

Chapter 1 非参数统计

Statistical Inference 的步骤

① 建立模型 $\{P_f : f \in F\}$ 参数 f 在参数空间 F 中。

用估计量 $T(Y)$ 来估计 f , 其中 Y 是我们的观测值

② 建立 test function $\Psi(Y)$ 以对提出的 f 进行假设检验

③ 建立关于参数 f 的置信区间 (confidence sets)

在本书中, 主要考虑参数空间 F 是无穷维的情况。

1.1 Statistical Sampling Models.

1.1.1

我们定义衡量两个 probability 距离远近的 metric. 其中 base probability space 为 (X, A)

① total variational metric

$$\|P - Q\|_{TV} = \sup_{A \in \mathcal{A}} |P(A) - Q(A)|$$

② bounded Lipschitz metric:

$$B_{(X, d)}(P, Q) = \sup_{f \in BL(1)} \left| \int_X f(dP - dQ) \right|, \text{ 其中 } BL(1) = \{f : \sup_{x \in X} |f(x)| + \sup_{x \neq y} \frac{|f(x) - f(y)|}{d(x, y)} \leq 1\}$$

③ kolmogorov distance:

$$\|F_p - F_q\|_\infty = \sup_{x \in \mathbb{R}} |F_p(x) - F_q(x)|. \text{ 其中 } F(x) = P(X \leq x).$$

④

$$\|f_p - f_q\|_1 = \int_R |f_p(x) - f_q(x)| dx$$

1.1.2 Indirect observations.

事实上能观测到的都是在原始数据上添加了扰动。

例如: $Y_i = X_i + \varepsilon_i$, ε_i 为观测误差, Y_i 为观测结果。

$$P_Y = P_X * P_\varepsilon, * \text{ 代表卷积}$$

1.2. Gaussian Models

1.2.2 Some nonparametric Gaussian Models

The Gaussian White Noise Model

$$[1.5] \quad dY(t) = dY_f^{(n)}(t) = f(t)dt + \frac{\sigma}{\sqrt{n}} dW(t), \quad t \in [0,1], \quad f \in L^2([0,1]).$$

其中 σ 为 dispersion parameter

dW : standard Gaussian white noise process.

上面式子可说成 We observe the function f in Gaussian model, at the noise level $\frac{\sigma}{\sqrt{n}}$.

在这里, 我们本该将 dW 视为 standard Brownian motion $\{W(t) : t \in [0,1]\}$ 的导数, 但是事实上可以这么考虑

$$[1.6] \quad g \mapsto \int_0^1 g(t) dY_f^{(n)}(t) = Y_f^{(n)}(g) \sim N(\langle f, g \rangle, \frac{\|g\|_2^2}{n}), \quad g \in L^2([0,1])$$

$$[1.7] \quad g \mapsto \int_0^1 g(t) dW(t) = W(g) \sim N(0, \|g\|_2^2), \quad g \in L^2([0,1])$$

若在 (1.7) 中, 取 g 为 L^2 中的有限个标准正交基 (e_k)

$$e_k \mapsto \int_0^1 e_k(t) dW(t) = Y_f^{(n)}(e_k) \sim N(\langle f, e_k \rangle, \frac{1}{n})$$

回顾 Kolmogorov consistency theorem, 知道

[For T indexed set. $s, t \in T$ 若 \exists Gaussian process $X(t)$, 使得

$$E(X(t)) = f(t), \quad E[(X(t) - f(t))(X(s) - f(s))] = \Phi(s, t), \quad \text{则 } X$$

the covariance and of this process and with f , its expectation.]

那么现在: $T = L^2([0,1])$ $Y(g) \sim N(\langle f, g \rangle, \frac{\|g\|_2^2}{n})$, $W(g) \sim N(0, \|g\|_2^2)$

$$\begin{aligned} E(Y(t) - f(t))(X(s) - f(s)) &= E(Y(t) - f(t))E(X(s) - f(s)) \\ &= 0 \end{aligned}$$

那我们该如何理解 model (1.5) 呢?

$$dY(t) = dY_f^{(n)}(t) = f(t)dt + \frac{\sigma}{\sqrt{n}} dW(t), \quad t \in [0, 1]$$

也就是说, 给一个 g , 我们就有一个 $N(\langle f, g \rangle, \frac{\|g\|^2}{n})$ 的正态分布. However, the natural sample space now is hard to work with.

Gaussian Sequence Space Model

回顾我们的 Gaussian process $\{Y_f^{(n)}(g) : g \in L^2\}$ just means, we observe $Y_f^{(n)}(g)$ for all $g \in L^2$ simultaneous. Now take $\{e_k : k \in \mathbb{Z}\}$ is orthonormal basis of L^2 . Y

$$(1.8) \quad Y_k = Y_{f,k}^{(n)} = \langle f, e_k \rangle + \frac{\sigma}{\sqrt{n}} g_k, \quad k \in \mathbb{Z}, n \in \mathbb{N}. \text{ 被称为 GSSM.}$$

$$g_k: \text{i.i.d r.v of law } W(e_k) \sim N(0, \|e_k\|_2^2) = N(0, 1)$$

(1.5) and (1.8) are observationally equivalent to each other.

(1.8) 的 special form: $Y_k = \theta_k + \frac{\sigma}{\sqrt{n}} g_k, k=1, 2, \dots, n.$

(1.8) 的进一步说明:

$\{e_k : k \in \mathbb{Z}\}$ is a sequence space isometry from L^2 to the sequence space ℓ_2 of all square-summable infinite sequence through the mapping $f \mapsto \langle f, e_k \rangle$. the law $\{Y_{f,k}^{(n)} : k \in \mathbb{Z}\}$ completely characterise the finite-dimensional distributions, and thus the law, of the process $Y_f^{(n)}$.

Q: 为何说 有限维分布 determines the law of process $Y_f^{(n)}$?

1.2.3 Equivalence of Statistical Experiences

The Le Cam Distance of Statistical Experiences.

记统计实验为 $\mathcal{E}^{(i)}$, 其中 $\mathcal{E}^{(i)}$ 由样本空间 \mathcal{Y}_i 与 \mathcal{Y}_i 上的测度 $P_f^{(i)}$ 构成.

定义损失函数 (用于观测 the performance of a decision procedure $T^{(i)}(\mathcal{Y}^{(i)}) \in \mathcal{T}$ based on observations $\mathcal{Y}^{(i)}$)

$L : \mathcal{F} \times \mathcal{T} \longrightarrow [0, \infty)$, 其中 \mathcal{T} 代表由决策规则所组成的集合.

$$f, T^{(i)}(\mathcal{Y}^{(i)}) \mapsto L(f, T^{(i)}(\mathcal{Y}^{(i)}))$$

举例来说, 如若 $T^{(i)}(\mathcal{Y}^{(i)})$ 表示对参数 f 本身的估计量, 此时 $\mathcal{F} = \mathcal{T}$. $L(f, T) = d(f, T)$, d 为参数空间上定义的某种距离.

在概率 $P_f^{(i)}$ 下, 记录此损失函数的数学期望为 $R^{(i)}(f, T^{(i)}, L) = E_{P_f^{(i)}}[L(f, T^{(i)}(\mathcal{Y}^{(i)}))]$

"Le Cam distance" (between 2 experiments)

$$\Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) = \max_f \left[\sup_{T^{(2)}} \inf_{T^{(1)}} \sup_{f, L: |L|=1} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, T^{(2)}, L)| \right],$$

$$\sup_{T^{(2)}} \inf_{T^{(1)}} \sup_{f, L: |L|=1} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, T^{(2)}, L)|.$$

其中 $|L| = \sup \{L(f, T) : f \in \mathcal{F}, T \in \mathcal{T}\}$

如若 $\Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) = 0$, 说明任意决策规则 $T^{(2)}$ in experiment $\mathcal{E}^{(1)}$ 都可被转换为实验 $\mathcal{E}^{(2)}$ 下的决策规则 $T^{(1)}$.

• propositions

- (1) $\mathcal{Y}^{(1)} = \mathcal{Y}^{(2)} = \mathcal{Y}$ (为完备可分度量空间时), 当 $P_f^{(1)}, P_f^{(2)}$ 有共同的 dominating measure μ on \mathcal{Y} 时, 有 $\Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) \leq \sup_{f \in \mathcal{F}} \int_{\mathcal{Y}} \left| \frac{dP_f^{(1)}}{du} - \frac{dP_f^{(2)}}{du} \right| du = \|P^{(1)} - P^{(2)}\|_{1, \mu, F}$

这里由于:

$$\Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) = \max_{\mathbb{F}} \left[\sup_{f \in \mathbb{F}} \inf_{T^{(1)} \in \mathcal{T}^{(1)}} \sup_{L \in \mathcal{L}} |R^{(1)}(f; T^{(1)}, L) - R^{(2)}(f; T^{(2)}, L)|, \right.$$

$$\left. \sup_{T^{(1)}} \inf_{T^{(2)} \in \mathcal{T}^{(2)}} \sup_{f \in \mathbb{F}, L \in \mathcal{L}} |R^{(1)}(f; T^{(1)}, L) - R^{(2)}(f; T^{(2)}, L)| \right]$$

而 $\sup_{T^{(2)}} \inf_{T^{(1)} \in \mathcal{T}^{(1)}} \sup_{f \in \mathbb{F}, L \in \mathcal{L}} |R^{(1)}(f; T^{(1)}, L) - R^{(2)}(f; T^{(2)}, L)|$

$$\leq \sup_{T^{(2)}} \sup_{f \in \mathbb{F}, L \in \mathcal{L}} |R^{(1)}(f; T^{(2)}, L) - R^{(2)}(f; T^{(2)}, L)|$$

$$= \sup_{T^{(2)}} \sup_{f \in \mathbb{F}, L \in \mathcal{L}} \left| \int_Y L(f, T^{(2)}(Y^{(2)})) dP_f^{(1)}(Y) - \int_Y L(f, T^{(2)}(Y^{(2)})) dP_f^{(2)}(Y) \right|$$

$$\leq \int_Y |L(f, T)| |dP_f^{(1)}(Y) - dP_f^{(2)}(Y)| \leq \|P^{(1)} - P^{(2)}\|_{1, u, F}.$$

以上我们须要找 $y^{(1)} = y^{(2)} = y$, 但 $y^{(1)} \neq y^{(2)}$ 时怎么办呢?

- two experiments are defined on different sample space.

构建同构 from $\mathcal{Y}^{(1)}$ to $\mathcal{Y}^{(2)}$, independent of f .

$$\text{使 } P_f^{(1)} = P_f^{(2)} \circ B, \quad \forall f \in \mathbb{F}$$

因此, Given $Y^{(2)} \in \mathcal{Y}^{(2)}$, $T^{(2)}(Y^{(2)}) = T^{(1)} \circ B^{-1}(Y^{(2)})$ in $\mathcal{E}^{(1)}$

此时 $\Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) = \Delta_F(\mathcal{E}^{(1)}, B^{-1}(\mathcal{E}^{(2)})) = 0$

- In absence of such a bijection:

运用充分统计量.

$\mathcal{Y}^{(1)}$ 为样本空间, $\mathcal{E}^{(1)}$ 为 experiment give rise to observations $Y^{(1)}$ in law $P_f^{(1)}$ on $\mathcal{Y}^{(1)}$. 设 $S: \mathcal{Y}^{(1)} \rightarrow \mathcal{Y}^{(2)}$ 使

$$Y^{(2)} = S(Y^{(1)}), \quad Y^{(2)} \sim P_f^{(2)} \text{ on } \mathcal{Y}^{(2)}$$

且 $S(Y^{(1)})$ 为 $Y^{(1)}$ 的充分统计量.

$$\text{则 } \Delta_F(\mathcal{E}^{(1)}, \mathcal{E}^{(2)}) = 0.$$

$$\text{证明: } \Delta_F(\xi^{(1)}, \xi^{(2)}) = \max \left[\sup_{T^{(2)}} \inf_{T^{(1)}} \sup_{f, L} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, T^{(2)}, L)|, \right.$$

$$\left. \sup_{T^{(1)}} \inf_{T^{(2)}} \sup_{f, L} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, T^{(2)}, L)| \right]$$

为此我们要先回顾一下充分统计量之定义:

直观上: 给定统计量的值 $T=t$, 它所对应的分布 $F_\theta(X|T=t)$ 是一个与参数 θ 无关的分布, 也就是说, 只要 T 值确定, 不论参数 θ 选择什么, 样本的概率分布都是一成不变的.

DEF (Sufficient statistics)

A statistic $t=T(x)$ is sufficient for underlying parameter θ . if the conditional prob distribution of data X , given the statistic $t=T(x)$, doesn't depend on the parameter θ .

Characteristic: ① The data processing inequality: $I(\theta; T(x)) = I(\theta; X)$

② Fisher-Neyman factorization theorem:

$f_\theta(x)$: prob density function, T is sufficient for $\theta \Leftrightarrow \exists g, h \geq 0$,

$$f_\theta(x) = h(x) g_\theta(T(x))$$

现在,

$$\sup_{T^{(2)}} \inf_{T^{(1)}} \sup_{f, L} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, T^{(2)}, L)|$$

$$\sup_{T^{(1)}} \sup_{f, L} |R^{(1)}(f, T^{(1)}, L) - R^{(2)}(f, S(Y^{(1)}, L))|$$

$$= \left| \int_{Y^{(1)}} L(f, T^{(1)}(Y^{(1)})) dP_f^{(1)}(Y) - \int_{Y^{(2)}} L(f, T^{(2)}(Y^{(2)})) dP_f^{(2)}(Y) \right|$$

$$\leq \left| \int_{Y^{(1)}} L(f, T^{(1)}(Y^{(1)})) dP_f^{(1)}(Y) - \int_{Y^{(2)}} L(f, T^{(1)}(S(Y^{(1)}))) g(S(Y^{(1)}, f)) h(Y^{(1)}) dy^{(1)} \right|$$

$$\leq \left| \int_{Y^{(1)}} L(f, T^{(1)}(Y^{(1)})) dP_f^{(1)}(Y) - \int_{Y^{(1)}} L(f, \tilde{T}^{(1)}(Y^{(1)})) g(S(Y^{(1)}, f)) h(Y^{(1)}) S(Y^{(1)}) dy^{(1)} \right|$$

$$P_f^{(1)}(Y) = f(Y^{(1)}) dy^{(1)}$$

= 0 (暂未想到证明)

Asymptotic Equivalence for Nonparametric Gaussian Regression Models

这里首先定义 $F(a, M) = \{f: [0, 1] \rightarrow \mathbb{R}, \sup_{x \in [0, 1]} |f(x)| + \sup_{0 < |x-y| \leq 1} \frac{|f(x) - f(y)|}{|x-y|^a} \leq M\}$

Theorem. (三个模型之等价性)

$$\mathcal{E}_n^{(1)}: Y_i = f(x_i) + \varepsilon_i, x_i = \frac{i}{n}, \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$\mathcal{E}_n^{(2)}: dY(t) = f(t)dt + \frac{\sigma}{\sqrt{n}} dW_t, t \in [0, 1], n \in \mathbb{N}$$

$$\mathcal{E}_n^{(3)}: Y_k = \langle f, e_k \rangle + \frac{\sigma}{\sqrt{n}} g_k, g_k \sim N(0, \|e_k\|_2^2) = N(0, 1)$$

(1) for F any family of bdd. functions on $[0, 1]$,

$$\Delta_F(\mathcal{E}_n^{(2)}, \mathcal{E}_n^{(3)}) = 0, \quad \Delta_F(\mathcal{E}_n^{(1)}, \mathcal{E}_n^{(2)}) \leq \sqrt{\frac{n\sigma^2}{2}} \sup_{f \in F} \|f - \pi_n(f)\|_2.$$

(2). $\mathcal{F} = F(a, M)$, 其中 $a > \frac{1}{2}$, $M > 0$, 则以上三个实验是在 Le Cam 距离的意义下渐近等价的.