I スコア関数とFisher情報量

abstract 情報幾何学の基礎を支えるスコア関数とFisher情報量を復習します。

1 Introduction

情報幾何学は、確率分布全体の集合を一つの空間(多様体)と見なし、そこに幾何学的な構造を与えて解析する学問分野です。

私たちが馴染んでいるユークリッド空間では、点と点の間に「距離」を定義したり、曲線の「長さ」を 測ったりすることができます。情報幾何学は、この考え方を統計モデルの世界に持ち込みます。

- 空間の「点」:一つ一つの確率分布(例:平均0,分散1の正規分布)
- 空間の「座標」:分布を特徴づけるパラメータ(例:正規分布の平均 μ と分散 σ^2)

この「統計多様体」と呼ばれる空間に、距離や角度、曲率といった幾何学的な"ものさし"を導入することで、これまでとは異なる視点から統計モデルの性質を深く理解することを目指します。

実はこのような視点は、すでに数理統計学を学んだことがある方なら、暗に使って来た考え方です。例えば、スコア関数、Fisher情報量そしてKullback-Leibler divergenceを使ったことがあれば、皆さんは暗に情報幾何学的な考え方に触れたことがあるでしょう。そこでこの講義は、これらの概念を情報幾何学を通して理解することを目指します。

2 各概念の基本的な性質

2.1 スコア関数

対数尤度関数の導関数をスコア関数といいます。

定義 スコア関数

パラメータ $\theta=(\theta_1,\cdots,\theta_n)$ を持つ確率分布 $p(x;\theta)$ に対して、そのスコア関数とは以下のように定義される。

$$S(heta) = rac{\partial \log p(x; heta)}{\partial heta}$$

なぜこのような量が重要なのか、現在は理解が難しいと思いますが、情報幾何学を学んでいくうちにその重要性がよくわかることでしょう。一旦、具体例を与えておきます。

例 正規分布のスコア関数

正規分布 $N(\mu, \sigma^2)$ のスコア関数は $S(\mu), S(\sigma^2)$ はそれぞれ以下のようになる。

$$S(\mu)=rac{x-\mu}{\sigma^2},\quad S(\sigma^2)=rac{(x-\mu)^2}{2\sigma^4}-rac{1}{2\sigma^2}$$

2.2 スコア関数の性質

スコア関数の重要な性質として、スコア関数の期待値は 0 になるというものがあります。

定理 スコア関数の期待値は0

以下の性質が成り立ちます。

$$\mathbb{E}\left[rac{\partial \log p(x; heta)}{\partial heta}
ight]=0$$

証明 対数微分から、左辺は

$$\mathbb{E}\left[rac{\partial \log p(x; heta)}{\partial heta}
ight] = \int p(x; heta) rac{rac{\partial p(x; heta)}{\partial heta}}{p(x; heta)} dx = \int rac{\partial p(x; heta)}{\partial heta} dx$$

と式変形できます。さらに、微分と積分の順序交換が可能だとすると、

$$\int rac{\partial p(x; heta)}{\partial heta} dx = rac{\partial}{\partial heta} \int p(x; heta) dx = rac{\partial 1}{\partial heta} = 0$$

です。■

2.3 Fisher情報量

Fisher情報量は各パラメータのスコア関数の内積として定義されます。この内積として定義されているという事実は非常に重要で、ここにFisher情報量と情報幾何学のつながりが垣間見えています。

定義 Fisher情報量

パラメータ $heta=(heta_1,\cdots, heta_n)$ を持つ確率分布 p(x; heta) に対して、Fisher情報量 I(heta) の (i,j) 成分を

$$I_{ij}(heta) = \mathbb{E}igg[rac{\partial \log p(x; heta)}{\partial heta_i} rac{\partial \log p(x; heta)}{\partial heta_j}igg]$$

で定義します。これはスコア関数の分散共分散行列に他なりません。

具体例として、正規分布のFisher情報量を求めておきます。

例 正規分布のFisher情報量

 $X \sim N(\mu, \sigma^2)$ とする。X が持つパラメータ (μ, σ^2) のFisher情報量は以下のように与えられる。

$$I(\mu,\sigma^2) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 1/2\sigma^4 \end{pmatrix}$$

2.4 Fisher情報量の性質

Fisher情報量が確率分布にとって本質的かどうかは、以下の定理を知るとひとまずの了解を得ることができるのではないでしょうか。

定理 Fisher情報量は全単射な変数変換に対して不変

確率変数 X がパラメータ θ によって定義される確率分布 $p(x;\theta)$ に従っているとする。変換 Y=f(X) が逆 関数 $g=f^{-1}$ を持つとき、X が持つFisher情報量 $I_X(\theta)$ と Y が持つFisher情報量 $I_Y(\theta)$ は等しい。

証明 変数変換の公式から以下のことがわかる点に注意します。

$$p(y; heta) = p(x(y); heta)rac{dx}{dy}, \quad \log p(y; heta) = \log p(x(y); heta) + \lograc{dx}{dy}$$

特に、dx/dy は θ に依存しないので、

$$rac{\partial \log p(y; heta)}{\partial heta} = rac{\partial \log p(x(y); heta)}{\partial heta}$$

です。X が持つFisher情報量を置換積分で式変形すると、以下のようになります。

$$egin{aligned} I_X(heta) &= \int p(x; heta) rac{\partial \log p(x; heta)}{\partial heta} dx \ &= \int p(x(y); heta) rac{\partial \log p(x(y); heta)}{\partial heta} rac{dx}{dy} dy \ &= \int p(y; heta) rac{\partial \log p(y; heta)}{\partial heta} dy \ &= I_Y(heta) \end{aligned}$$

以上で定理の主張が確認できました。■

3 Fisher情報量の重要性

Fisher情報量は、統計的推定の理論において中心的な役割を果たします。その代表例がCramer-Raoの下限です。もう一つの重要な内容としてLehmann-Scheffeの定理がありますが、こちらは後ほど詳しく解説します。

3.1 Cramer-Raoの下限

Cramer-Raoの下限は、不偏推定量が持ちうる分散の下限を与えます。

定理 Cramer-Raoの下限

パラメータ heta の任意の不偏推定量 $\hat{ heta}$ (すなわち $\mathbb{E}[\hat{ heta}]= heta$ を満たす推定量) の分散は次を満たします。

$$\mathbb{V}[\hat{ heta}] \geq rac{1}{I(heta)}$$

ここで $I(\theta)$ はパラメータ θ に関するFisher情報量です。

この不等式は、Fisher情報量が大きいほど推定量の分散が小さくなる(=より高精度な推定が可能になる)ことを示します。

3.2 Cramer-Raoの下限の証明

不偏推定量 $\hat{\theta}$ をスコア関数との関係で特徴づけてみましょう。不偏性の定義は

$$\mathbb{E}[\hat{ heta}] = \int \hat{ heta}(x) p(x; heta) \, dx = heta$$

でした。ここで両辺を θ で微分してみましょう。微分と積分の順序交換が可能だとすると、

$$\int \hat{\theta}(x) \frac{\partial p(x;\theta)}{\partial \theta} \, dx = 1$$

が得られます。ここで、 $\dfrac{\partial p}{\partial heta} = p \dfrac{\partial \log p}{\partial heta}$ を用いると

$$\mathbb{E}\left[\hat{\theta}\frac{\partial \log p}{\partial \theta}\right] = 1$$

が得られ、これが不偏推定量とスコア関数との間の関係です。特に、スコア関数の期待値が0であることを思い出せば、

$$Covig(\hat{ heta}, rac{\partial \log p}{\partial heta}ig) = \mathbb{E}ig[\hat{ heta} rac{\partial \log p}{\partial heta}igg] - \mathbb{E}[\hat{ heta}]\mathbb{E}igg[rac{\partial \log p}{\partial heta}igg] = 1$$

が成り立ちます。

あとは、コーシー・シュワルツの不等式 $Cov(X,Y)^2 \leq \mathbb{V}[X]\mathbb{V}[Y]$ を適用すると

$$1^2 \leq \mathbb{V}[\hat{ heta}] \, \mathbb{V} igg[rac{\partial \log p}{\partial heta} igg]$$

が得られ、
$$\mathbb{V}\left[rac{\partial \log p}{\partial heta}
ight] = \mathbb{E}\left[\left(rac{\partial \log p}{\partial heta}
ight)^2
ight] = I(heta)$$
 より

$$\mathbb{V}(\hat{ heta}) \geq rac{1}{I(heta)}$$

が従います。これはCramer-Raoの下限です。