手法 (empirical Bayes approach)”と呼ばれる一連の理論の展開のきっかけを与えたことに在ろう。経験的ベイズという言葉は [1], [2] の論文に於ては用いられていないが、彼等の思考の根底にそのような概念が潜んでいたことは、Robbins が続いて発表した論文 [3], [4] から十分推測できる。尚これら一連の論文の扱っている主題が同じものであることは、§4

の一般的定式化(詳しくは拙著[5])を眺められれば明らか
となろう。

従来ノベース推論は、純理論統計家(統計学の中に問題を
求めそれを数学的に処理している研究家)の間で諸概念(許
容性とか完全性)の説明の爲に便宜的に用いられていたが、
実践家や実践的統計家からは(先験分布を勝手に与えて了う
ことから)実用上問題にならぬとそっぽを向かれてきた。純
理論家の多くは実際問題を直視せず、Robbinsの提唱したこの
経験的ベースの立場を受け容れようとしない。理論家と実践
家との間のこのような溝を埋めるのに“経験的ベース”は恰
好な主題であると考えられるので、ぜひ論文[1],[3]の一読
をお勧めする。これらに関連した文献については[5]の参考
文献を参照されたい。

参 考 文 献

- [1] H. Robbins, Proc. 2nd Berkeley Symp. (1951), 131—148.
- [2] ———, Ann. Math. Statist. 26 (1955), 37—51.
- [3] ———, Proc. 3rd Berkeley Symp. (1955), 157—163.
- [4] ———, Ann. Math. Statist. 33 (1964), 1—20.
- [5] G. Suzuki, Research Memorandum No. 19 (1968)
(Inst. Statist. Math.)

§2. Robbins の論文 [1], [2] のスケッチ

H. Robbins は論文 [1] に於て次のような問題を提示した。観測可能な確率変数 X は、分散が 1, 平均が θ の正規分布に従うものとする。ここでパラメータ θ の値は未知であるが、 -1 か $+1$ のいづれかの値しかとり得ないものとする。今、 n 個の独立な (無関係な) 観測 (X_1, X_2, \dots, X_n) が与えられているとする。ここで X_i は正規分布 $N(\theta_i, 1)$ に従うことが判っているが、 θ_i の値は -1 か $+1$ かは判らない。しかし、その値がいづれであるかの決定 (decision) を強いられている。目的は間違つた決定 (θ の値が $+1$ のときに -1 と判定する場合と丁度反対の場合の 2 種類ある) を下す割合の期待値をできるだけ小さくするような決定法則を見出すことである。(2 種類の失敗を同じ比重で扱うという仮定や分布が正規族であるといった仮定は除き得ることを後の論文 [2] で示している。)

このような複合型決定問題、即ち観測値ベクトル $X = (X_1, X_2, \dots, X_n)$ からパラメータベクトル $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ を推定する問題に於ける決定法則 (decision rule) D とは、 $u_i(X)$ の確率で $\theta_i = +1$ ($1 - u_i(X)$ の確率で $\theta_i = -1$) と定めるような n 組の函数 $u = (u_1, u_2, \dots, u_n)$ も対応させることである。さて、真のパラメータベクトルが $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ のとき、

このような決定 D を用いてこおむるリスク (向違った分類の割合の期待値) は

$$(1) \quad r(D, \underline{\theta}) = E \left\{ \frac{1}{n} \sum_{i=1}^n \left[\frac{1+\theta_i}{2} - \text{sgn}(\theta_i) u_i(X) \right] \mid \underline{\theta} \right\}$$

$$= P(\underline{\theta}) - \frac{1}{n} \sum_{i=1}^n \text{sgn}(\theta_i) \int u_i(x) \phi(x, \underline{\theta}) dx$$

と表わせる。但し、

$$P(\underline{\theta}) = \frac{1}{2^n} \prod_{i=1}^n (1 + \theta_i),$$

$$\phi(x, \underline{\theta}) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \theta_i)^2}.$$

従って真のパラメータベクトル $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$ に対する最尤推定量は

$$\hat{\theta}_i = \text{sgn}(X_i), \quad i = 1, 2, \dots, n$$

で与えられる。この決定法則も \tilde{D} と表わそう。対応する u -函数は

$$u_i(x) = u(x_i) = \begin{cases} 1, & x_i > 0, \\ 0, & x_i < 0, \end{cases} \quad i = 1, 2, \dots, n$$

であるから、容易に

$$r(\tilde{D}, \underline{\theta}) = \Phi(-1) = 0.1587$$

となることが示され (Φ は標準正規分布の分布函数)、而もこの法則が唯一のミニマックス解となることも判る。

ある決定法則 D が simple であるというのは、対応する u -函数が

$$u_i(x) = u(x_i), \quad i=1, 2, \dots, n$$

の形をしている場合である。simpleな法則 D に対するリスク
(1) は

$$(2) \quad r(D, \underline{\theta}) = p(\underline{\theta}) - \int \{p(\underline{\theta})\varphi(x-1) - [1-p(\underline{\theta})]\varphi(x+1)\} u(x) dx$$

の如く $p(\underline{\theta})$ の一次関数として書き表わせる。(φ は標準正規分布の密度関数) 任意の λ ($0 \leq \lambda \leq 1$) に対して

$$u(x) = u_\lambda(x) = \begin{cases} 1 & \lambda\varphi(x-1) - (1-\lambda)\varphi(x+1) > 0 \\ 0 & \lambda\varphi(x-1) - (1-\lambda)\varphi(x+1) \leq 0 \end{cases}$$

と置く。この u_λ に対応する決定法則 D_λ は

$$\theta_i = \text{sgn} \left(X_i - \frac{1}{2} \log \frac{1-\lambda}{\lambda} \right), \quad i=1, 2, \dots, n$$

で、リスクは

$$r(D_\lambda, \underline{\theta}) = p(\underline{\theta}) \alpha(\lambda) + [1-p(\underline{\theta})] \beta(\lambda)$$

で与えられる。ここで

$$\alpha(\lambda) = \Phi \left(-1 + \frac{1}{2} \log \frac{1-\lambda}{\lambda} \right), \quad \beta(\lambda) = \Phi \left(-1 - \frac{1}{2} \log \frac{1-\lambda}{\lambda} \right).$$

特に $D_{\frac{1}{2}} \equiv \tilde{D}$ であるから、 $p(\underline{\theta}) = \frac{1}{2}$ でない限り確に \tilde{D} より良い simple な決定法則 $D_{p(\underline{\theta})}$ がある筈で、そのリスクは

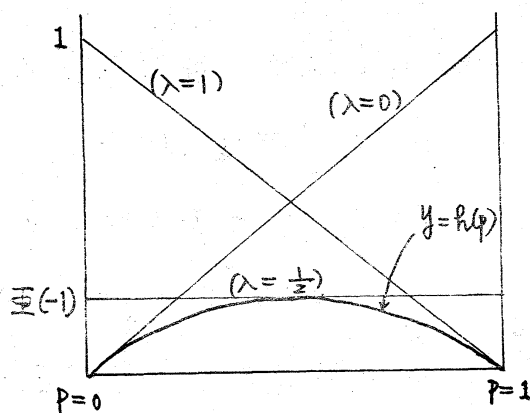
$$(3) \quad h(p(\underline{\theta})) = p(\underline{\theta}) \alpha(p(\underline{\theta})) + [1-p(\underline{\theta})] \beta(p(\underline{\theta}))$$

である。実際曲線 $h(p)$ は、直線群

$$\{ y = \alpha(\lambda)p + \beta(\lambda)(1-p); 0 \leq \lambda \leq 1 \}$$

の包絡曲線で、直線 $y = \alpha(\frac{1}{2})p + \beta(\frac{1}{2})(1-p) \equiv \Phi(-1)$ の下方に在ること判る。

しかし $p(\theta)$ の値が既知というのは稀であつて、多くの場合このような決定法則を得ることはできない。唯 $p(\theta)$ も推定することが可能なときには $\lambda = \hat{p}$ ($p(\theta)$ の推定量) と置いた決定法則 R_λ が利用でき、而も \hat{p} が比較的



良い推定量であれば、この決定法則の方が \tilde{D} より (少くともサンプル数が大きいとき) 良いと期待できるかも知れない。但し推定量 \hat{p} が確率変数、従つて λ が α の函数になつてゐるから、このようにして得られる決定法則はもはや simple ではなくなる。

さて、 α の最尤推定量を用いて $p(\theta)$ に対する最尤推定量

$$\hat{p} = \hat{p}(\alpha) = \frac{1}{n} \{ \text{no. of } i \text{ for which } X_i > 0 \}$$

を得る。これは

$$E\{\hat{p} | \alpha\} = [1 - 2\alpha(-1)] p(\theta) + \alpha(-1)$$

であるから $p(\theta) = \frac{1}{2}$ のとき以外不偏でない。そこで

$$Z = [\hat{p} - \alpha(-1)] / [1 - 2\alpha(-1)]$$

と修正すれば不偏推定量になる。Z の分散は

$$V\{Z | \alpha\} = \frac{\alpha(-1)[1 - \alpha(-1)]}{n[1 - 2\alpha(-1)]^2} \doteq \frac{0.29}{n}.$$

更に小さな分散の不偏推定量は

$$V = \frac{1}{2}(1 + \bar{X}), \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

で与えられる。実際初等的計算によって

$$E\{V|\theta\} = P(\theta), \quad V\{V|\theta\} = \frac{1}{4n} = \frac{0.25}{n}$$

となることも示される。唯 $0 \leq P(\theta) \leq 1$ であるから、これを

$$\tilde{V} = \begin{cases} 0 & V \leq 0 \\ V & 0 < V < 1 \\ 1 & V \geq 1 \end{cases}$$

のように修正した方が妥当であろう。但しこの \tilde{V} には多少の偏りがある。さてこの \tilde{V} を $P(\theta)$ に対する推定量としてとったとき、対応する決定法則 D^* は

$$(4) \quad \theta_i = \text{sgn}(X_i - X^*), \quad i = 1, 2, \dots, n$$

で与えられる。ここで

$$X^* = \begin{cases} +\infty, & \bar{X} \leq -1, \\ \frac{1}{2} \log \frac{1-\bar{X}}{1+\bar{X}}, & -1 < \bar{X} < 1, \\ -\infty, & \bar{X} \geq 1. \end{cases}$$

この D^* は勿論 simple ではないが、リスクは

$$r(D^*, \theta) = h(P(\theta), n)$$

で与えられる。ここで

$$h(p, m) = p \Phi(-2p\sqrt{m}) + (1-p) \Phi(-2(1-p)\sqrt{m}) \\ + \int_{-2(1-p)\sqrt{m}}^{2p\sqrt{m}} g(x; p, m) \varphi(x) dx,$$

$$g(x; p, n) = p \Phi \left[\sqrt{\frac{n}{n-1}} \left(-1 - \frac{x}{\sqrt{n}} + \frac{1}{2} \log \frac{1-p+x/2\sqrt{n}}{p-x/2\sqrt{n}} \right) \right] \\ + (1-p) \Phi \left[\sqrt{\frac{n}{n-1}} \left(-1 - \frac{x}{\sqrt{n}} - \frac{1}{2} \log \frac{1-p+x/2\sqrt{n}}{p-x/2\sqrt{n}} \right) \right].$$

明らかに

$$\lim_{n \rightarrow \infty} h(p, n) = h(p).$$

$h(p)$ は (3) で与えられた函数であるから、

$$\lim_{n \rightarrow \infty} r(D^*, \underline{\theta}) = h(p(\underline{\theta})) \leq \Phi(-1).$$

従つて (4) で与えられる (non-simple な) 法則 D^* がミニマックスな決定法則 \tilde{D} より劣るとは何人も断言できないだろう。

Hannan - Robbins は論文 [2] に於て、この議論をもつと一般的に記述した。まず θ のとり得る値は便宜上 0 又は 1 とし (パラメータ空間は勝手な抽象空間の部分集合で良く元の個数のみが本質的である)、2種類の同違つた判断にウェイトをつけた。即ち $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$ に対して $\underline{d} = (d_1, d_2, \dots, d_n)$ という決定もする ($d_i, 1-d_i$ の確率で夫々 $\theta_i = 1, \theta_i = 0, i=1, 2, \dots, n$ と判断する) ときの損失を

$$W(\underline{d}, \underline{\theta}) = \frac{1}{n} \sum_{i=1}^n [a \theta_i (1-d_i) + b (1-\theta_i) d_i]$$

で与える。(a, b は適当な正数で、 $a = b = 1$ の場合が [1] で扱われている。)

前と同様に観測値ベクトル $\underline{X} = (X_1, X_2, \dots, X_n)$ が与えられている。各 X_i は独立で分布函数として $F(\cdot, \theta_i)$ をもっている、

但し $\theta_i = 0$ 又は 1 , $i=1, 2, \dots, n$ とする。

$$\mu_0(B) = \int_B dF(x, \theta), \quad \theta = 0, 1,$$

$$\mu(B) = \mu_0(B) + \mu_1(B), \quad B \text{ は } n\text{-次元ボレル集合}$$

によって有限な測度 μ を定義し、 μ_θ の μ に因する一般化された密度を $f(\cdot, \theta)$ と表わす ($\theta = 0, 1$)。真のパラメータベクトルが $\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$ のとき決定函数 $\underline{u} = (u_1, u_2, \dots, u_n)$ を用いて ($u_i(x)$, $1-u_i(x)$ の確率で夫々 $\theta_i = 1$, $\theta_i = 0$ と判断して) におむるリスクは

$$\begin{aligned} r(\underline{u}, \underline{\theta}) &= E \{ W(\underline{u}(X), \underline{\theta}) | \underline{\theta} \} \\ &= \frac{1}{n} \int \left[\sum_{i=1}^n [a\theta_i(1-u_i(x)) + b(1-\theta_i)u_i(x)] f(x, \underline{\theta}) \right] d\mu^n \end{aligned}$$

である。ここで $f(x, \underline{\theta}) = \prod_{i=1}^n f(x_i, \theta_i)$ は真のパラメータが $\underline{\theta}$ のときの直積測度 μ^n に因する $X = (X_1, X_2, \dots, X_n)$ の一般化された直積密度である。特に simple な $\underline{u} = (u_1, u_2, \dots, u_n)$ に対しては、

$$(5) \quad r(\underline{u}, \underline{\theta}) = p(\underline{\theta}) - \int [ap(\underline{\theta})f(x, 1) - b(1-p(\underline{\theta}))f(x, 0)] u(x) d\mu$$

で与えられる。但し

$$p(\underline{\theta}) = \frac{1}{n} \sum_{i=1}^n \theta_i.$$

従って

$$u_{\lambda}(x) = \begin{cases} 1, & x \in C[\lambda] = \{x : a\lambda f(x, 1) > b(1-\lambda)f(x, 0)\}, \\ 0, & x \notin C[\lambda] \end{cases}$$

と置くとき、(5) は simple な $u_{p(\underline{\theta})}$ によって最小値

$$\phi(p(\underline{\theta})) = p(\underline{\theta}) - ap(\underline{\theta})\mu_1(C[p(\underline{\theta})]) + b(1-p(\underline{\theta}))\mu_0(C[p(\underline{\theta})])$$

が与えられる。

さて $p(\theta)$ の値が未知の場合このような simple な決定函数は定められないが、 $p(\theta)$ に対するある一致推定量 $\hat{p}_n = \hat{p}_n(x)$ が与えられているとき ([2] では $p(\theta)$ に対する不偏推定量のある部分クラスの中で分散を最小にするものの求め方についても記されている)、

$$(6) \quad u_i^*(x) = \begin{cases} 1, & x_i \in C[\hat{p}_n(x)], \\ 0, & x_i \in C[\hat{p}_n(x)], \quad i=1, 2, \dots, n \end{cases}$$

によって定義される non-simple な函数 $\underline{u}^* = (u_1^*, u_2^*, \dots, u_n^*)$ に対して次のような結果が成立する。

“任意の $\varepsilon > 0$ に対し

$$Pr\{W(\underline{u}^*(x), \theta) \leq \phi(p(\theta)) + \varepsilon, \text{ for all } n \geq n_0(\varepsilon)\} \geq 1 - \varepsilon$$

となるような整数 $n_0(\varepsilon)$ が得られる。従って、

$$r(\underline{u}^*, \theta) \leq \phi(p(\theta)) + \varepsilon, \text{ for all } n \geq n_1(\varepsilon)$$

となるような整数 $n_1(\varepsilon)$ も選ぶことができる。”

即ち (6) で与えられる決定函数は、サンプル数が充分大であれば、simple な決定函数のクラスの中に於る最良なものと同じ程度に良い決定を為し得ることが示された。彼等は尚、simple なものも含むより広い (各段階での決定を置換してもリスクを変えないようなものの) クラスの中での最良のものと比較しても決定函数 (6) が見劣りしないことを示している。

§3. Robbins の論文 [3], [4] のスケッチ

X を離散的な値をとる確率変数 (連続的な値をとる場合も
平行した議論ができる)、 Λ もある (未知の) 先験分布

$$G(\lambda) = \Pr\{\Lambda \leq \lambda\}$$

に従う確率変数とする。 $\Lambda = \lambda$ を与えたときの X の条件付分
布

$$p(x|\lambda) = \Pr\{X = x | \Lambda = \lambda\}$$

は判っているものとする。従って X の条件付でない分布は

$$P_G(x) = \Pr\{X = x\} = \int p(x|\lambda) dG(\lambda)$$

で与えられる。 Λ に対する $\varphi(X)$ の形の推定量の自乗平均誤
差は

$$\begin{aligned} E\{\varphi(X) - \Lambda\}^2 &= E\{E\{\varphi(X) - \Lambda\}^2 | \Lambda\} \\ &= \sum_x \int p(x|\lambda) [\varphi(x) - \lambda]^2 dG(\lambda) \end{aligned}$$

で与えられ、 X を与えたときの Λ の条件付期待値

$$(1) \quad \varphi_G(X) = E\{\Lambda | X\}$$

によって最小になる。ところで

$$\varphi_G(x) = \int p(x|\lambda) \lambda dG(\lambda) / \int p(x|\lambda) dG(\lambda)$$

であるから、分布 G が判ればベイズ推定量 (1) を作ることが
できる。しかし多くの場合分布 G は未知であるから、このよ
うな推定量を作ることはできない。それで従来は、不偏性と
かミニマックス基準とかの G に関する情報も必要としない判

約条件のみを課して推定量を定めてきた。

しかし次のような場合には事情が一寸異なる。確率変数の組の同一分布に従う独立系列 $\{(X_i, \Lambda_i), i=1, 2, \dots\}$ が与えられているとする。各 (X_i, Λ_i) の分布の構造は最初に述べた確率変数の組 (X, Λ) と同一で、 Λ_i の方は観測可能でない確率変数とする。 X_1, X_2, \dots, X_n に関する観測値を全部利用して、 Λ_n に対する推定量を作りたい。Robbins [3] の提唱した問題というのは、次のようなものである。

“ X_1, X_2, \dots, X_{n-1} に関する情報も函数形を定めるのに用いて、(1) で与えられる φ_G に確率収束するようなランダムな函数列 $\{\hat{\varphi}_n(\cdot) = \hat{\varphi}_n(X_1, X_2, \dots, X_{n-1}; \cdot), n=2, 3, \dots\}$ を構成できるか？”

一般的には否定的答えしか出ない。何故ならそのようなことが可能となるのは、原則として未知の先験分布 G に対する一致推定量 G_n が得られねばならず、而も Λ_i を直接観測できないうえに X_i の観測値を通して為されなければならない。これは *mixture* より *mixing* を求める問題、言い換えるとヤークフッドホルム型積分方程式を解く問題に帰着され、無条件では解けないからである。(少なくとも *mixture* と *mixing* の間の一対一対応 — Teicher: *Ann. Math. Stat.* 34(1963), 1265-9 — の条件が必要である。唯一対一対応がないからといってせじを投げるのも早計である — Robbins [4; §5].)

Robbins [3, 4] は核 $p(x|\lambda)$ が指数型分布族の場合に、
 ては肯定的答えることを例示した。まず条件付密度が

$$p(x|\lambda) = \lambda^x f(x) h(\lambda), \quad x = 0, 1, 2, \dots$$

の形をしている場合 (ポアソン分布や負の二項分布) を考へる。

$$\begin{aligned} \int p(x|\lambda) \lambda dG(\lambda) &= \int \lambda^{x+1} f(x) h(\lambda) dG(\lambda) \\ &= \frac{f(x)}{f(x+1)} P_G(x+1) \end{aligned}$$

であるから、(1) の函数に対して

$$\varphi_G(x) = f(x) P_G(x+1) / f(x+1) P_G(x)$$

なる関係が得られる。そこで各 m に対し

$$P_m(x) = \frac{1}{n} \{ \text{no. of } X_1, \dots, X_n \text{ which are equal to } x \}, \quad x = 0, 1, \dots$$

を定義すると、これは $P_G(\cdot)$ に対する一致推定量となる。

従つて

$$\begin{aligned} \hat{\varphi}_m(x) &= \hat{\varphi}_m(X_1, X_2, \dots, X_{n-1}; x) \\ (2) \quad &= f(x) P_{m-1}(x+1) / f(x+1) P_{m-1}(x) \end{aligned}$$

がベイズ推定量 (1) に対する一つの一致推定量を与える；

$$(3) \quad E[\hat{\varphi}_m(X_m) - \Lambda_m]^2 \longrightarrow E[\varphi_G(X) - \Lambda]^2.$$

特にポアソン核

$$p(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, \dots; \lambda > 0$$

で先験分布 G がガンマタイプ°

$$G'(\lambda) = \frac{h^b}{\Gamma(b)} \lambda^{b-1} e^{-h\lambda}, \quad \lambda, b, h > 0$$

の場合を考へてみる。ベイズ推定量 (1) は

$$\varphi_G(X) = [X+b]/[1+r]$$

で、その自乗平均誤差は

$$E[\varphi_G(X) - \Lambda]^2 = b/r(1+r).$$

一方通常の (最小分散不偏) 推定量は X 自身で、誤差は

$$E[X - \Lambda]^2 = E\Lambda = b/r$$

と計算される。従ってこのような場合、(2) で構成される推定量 $\hat{\varphi}_m(X_m)$ が通例のポアソン推定量 X_m より良くなる場合が十分期待できる (少なくとも漸近的には (3) の関係から前者の方が後者より efficient になる)。

最後にラプラスアンタイプ

$$p(x|\lambda) = e^{-\lambda x} f(x) h(\lambda)$$

の場合も考える。微分と積分の演算の交換が許されるものとする

$$\begin{aligned} \varphi_G(x) P_G(x) &= \int p(x|\lambda) \lambda dG(\lambda) \\ &= f(x) \int \lambda e^{-\lambda x} h(\lambda) dG(\lambda) \\ &= f(x) \frac{d}{dx} \{P_G(x)/f(x)\} \end{aligned}$$

であるから

$$\varphi_G(x) = \frac{d}{dx} \log \left\{ \frac{P_G(x)}{f(x)} \right\}$$

なる関係を得る。この場合にも $P_m(x)$ を適当な係数 α を挿入するように修正した $\hat{P}_m(x)$ で $P_G(x)$ を近似することにより、(3) の関係が成立するような (2) に対応する $\hat{\varphi}_m(\cdot)$ が得られる。

Robbins は更に [4] に於て 行動空間 (action space) を導入して、このようなベイズ推定やベイズ型仮説検定問題等を含むような一般論を展開した。しかし [1], [2] の複合型決定問題も含む程一般的ではない。彼は漸近的にベイズ決定法則に (平均的に) 近づくような (例えば (3) の関係が成立するようなもので、かかる性質を *asymptotically optimal* と呼んでいゝる) ランダムな決定函数を得る為の条件を示しているが、それは殆んど定義を少し書き直したようなものに過ぎない。更に先験分布 (mixing) の推定問題にも触れているが、一般的な場合にはやはり見るべき結果を与えていない。唯パラメータ空間が有限個の元から成る場合については、有限個の函数のクラス $\{f(\cdot, \theta_1), \dots, f(\cdot, \theta_k)\}$ から互に直交する成分をとり出し、これを用いて先験分布も推定するという有効な手法を示した。論文 [4] で定式化されたものは、次の §4 で述べるより一般の問題に含まれるので、これ以上の直接的紹介は省略する。

§4. 一般的定式化

\mathcal{X} , \mathcal{Y} , A を夫々標本, パラメータ, 行動空間とする。更に \mathcal{X} 上の確率測度の集合

$$\{ \nu(\cdot, \theta) : \theta \in \mathcal{Y} \}$$

が与えられていて、ある σ -有限な測度 μ によって dominate されている。 μ に関する R-N. 微係数を

$$\frac{dV(\cdot, y)}{d\mu} = f(\cdot, y)$$

と表わす。 $w(a, y)$ は与えられた $A \times Y$ 上の有界可測函数とする (非負値でなくとも構わない)。 \mathcal{A} から A の中への写像 d を決定函数と呼び、このような d の集合を D で表わす。 $D \times Y$ 上にリスク函数

$$(1) \quad r(d, y) = \int_{\mathcal{A}} w(d(x), y) f(x, y) \mu(dx)$$

が定義される。 これを用いて Y 上の距離

$$\delta(y_1, y_2) = \sup_{d \in D} |r(d, y_1) - r(d, y_2)|$$

を定義する。 この δ の意味での Y の閉集合をすべて含むような最小の σ -体を \mathcal{Y} と表わす。 全く同様に \mathcal{D} も定義される。 空間 Y が可分ならば、函数 (1) は $\mathcal{D} \times \mathcal{Y}$ -可測である。

\mathcal{M}_0 を Y 上の確率測度の全体とし、 $D \times \mathcal{M}_0$ 上に

$$r_i(d, \nu) = \int_Y r(d, y) \nu(dy)$$

なるリスク函数、更に \mathcal{M}_0 上にも

$$r_i(\nu) = \inf_{d \in D} r_i(d, \nu)$$

なるベースリスクを定義する。 後で \mathcal{M}_0 の元を推定する問題を考えるが、^{直接}得られる推定量は一般に確率 1 で \mathcal{M}_0 に属さない。

\mathcal{M}_0 へ入るよう修正できるが、リスクを最小にするという観点からはそのような方法が最良とは言えない。 そこでベースの

概念をもっと拡張して置く。 \mathcal{M} を \mathcal{Y} 上の有限な有符号測度の全体とする。任意の $\nu \in \mathcal{M}$ に対して、 $\nu = \nu^+ - \nu^-$ (ν^+ , ν^- は有限測度) なる表現が一通りにできるから、 $D \times \mathcal{M}$ 上のリスク関数を

$$r^*(d, \nu) = \int_{\mathcal{Y}} r(d, y) |\nu|(dy) = |\nu|(\mathcal{Y}) r_1(d, \tilde{\nu})$$

で与えることにする。但し

$$|\nu| = \nu^+ + \nu^-, \quad \tilde{\nu} = [|\nu|(\mathcal{Y})]^{-1} \nu \in \mathcal{M}_0.$$

$|\nu| = |\nu'|$ となる ν (同一リスクも与える ν) の集り全体で \mathcal{M} の一個の元と見做すことにする。 \mathcal{M} 上のベースリスク及び距離も次のように定義して置く。

$$r^*(\nu) = \inf_{d \in D} r^*(d, \nu),$$

$$\delta^*(\nu_1, \nu_2) = \sup_{d \in D} |r^*(d, \nu_1) - r^*(d, \nu_2)|.$$

次に \mathcal{Y} 上のある確率測度 m を一選んで固定する。任意の $p \geq 1$ に対し、 $L^p = L^p(\mathcal{Y}, m)$ を

$$\|g\|_p = \left\{ \int_{\mathcal{Y}} |g(y)|^p m(dy) \right\}^{1/p} < \infty$$

となるような \mathcal{Y} 上の実数値函数 g の全体として

$$\mathcal{M}^p = \mathcal{M}^p(\mathcal{Y}, m) = \left\{ \nu \in \mathcal{M} : \nu \ll m, \frac{d\nu}{dm} \in L^p \right\}$$

なる \mathcal{M} の部分クラスを定義する。 \mathcal{M}_0^p についても同様に定義する。 w は有界であるから

$$\bar{w}(y) = \sup_{a \in A} |w(a, y)|$$

なる函数が定義できて、 $\bar{w} \in L^1 = L^1(\mathcal{Y}, m)$, $\frac{1}{p} + \frac{1}{q} = 1$, である。

\mathcal{M}^p に次のような距離を入れる。

$$\delta_p(\nu_1, \nu_2) = \|\bar{w}\|_q \cdot \left\| \frac{d\nu_1}{d\mu} - \frac{d\nu_2}{d\mu} \right\|_p.$$

Schwartz の不等式から明らかのように

$$\delta^*(\nu_1, \nu_2) \leq \delta_p(\nu_1, \nu_2)$$

である。即ち δ_p の方が δ^* より強い位相である。

\mathcal{M}^p より D の中への可測写像 $V = V(\nu) \in \text{decision procedure}$ と呼ぶ。ある $\varepsilon \geq 0$ に対して

$$r^*(V(\nu), \nu) - r^*(\nu) \leq \varepsilon, \quad \text{for all } \nu \in \mathcal{M}^p$$

となる procedure V を ε -optimal といい、 V^ε と表わす。

\mathcal{H} の σ -体 \mathcal{G} に、確率測度 m を用いて

$$p(E, F) = m(E \Delta F) = m(E - F) + m(F - E), \quad E, F \in \mathcal{G}$$

な距離を導入する。以下距離空間 (\mathcal{G}, p) が可分であると仮定する。位相解析の定理 (例えば Zaanen: Linear Analysis p.74) によつて、この仮定の下で $(\mathcal{M}^p, \delta_p)$ が可分になる。(この距離の意味での可分性のみからではこの性質は導れない。)

次に \mathcal{M}^p に於ける任意の系列 $\underline{z} = \{z_n\}$ を考へて、

$$z_n^{(0)} = z_n, \quad n = 1, 2, \dots$$

$$z_n^{(i)} = \frac{1}{n} \sum_{k=1}^n z_k^{(i-1)}, \quad n = 1, 2, \dots; \quad i = 1, 2, \dots$$

なる系列を順次定義する。ある i に対し

$$\lim_{n \rightarrow \infty} \delta_p(z_n^{(i)}, z) = 0$$

となる $z \in \mathcal{M}^p$ が唯一つ存在するとき、 $\underline{z} = \{z_n\}$ は (H, i) -収斂

すると言ひ、 $\underline{z} \in H^{(i)}$ と書くことにする。明らかに $H^{(0)} \subset H^{(1)} \subset \dots$ である。i 次の極限も $\underline{z}^{(i)}$ と表わす。

さて以上の準備の下で、一般的な決定問題を定式化しよう。 (Ω, \mathcal{G}, P) をき々の基本確率空間とする。この空間より標本空間 ω 及びパラメータ空間 \mathcal{Y} の中への可測写像 $X = X(\omega)$, $Y = Y(\omega)$ を夫々ランダムサンプル, ランダムパラメータと呼ぶ。 $\{(X_n, Y_n), n=1, 2, \dots\}$ をランダムサンプル, ランダムパラメータの組の独立な系列とする。 Y_n の分布は $\mu_n \in \mathcal{M}^{\mathcal{Y}}$ (同一分布でなくともよい), $Y_n = y$ を与えたときの X_n の分布は $\nu(\cdot, y)$ (ν は μ に関する微分 $f(\cdot, y)$) で与えられる。次に、 $\underline{z} = \{z_n\} \in H^{(i)}$ とする最小の整数を i_0 とする (かかる有限な i_0 は存在するものと仮定して)。各 $d \in D$ に対し次のような確率変数列を順次定義する。

$$W_n^{(0)}(d) = w(d(X_n), Y_n), \quad n=1, 2, \dots,$$

$$W_n^{(i)}(d) = \frac{1}{n} \sum_{k=1}^n W_k^{(i-1)}(d), \quad n=1, 2, \dots; \quad i=1, \dots, i_0.$$

容易に判るように

$$(2) \quad E\{W_n^{(i_0)}(d)\} = r^*(d, \underline{z}_n^{(i_0)}).$$

この値を最小にするような $d \in D$ を見出すことが望ましいが、分布 $\mu_n^{(i_0)}$ に関する完全な情報を与えられていない場合には不可能である。しかし、 n 着目の決定を為す前に X_1, X_2, \dots, X_n

に関する観測が与えられている場合は常である。そこでこの
 ような先験情報を活用して

$$\delta_p(\hat{\zeta}_m, \zeta_m^{(i)}) \xrightarrow{P} 0$$

となる推定量 $\hat{\zeta}_m = \hat{\zeta}_m(X_m)$ ($\zeta^{(i)}$ に対する一致推定量) が与えら
 れたとすると、(2) を最小にする $d(\zeta_m^{(i)})$ に対するベイズ) と
 同じ程度良い "ランダム" な決定法則を作ることができるのでは
 ないかと期待するのは自然であろう。実際次の結果が得られ
 る。

定理 1. $L^p = L^p(Y, m)$ に於て $g_m^{(i)} = \frac{d\zeta_m^{(i)}}{dm}$ に対する推定量 $\hat{g}_m =$
 $\hat{g}_m(X_m)$ が、関係

$$(3) \quad E\{\|\hat{g}_m - g_m^{(i)}\|_p\} \leq K m^{-\Delta} \quad (K, \Delta \text{ はある正定数})$$

を満足するものとする。

$$\hat{\zeta}_m(E) = \int_E \hat{g}_m(y) m(dy), \quad E \in \mathcal{Y}$$

によって \mathcal{M}^p に於ける推定量を定義する。 $\varepsilon_m = \|\bar{w}\|_q m^{-\Delta}$ と置いて

$$\hat{d}_m = v^{\varepsilon_m}(\hat{\zeta}_m), \quad m = 1, 2, \dots$$

により "ランダム" な決定函数の系列を定義すると、

$$E\{r^*(\hat{d}_m, \zeta_m^{(i)})\} \leq r^*(\zeta_m^{(i)}) + (K+1)\|\bar{w}\|_q m^{-\Delta}.$$

従つて、任意に固定された $\varepsilon > 0$ に対し

$$P\{r^*(\hat{d}_m, \zeta_m^{(i)}) - r^*(\zeta_m^{(i)}) \geq \varepsilon\} \leq \varepsilon^{-1}(K+1)\|\bar{w}\|_q m^{-\Delta}.$$

即ち

$$(4) \quad r^*(\hat{d}_m, \zeta_m^{(i)}) - r^*(\zeta_m^{(i)}) \xrightarrow{P} 0.$$

[注1.] ε -optimal な decision procedure の作り方について簡単に触れて置く。距離空間 $(\mathcal{M}^P, \delta_P)$ が可分であるから、任意の $\varepsilon > 0$ に対して次の関係を満たすような高々可附番個の \mathcal{M}^P の (空でない) 可測部分集合の系列 $\{\mathcal{M}_n, n=1, 2, \dots\}$ がとれる:

$$\bigcup_{n=1}^{\infty} \mathcal{M}_n = \mathcal{M}^P, \quad \mathcal{M}_i \cap \mathcal{M}_j = \phi \quad (i \neq j),$$

$$\zeta_1, \zeta_2 \in \mathcal{M}_n \text{ なら } \delta_P(\zeta_1, \zeta_2) \leq \frac{\varepsilon}{3}, \quad n=1, 2, \dots$$

次に各 \mathcal{M}_n から任意に一つ ζ_n をとり出し、

$$r^*(d_n, \zeta_n) - r^*(\zeta_n) \leq \frac{\varepsilon}{3}$$

となるような $d_n \in D$ を一つ選んで

$$v(\zeta) = d_n, \quad \text{for all } \zeta \in \mathcal{M}_n, \quad n=1, 2, \dots$$

なる procedure を定義すると、これが ε -optimal なものであることは容易に証明され、而も v の値域が D の元の高々可附番個のものから成っていることも明らかである。

[注2.] (4) の関係を成立せしめるようなランダムな決定列 $\hat{d} = \{\hat{d}_n\}$ も $\zeta = \{\zeta_n\} \in H^{(i_0)}$ に対して (H, i_0) -optimal であると呼ぶことにすれば、§2 で紹介した複合型決定問題 (最適 simple 決定法則に近似も与える問題) は $(H, 1)$ -optimal なものを見出す問題に還元される。尚その際、 Y_n (Nature) の与える分布 (strategy) ζ_n は実 θ_n で全質量を与え、一実分布 (pure strategy) であると考えれば完全に一致する。又 §3 で記述したベイズ推論の場合には、 $(H, 0)$ -optimal なものを

見出す問題になる(この場合各 g_n は同一分布 $g = \frac{d\eta}{d\alpha_m}$ に従っているものと考えれば良い)。

定理 1 によって、(3) の関係を満たすような推定量 $\hat{g}_n = \hat{g}_n(X_n)$ を作れば問題は完結する。そこで $g_m^{(1)}$ ($i_0 = 1$ の場合) に対する推定を考えてみよう。(一般の i_0 にしたましても平行した議論が可能であり、又 $g_m = g$ ($m = 1, 2, \dots$) と置くことにより未知の唯一の先験分布 g の推定に於る結果も含むことになる。) 更に(直交展開が使える等の)種々の観点から、 $p=2$ の場合を考えるのが妥当であろう。

即ち問題は、独立で観測可能な確率変数の系列 $\{X_n, n=1, 2, \dots\}$ が与えられていて、各 X_n は(必ずしも同一でない)分布

$$(5) \quad f(x, g_n) = \int_{\mathcal{Y}} f(x, y) g_n(y) m(dy), \quad n=1, 2, \dots$$

に従っているものとする。更に核函数 $f(\cdot, \cdot)$ は直積空間 $\mathcal{X} \times \mathcal{Y}$ で可測で、Identifiability の条件

$$m\{y \in \mathcal{Y} : g_0(y) \neq g_j(y)\} > 0 \text{ ならば } \mu\{\alpha \in \mathcal{X} : f(\alpha, g_0) \neq f(\alpha, g_j)\} > 0$$

を満しているとする。

最初にパラメータ空間が有限個から成る場合 $\mathcal{Y} = \{y_1, y_2, \dots, y_k\}$ について考える。 m は各事に $1/k$ の質量を与える離散型測度にとる。

$$g_m(y_j) = R \alpha_j^{(m)}, \quad j=1, 2, \dots, k; \quad n=1, 2, \dots$$

$$f(x, y_j) = \beta_j(x), \quad j=1, 2, \dots, k$$

と置くと、 X_m の分布 (5) は

$$(6) \quad f(x; \underline{\alpha}^{(m)}) = \sum_{j=1}^k \alpha_j^{(m)} \beta_j(x)$$

と表わされる。

$\mathcal{P}_k = \{X_1, X_2, \dots, X_k\}$ が $\beta_k = (\beta_1, \dots, \beta_k)$ に関する測度空間

(\mathcal{X}, μ) の正則分割であるとは

a) 各 X_i は μ -可測, $\mathcal{X} = \bigcup_{i=1}^k X_i$, $\mu(X_i \cap X_j) = 0$ ($i \neq j$)

b) 行列 $D_k = D_k(\mathcal{P}_k; \beta_k) = \left[\int_{\mathcal{X}} \beta_i(x) \beta_j(x) \mu(dx) \right]$ が正則、但し β_j は集合 X_j の特性函数である。

このような正則分割を用いて

$$(7) \quad \varphi_i(x) = \sum_{j=1}^k d_{ij}^* \beta_j(x), \quad i=1, 2, \dots, k$$

なる函数が定義できる。 d_{ij}^* は D_k の逆行列の (i, j) -元とする。

定理 2. 独立な確率変数列 $\{X_m, m=1, 2, \dots\}$ は夫々 (6) の分布に従うものとする。 (7) の函数を用いて

$$\hat{\alpha}_i^{(m)} = \hat{\alpha}_i^{(m)}(X_m) = \frac{1}{n} \sum_{k=1}^m \varphi_i(X_k), \quad i=1, 2, \dots, k$$

なる推定量を作れば、

$$E\{\hat{\alpha}_i^{(m)}\} = \bar{\alpha}_i^{(m)} = \frac{1}{n} \sum_{k=1}^m \alpha_i^{(k)}, \quad i=1, 2, \dots, k,$$

$$(8) \quad \sum_{i=1}^k V\{\hat{\alpha}_i^{(m)}\} = E\left\{ \sum_{i=1}^k (\hat{\alpha}_i^{(m)} - \bar{\alpha}_i^{(m)})^2 \right\} \leq d_k^{*2} k/m,$$

但し $(d_k^*)^{-1}$ は行列式 $|D_k|$ の絶対値である。

(8) の不等式より

$$E\{\|\hat{g}_m - g_m^{(1)}\|_2^2\} = E\left\{ k \sum_{i=1}^k (\hat{\alpha}_i^{(m)} - \bar{\alpha}_i^{(m)})^2 \right\} \leq k^2 d_k^{*2} / m$$

であるから、結局

$$E\{\|\hat{g}_n - g_n^{(0)}\|_2\} \leq [E\{\|\hat{g}_n - g_n^{(0)}\|_2^2\}]^{1/2} \leq k d_k^* n^{-1/2}.$$

従ってこの場合には、 $K = k d_k^*$, $\Delta = \frac{1}{2}$ として (3) の関係も満たすような推定量を作れた。

[注3.] (7) で与えられる φ_i 及び (6) に従う確率変数 X_n に対して

$$E\{\varphi_i(X_n)\} = \alpha_i^{(0)}$$

が成立するが、Robbins [4; §7] は一次独立な函数系 $\{\beta_1, \dots, \beta_k\}$ より正規直交系 $\{\beta_1^*, \dots, \beta_k^*\}$ を作りだして、 $\beta_i^*(X_n)$ を $\alpha_i^{(0)}$ の不偏推定として用いた。 φ_i と β_i^* との間には

$$\varphi_i(x) = \beta_i^*(x) + \tilde{\varphi}_i(x) \quad \text{a.e. } \mu, \quad i=1, 2, \dots, k$$

なる関係が成立する。ここで $\tilde{\varphi}_i$ は $\{\beta_1, \dots, \beta_k\}$ によつて span される $L^2(\mathcal{X}, \mu)$ の部分空間に直交する成分である。つまり、Robbins の直交化の方法は $\tilde{\varphi}_i(x) = 0$ a.e. μ , $i=1, 2, \dots, k$ なる $\varphi_R = (\varphi_1, \dots, \varphi_k)$ を用いたに過ぎない。如何なる φ_i をとるのが (分散も小さくする等の意味に於て) 最適かという問題は、そう単純には解決されそうにない。従つて作り方が単純であるという見地からは、吾々の提示した方法のオにヤリ利があると言えるかも知れない。^又完全に直交化する方法は、解析的取扱いを要易にするが、有限でない場合の吾々の truncation に依る方法に対しては用いることができない。

最後に一般の $L^2 = L^2(Y, m)$ の場合の結果を述べよう。この場合には少々きつい条件が附加されるので、実用上未だ難点が残ろう。 $f_m \in L^2$ であるから適当な正規直交系 $\{\psi_j, j=1, 2, \dots\}$ を用いて

$$f_m(y) = \sum_{j=1}^{\infty} \alpha_j^{(m)} \psi_j(y), \quad m=1, 2, \dots$$

の如く展開できる。従つて X_m の分布 (5) は

$$(9) \quad f(x; \alpha^{(m)}) = \sum_{j=1}^{\infty} \alpha_j^{(m)} \beta_j(x), \quad m=1, 2, \dots$$

$$\beta_j(x) = \int_{\mathcal{Y}} f(x, y) \psi_j(y) m(dy), \quad j=1, 2, \dots$$

と表わせる。

次のような条件を置く。ある非負値定数 γ 及び α の分割列 $\{P_k, k=2, 3, \dots\}$ が与えられていて次の条件を満す。即ち、各分割 P_k が $\beta_k = (\beta_{k1}, \dots, \beta_{kr})$ に関する正則分割で、而も整数 k_0 と正定数 C とが存在して

$$d_k^* = |\det D_k(P_k; \beta_k)|^{-1} \leq C k^\gamma, \quad k = k_0, k_0+1, \dots$$

更に (9) に表れる係数が、ある $\delta > \gamma$ に対し次の条件を満す。

$$|\alpha_j^{(m)}| = \left| \frac{1}{\pi} \sum_{l=1}^m \alpha_j^{(l)} \right| \leq B j^{-(\frac{3}{2} + \delta)}, \quad j=1, 2, \dots \quad (B \text{ は正定数}).$$

このような仮定の下で、次の一般的命題が成立する。

定理 3. 各 m に対し

$$k = k(m) = \left[m^{\frac{1}{1+2\delta}} \right] \quad ([\cdot] \text{ はガウス記号})$$

と置く。(7) で与えらるる $\varphi_n = (\varphi_1, \dots, \varphi_n)$ を用いて

$$\hat{f}_m(y) = \sum_{j=1}^{k(m)} \hat{\alpha}_j^{(m)} \psi_j(y) = \sum_{j=1}^{k(m)} \left[\frac{1}{k(m)} \sum_{l=1}^m \varphi_j(X_l) \right] \psi_j(y)$$

なる推定量を作ると

$$E\{\|\hat{g}_n - g_n^{(1)}\|_2^2\} = E\left\{\sum_{j=1}^R (\hat{\alpha}_j^{(n)} - \alpha_j^{(n)})^2 + \sum_{j=R+1}^{\infty} \alpha_j^{(n)^2}\right\} = O(n^{-\frac{2(\delta-\gamma)}{1+2\delta}})$$

従って

$$E\{\|\hat{g}_n - g_n^{(1)}\|_2\} = O(n^{-\frac{\delta-\gamma}{1+2\delta}}).$$

[注4.] 条件式 (10) は技巧的でやゝ不自然な感じを与えるので、各 β_i が非負値函数 (例えば密度函数) である場合に

一つの十分条件を示してみよう。 $\mathcal{X} = \sum_{i=1}^{\infty} \mathcal{X}_i$, $\mu(\mathcal{X}_i \cap \mathcal{X}_j) = 0$,

$i \neq j$ なる可附番個の分割が次の条件を満たしているとする。

$$\det M_{k_0} \geq \frac{1}{c k_0^\delta}$$

$$\int_{\mathcal{X}_k} \beta_k(x) \mu(dx) \geq 1 - \frac{1}{2c} \left\{ \frac{1}{(k-1)^\delta} - \frac{1}{k^\delta} \right\}, \quad k = k_0+1, k_0+2, \dots$$

但し

$$M_k = \left(\int_{\mathcal{X}_j} \beta_i(x) \mu(dx) \right)_{i,j=1,2,\dots,k}, \quad k = k_0, k_0+1, \dots$$

となる正定数 c 及び整数 k_0 が存在する。”

まず Hadamard の不等式を用いて

$$(11) \quad \det M_k \geq \frac{1}{c k^\delta}, \quad k = k_0, k_0+1, \dots$$

なる関係が示される。

$$P_k^* = \left\{ \mathcal{X}_1, \dots, \mathcal{X}_{k-1}, \sum_{i=k}^{\infty} \mathcal{X}_i \right\}, \quad k = 2, 3, \dots$$

なる分割列を考えると、 β_j の非負値性と M_k の定義より

$$\det D_k(P_k^*) \geq \det M_{k-1} - \frac{1}{c} \left\{ \frac{1}{(k-1)^\delta} - \frac{1}{k^\delta} \right\}.$$

これと (11) の不等式より (10) の関係を得る。