

適応的サンプリングによる標本推定法の設計

- 発見科学からのアプローチ -

荒木 純道

- 目次 -

- ◆ 問題の背景
- ◆ 数学的準備
- ◆ 一括サンプリングとその問題点
- ◆ 適応サンプリング
- ◆ まとめ

確率事象の発生確率の標本推定

$$P \cong \frac{\sum_{i=1}^n X_i}{n} \stackrel{\Delta}{=} \tilde{P}$$

P : 真の発生確率 (未知の定数) \tilde{P} : 推定発生確率 (確率変数)

ただし, X_i : 2値 (0, 1) の確率変数

事象が発生した時: $X_i = 1$

発生しなかった時: $X_i = 0$

問題点

- ◆ P が低い時にどれだけの標本数 n を用意しなければならないか？
- ◆ 推定値 \tilde{p} と真値 P との(絶対, 相対)誤差はどのくらいか？
- ◆ 推定値 \tilde{p} の信頼性は？

発見科学(機械学習, データマイニング)の 立場と目標

◆ 確率事象の発生と確認は容易で, サンプル数は必要なだけ十分得られる.

cf. 古典的統計推定理論とは前提が少し違う.

◆ 推定結果が, 指定された信頼度と相対誤差の規範を満足するアルゴリズムを設計提供する.

標本推定の設計規範

$$\Pr\left[|\tilde{P} - P| < \varepsilon P\right] < 1 - \delta \quad (1)$$

信頼度: $1 - \delta$ ($0 < \delta < 1$)

相対誤差: ε ($0 < \varepsilon < 1$)

$\delta \rightarrow 0, \varepsilon \rightarrow 0$ が望ましい

数学的準備 (Chernoff 限界)

$\bar{X}_i = P$ の時, 任意の $\varepsilon (0 < \varepsilon < 1)$ に対して

$$\Pr[\tilde{P} > (1 + \varepsilon)P] < \exp(-nP\varepsilon^2/3) \quad (2a)$$

$$\Pr[\tilde{P} < (1 - \varepsilon)P] < \exp(-nP\varepsilon^2/2) \quad (2b)$$

$$\text{ただし, } \tilde{P} = \sum_{i=1}^n X_i / n$$

Chernoffのアイデア: Min - Max原理

$$\Pr[X > C] = \Pr[f(X:t) > f(C:t)]$$

X : 乱数 C : 定数

ただし, $f(x:t)$: パラメタ t を含む x に関する単調増加関

例 $f(x:t) = e^{tx} \quad (t > 0)$

次に, 確率上限値をパラメタ t に関して最小化をはかり, 上界評価をより正確にする.

Markovの不等式

$$\Pr[X > C\bar{X}] \leq \frac{1}{C}$$

ただし, X : 正值乱数 \bar{X} : X の平均値 C : 正の定数

$$\begin{aligned}\Pr[X > C\bar{X}] &= \int_{C\bar{X}}^{\infty} P(X) dX \leq \frac{1}{C\bar{X}} \int_{C\bar{X}}^{\infty} XP(X) dX \\ &\leq \frac{1}{C\bar{X}} \int_0^{\infty} XP(X) dX = \frac{1}{C}\end{aligned}$$

一括サンプリング(固定サンプル数)とその問題点

規範(1)を満足するためには(2)から

$$\exp(-nP\varepsilon^2/3) < \frac{\delta}{2}$$
$$\therefore n > \frac{3}{P\varepsilon^2} \ln\left(\frac{2}{\delta}\right) \quad (3)$$

のサンプル数があればよい。

しかし、 ε , δ は指定できるが、 P は未知であるので(3)からでは必要なサンプル数 n は決まらない。

適応サンプリングの原理

- ◆ nP ほぼ事象の「発生回数」 $\frac{3}{\varepsilon^2} \ln\left(\frac{2}{\delta}\right)$
- ◆ サンプル n 数を事前に決めるのではなく、事象の発生回数の上限を指定すればよい。

While 文構造にする。(くり返し回数は確率変数)

cf. 一括サンプリングは For 文構造

(くり返し回数は確定数)

[定理 1]

アルゴリズム A の出力結果 \tilde{p} は, (1) の規範を必ず満足

begin

$m \leftarrow 0; \quad n \leftarrow 0;$

ただし,

while $m < a$ do

$m \leftarrow m + x_i; \quad n \leftarrow n + 1;$

$$a = \frac{3(1 + \varepsilon)}{\varepsilon^2} \ln \left(\frac{2}{\delta} \right)$$

$\tilde{p} = m/n$

end

[定理 2]

アルゴリズム A のくり返し回数 n は

$$\Pr \left[n \leq \frac{3(1+\varepsilon)}{(1-\varepsilon)p\varepsilon^2} \ln \left(\frac{2}{\delta} \right) \right] > 1 - \frac{\delta}{2}$$

を満足する $(1+\varepsilon)/(1-\varepsilon)$ 倍だけ最適値より多い。

なお, P が既知の時, 必要くり返し回数は

$$n = \frac{3}{p\varepsilon^2} \ln \left(\frac{2}{\delta} \right)$$

数値例

$$\delta = 0.01 \quad \varepsilon = 0.1$$

$$a = \frac{3(1 + \varepsilon)}{\varepsilon^2} \ln\left(\frac{2}{\delta}\right) = 1747.9$$

◆ a の計算は容易

◆ $\delta \rightarrow 0, \quad \varepsilon \rightarrow 0, \quad a \rightarrow \infty$

◆ $a \propto \varepsilon^{-2}, \quad a \propto \ln\left(\frac{1}{\delta}\right)$

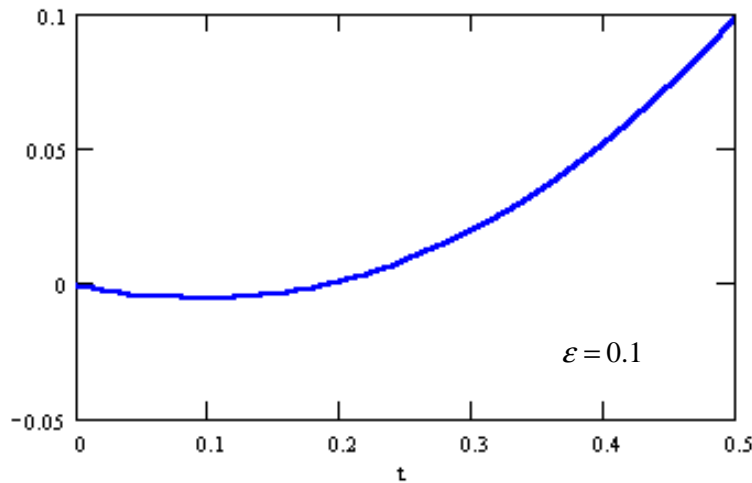
まとめ

- ◆ (δ, ε) 設計規範を満足する適応サンプリング
アルゴリズムを示した.
- ◆ 必要な事業発生回数は, 相対誤差 ε , 信頼度 δ に対して
 $\varepsilon^{-2}, \ln\left(\frac{1}{\delta}\right)$ に比例.
- ◆ アルゴリズムはwhile文構造で容易に記述.
- ◆ Chernoff限界ではなく, 中心極限定理を用いれば, もっと効率的なアルゴリズムが設計できる. (2 ~ 3倍早い)
- ◆ さあ, みんなでアルゴリズム A を使ってみて下さい!

附 録

$$\begin{aligned}\Pr[\tilde{P} > (1 + \varepsilon)P] &= \Pr\left[\sum_{i=1}^n X_i > (1 + \varepsilon)nP\right] \\ &= \Pr\left[e^t \sum_{i=1}^n X_i > e^{(1+\varepsilon)nPt}\right] \quad \forall t > 0 \text{ に対して} \\ &< \frac{E\left[e^t \sum_{i=1}^n X_i\right]}{e^{(1+\varepsilon)nPt}} \quad (\rightarrow \text{Markov不等式}) \\ &= \frac{\{E[e^{tX_i}]\}^n}{e^{(1+\varepsilon)nPt}} \quad (\rightarrow \{X_i\} \text{のiid性}) \\ &= \frac{(e^t P + 1 - P)^n}{e^{(1+\varepsilon)nPt}} \quad (\rightarrow e^x \geq 1 + x) \\ &< \frac{e^{nP(e^t - 1)}}{e^{nP(1+\varepsilon)t}} \\ &= e^{nP} [e^t - 1 - (1 + \varepsilon)t]\end{aligned}$$

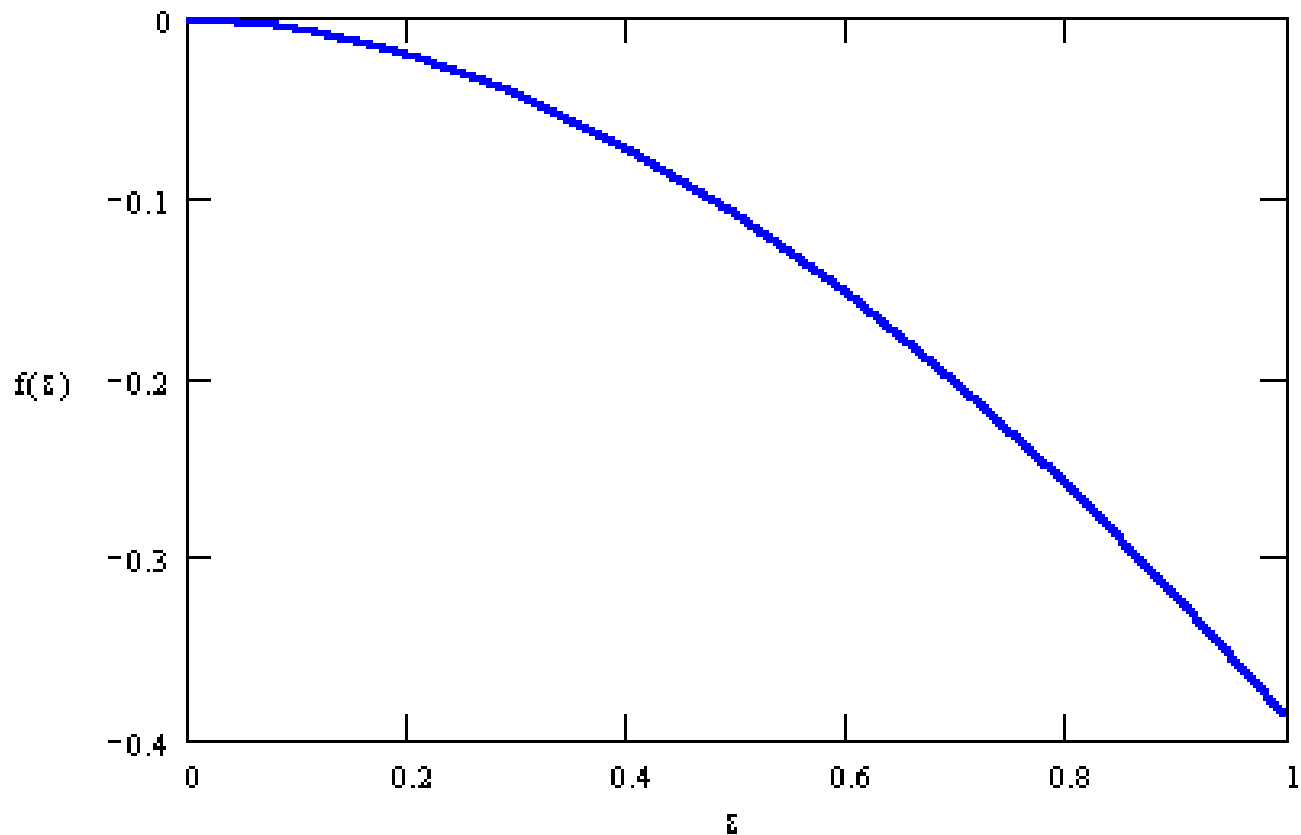
$f(t) = e^t - 1 - (1 + \varepsilon)t$ のグラフ



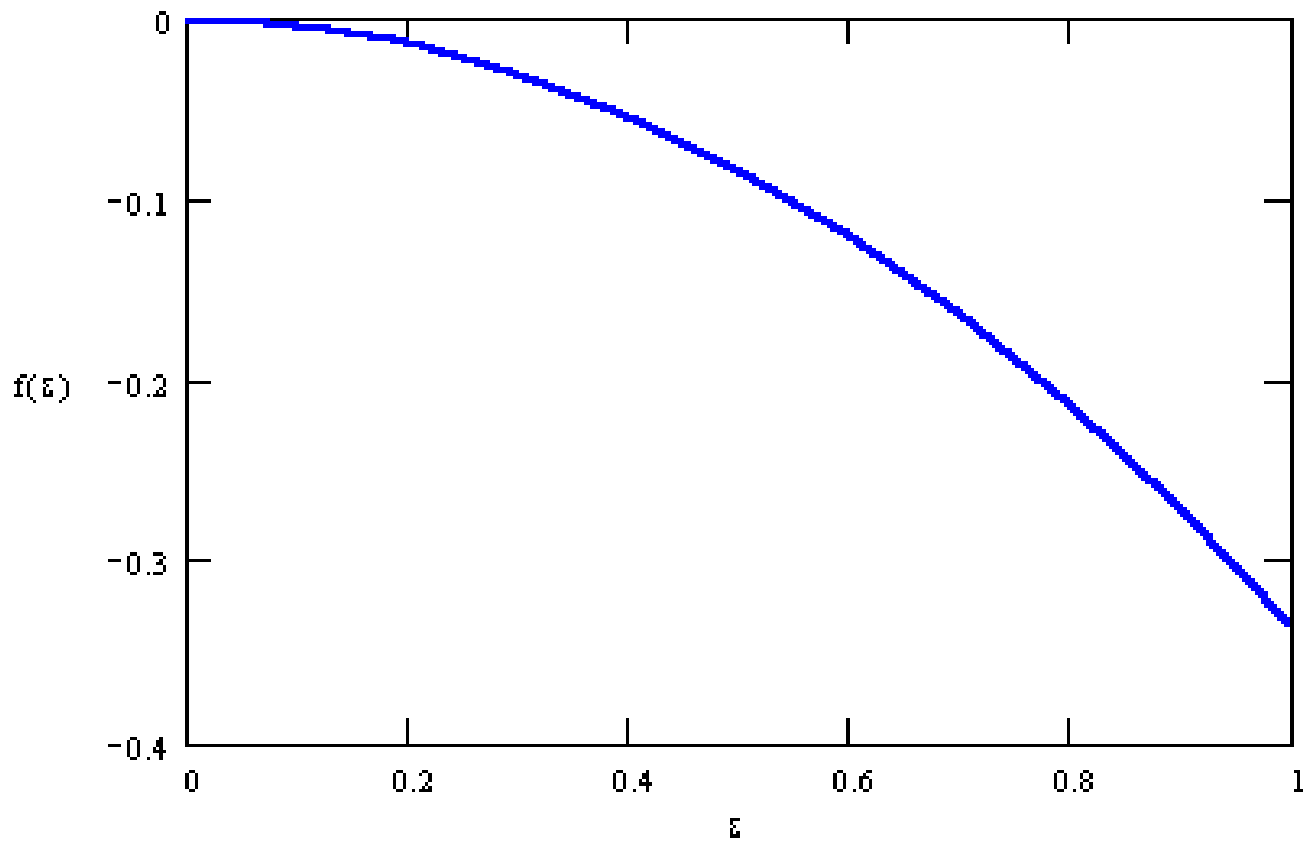
$$\begin{aligned} \text{Min } f'(t) &= e^t - (1 + \varepsilon) = \\ &\therefore t_{\min} = \ln(1 + \varepsilon) \end{aligned}$$

$$\begin{aligned} \text{Min } f(t_{\min}) &= (1 + \varepsilon) - 1 - (1 + \varepsilon)\ln(1 + \varepsilon) \\ &= \varepsilon - (1 + \varepsilon)\ln(1 + \varepsilon) \end{aligned}$$

$f(\varepsilon) = \varepsilon - (1 + \varepsilon)\ln(1 + \varepsilon)$ のグラフ



$f(\varepsilon) = -\varepsilon^2/3$ のグラフ



$$\therefore \varepsilon - (1 + \varepsilon) \ln(1 + \varepsilon) < -\varepsilon^2 / 3$$

$$\therefore \Pr[\tilde{P} > (1 + \varepsilon)P] < e^{-nP\varepsilon^2/3}$$

同様に $-\varepsilon - (1 - \varepsilon) \ln(1 - \varepsilon) < -\varepsilon^2 / 2$ だから

$$\Pr[\tilde{P} < (1 - \varepsilon)P] < e^{-nP\varepsilon^2/2}$$