

Conditional probability and I-divergence

東工大・理 菅原 昭博 (Akihiro Sugawara)

確率測度 P, Q に対して I-divergence (Kullback-Leibler information number) $I(P|Q)$ は

$$I(P|Q) = \begin{cases} \int \log \frac{dP}{dQ} \cdot dP & , \text{ if } P \ll Q \\ +\infty & , \text{ otherwise.} \end{cases}$$

と定義される。また、確率測度の族 \mathcal{P} と確率測度 Q に対しては

$$I(\mathcal{P}|Q) = \begin{cases} \inf \{ I(P|Q) ; P \in \mathcal{P} \} & , \mathcal{P} \neq \emptyset \\ +\infty & , \mathcal{P} = \emptyset \end{cases}$$

と定義する。

この小論では、標本の数が大きい観測において、二つの確率測度の分離度とされる I-divergence が、どのような向きをするかに視点をあわせる。特に、標本平均による条件付

確率測度と I -divergence の最小化がいかなるかわりを持つかが示される。

§1. I -divergence と仮説検定

仮説検定における才一種の誤りの確率が, I -divergence によって説明されることをみる。

(X, \mathcal{F}) を標本空間 (可測空間), その上の確率分布の族を $\{P_\theta; \theta \in \Theta\}$ とする。母数空間を $\Theta = \{\theta_0, \theta_1\}$ として, 次の単純仮説検定を考える。

$$\begin{cases} H_0; \theta = \theta_0 & (\text{帰無仮説}) \\ H_1; \theta = \theta_1 & (\text{対立仮説}) \end{cases}$$

仮説検定では, (X, \mathcal{F}) 上の $[0, 1]$ への可測写象を検定 (test) と呼び, 検定 φ が与えられたとき,

$$\alpha(\varphi) = E_{\theta_0}(\varphi), \quad \beta(\varphi) = E_{\theta_1}(1 - \varphi)$$

を, それぞれ才一種 (才二種) の誤りの確率という。一般に有意水準 α の検定とは $\alpha(\varphi) \leq \alpha$ の検定で, その中で $\beta(\varphi)$ を最小にするものを Neyman-Pearson の検定基準を満たすという。検定 φ は, $x \in X$ が観測の結果得られたとき, 確率 $\varphi(x)$ で仮説 H_0 を棄却, 確率 $1 - \varphi(x)$ で採択することを意

味する。

n 回の観測を行なうということは、上の仮説検定を標本空間を $(X^{(n)}, \mathfrak{F}^{(n)})$, その上の確率分布を $\{P_\theta^{(n)}; \theta \in \Theta\}$

$$X^{(n)} = \prod_{i=1}^n X, \quad \mathfrak{F}^{(n)} = \otimes_{i=1}^n \mathfrak{F}$$

$$P_\theta^{(n)} = \otimes_{i=1}^n P_\theta$$

として考えることを意味する。

標本空間 $(X^{(n)}, \mathfrak{F}^{(n)})$ における検定の全体を Φ_n とする、すなわち、

$$\Phi_n = \{ \varphi_n; X^{(n)} \rightarrow [0, 1], \text{可測} \}.$$

また、 α ($0 < \alpha < 1$) に対して、

$$\beta_n(\alpha) = \inf \{ E_{\theta_1}(1 - \varphi_n); \varphi_n \in \Phi_n, E_{\theta_0}(\varphi_n) \leq \alpha \}$$

とする。Neyman-Pearsonの定理によつて、 $\varphi_n^* \in \Phi_n$ で

$$\beta_n(\alpha) = E_{\theta_1}(1 - \varphi_n^*), \quad E_{\theta_0}(\varphi_n^*) = \alpha$$

を満たすものが存在することを注意する。

定理1. 任意の α , $0 < \alpha < 1$ に対して、

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n(\alpha) = -I(\theta_0 | \theta_1).$$

ここで, $I(\theta_0 | \theta_1) = I(P_{\theta_0} | P_{\theta_1})$ である。

この定理は統計で Stein's Lemma として知られているものである。証明は Bahadur [1] が詳しい。これによって, 才一種の誤りの確率を一定にしたときの, 才二種の誤りの確率の限界が I-divergence によって与えられることが示されたことになる。

§ 2. Large deviations

$\{f_i; i=1, 2, \dots\}$ を確率空間 (X, \mathcal{A}, P) 上の独立な同一分布を持つ確率変数列とする。このとき, 確率

$$P\left\{\frac{1}{n}\sum_{i=1}^n f_i \geq a\right\}$$

に注目する。平均 $E(f_1)$ が存在するとき, 大数の法則により,

$$P\left\{\frac{1}{n}\sum_{i=1}^n f_i \geq a\right\} \rightarrow 0, \quad \forall a > E(f_1).$$

である。この収束の速さと I-divergence の関係をみる。

f を (X, \mathcal{A}, P) 上の確率変数とする。 f の m.g.f. (moment generating function) を

$$\varphi(t) = \int \exp(t f(x)) P(dx), \quad -\infty < t < +\infty.$$

で定義する。 φ は \mathbb{R} 上の convex function で $\{t: \varphi(t) < \infty\}$ は 0 を含む区間である。 また、 $P\{f=0\} < 1$ かつ、 $\varphi(b) < \infty$ ($\exists b > 0$) であれば、

- (i) φ : strictly convex, continuous on $[0, b]$
- (ii) $\varphi \in C^\infty(0, b)$
- (iii) φ' : strictly increasing on $(0, b)$
- (iv) $\varphi'(0+) = E(f)$, $\varphi'(b-) = E(f \cdot \exp bf)$.

以上はよく知られた m.g.f の性質である。

Chernoff [3] は m.g.f を用いて次を示した。

定理 2. $\{f_i; i=1, 2, \dots\}$ は確率空間 (X, \mathcal{F}, P) 上の独立同一分布列で、 f_i の m.g.f φ において、 $\varphi(t) < +\infty$ ($\forall t \in [0, b)$) とする。 このとき、

$$a \in \left\{ \frac{\varphi'(t)}{\varphi(t)} : t \in (0, b) \right\}$$

に対して

$$P\left\{ \frac{1}{n} \sum_{i=1}^n f_i \geq a \right\} \leq \rho(a)^n,$$

$$\lim_{n \rightarrow \infty} \left(P\left\{ \frac{1}{n} \sum_{i=1}^n f_i \geq a \right\} \right)^{\frac{1}{n}} = \rho(a)$$

ここで、

$$\rho(a) = \inf_t \varphi(t) \exp(-ta)$$

である。

上記の $\rho(a)$ と I-divergence を結びつけた次の結果は、
Samov [7], Hoeffding [6] による。

定理3. 確率分布の族 \mathcal{P} を

$$\mathcal{P} = \left\{ Q : \int f_i(x) Q(dx) \geq a \right\}$$

とすれば,

$$\rho(a) = \exp(-I(\mathcal{P} | P))$$

である。

さて、上記の二つをまとめると、

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log P \left\{ \frac{1}{n} \sum_{i=1}^n f_i \geq a \right\} = -I(\mathcal{P} | P)$$

となる。Groenboom, Oosterhoff and Ruymgaart [5] は、この型の定理について、より一般的な場合について調べている。

§3. Conditional probability and I-divergence.

前節までは標本の数が大きいときのある種の確率の収束の速さが、I-divergence で表わされることをみた。以下では、前の結果が条件付確率を持ち出すことにより、見通しがよくなることをみる。

(X, \mathfrak{A}, P) を確率空間とする。 $A \in \mathfrak{A}$, $P(A) \neq 0$

による条件付確率測度を

$$P_A(B) = P(B|A), \quad B \in \mathfrak{X}$$

と表記する。

条件付確率測度は I -divergence で特徴付けできる。
すなわち次の定理が成立する。

定理4. $A \in \mathfrak{X}$, $P(A) \neq 0$ に対して

$$\mathcal{P}_A = \{ Q; Q(A) = 1 \}$$

とすれば, P_A は \mathcal{P}_A の中で I -divergence を最小にする。

$$I(P_A | P) = I(\mathcal{P}_A | P).$$

これは, $I(P_A | P) = -\log P(A)$ に注意すれば, 容易に証明
することができる(文献[9])。

前節の定理を条件付確率を用いて書きかえる。ただし
簡単のため, 有界性の条件を確率変数に付加する。

定理5. (X, \mathfrak{X}, P) を確率空間, $\{x_i; i=1, 2, \dots\}$
を, その上の自分自身に値をとる確率変数列で, 独立同一分
布 ($P_{x_i^{-1}} = P$) とする。 f を X 上の実確率変数で

$$\text{ess. inf } f(x_i) < a < \text{ess. sup } f(x_i)$$

とし、以後 $f(x_i)$ を f_i と書く。

$$\mathcal{P} = \{Q : \int f_i(x) Q(dx) \geq a\}$$

$$P_n(\cdot) = P(\cdot \mid \frac{1}{n} \sum_{i=1}^n f_i \geq a)$$

とすれば、

$$\lim_{n \rightarrow \infty} \frac{1}{n} I(P_n | P) = I(\mathcal{P} | P)。$$

これは、 $I(P_n | P) = -\log P\{\frac{1}{n} \sum_{i=1}^n f_i \geq a\}$ に注意すれば前定理そのものであるが、これによって、 n 回の観測によって得られる情報量の平均値は、 n が十分大きいとき、 $I(\mathcal{P} | P)$ で近似されることがわかる。

ところで、我々は、 $P^* \in \mathcal{P}$ が存在して

$$I(P^* | P) = I(\mathcal{P} | P)$$

を示すことができる。では、 P_n と P^* の関係はどうなるのであろうか？

定理6. 任意の $A \in \mathcal{E}$ に対して

$$P_n(A) \rightarrow P^*(A) \quad (n \rightarrow \infty)。$$

証明の概略は、 (X, \mathcal{E}) 上の有界な実確率変数 g を持て来る。 $g(x_i)$ を g_i と書き、 \bar{g} を

$$\bar{q} = \int q(x) P^*(dx).$$

として, 任意の $\varepsilon > 0$ に対して

$$P_m\left(\left|\frac{1}{m}\sum_{i=1}^m g_i - \bar{q}\right| \geq \varepsilon \mid \frac{1}{m}\sum_{i=1}^m f_i \geq a\right) \rightarrow 0$$

を示す, これから,

$$\int q P_m(dx) \rightarrow \int q P^*(dx).$$

となり, 目的は達せられる。詳細は文献[9]にゆずる。

これにより, I-divergence を最小化することによって得られた確率測度 P^* は, 標本平均による条件付確率測度の極小を意味していることが分かる。

REFERENCES.

- [1] Bahadur, R.R.; 'Some Limit Theorems in Statistics', SIAM, Philadelphia (1971).
- [2] Bahadur, R.R. and Zabell, S.L.; Large deviations of the sample mean in general vector spaces, Ann. Probability 7 587-621(1979).
- [3] Chernoff, H.; A measure of asymptotic efficiency for tests of a hypothesis based on sums of observations, Ann. Math. Statist. 23, 493-507(1952).

- [4] Csiszar, I.; I-divergence geometry of probability distributions and minimization problems, Ann. Probability 3, 146-158(1975).
- [5] Groeneboom, R., Oosterhoff, J. and Ruymgaart, F.H.; Large deviation theorems for empirical probability measures, Ann. Probability 7, 553-586(1979).
- [6] Hoeffding, W.; On probability of large deviations, Proc. Fifth Berkeley Symp. Math. Statist. Prob. 1, 203-219 (1967).
- [7] Sanov, I.N.; On the probability of large deviations of random variables, Sel. Transl. Math. Statist. Prob. 1, 213-244(1967).
- [8] Sugawara, A.; A mathematical information channel with a non-commutative intermediate system (submitted).
- [9] Sugawara, A.; Conditional probability and I-divergence, (under preparation).