

MATHEMATICAL TRIPOS Part III

Thursday, 31 May, 2018 1:30 pm to 3:30 pm

PAPER 210

TOPICS IN STATISTICAL THEORY

*Attempt no more than **THREE** questions.*

*There are **FOUR** questions in total.*

The questions carry equal weight.

STATIONERY REQUIREMENTS

Cover sheet

Treasury Tag

Script paper

SPECIAL REQUIREMENTS

None

<p>You may not start to read the questions printed on the subsequent pages until instructed to do so by the Invigilator.</p>

1

Let X_1, \dots, X_n be independent and identically distributed real-valued random variables with density f . Given a point $x \in \mathbb{R}$ define the k -nearest neighbour distance of x , denoted $\rho_{(k)}(x)$.

For $r \geq 0$ and $x \in \mathbb{R}$ write $p_x(r) = \int_{x-r}^{x+r} f(y) dy$. Show that the random variable defined by $P = p_x(\rho_{(k)}(x))$ has Beta density

$$B_{k, n+1-k}(s) = \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n+1-k)} s^{k-1}(1-s)^{n-k}$$

for $s \in (0, 1)$. Calculate $\mathbb{E}P$ and $\mathbb{E}P^2$.

Henceforth suppose that f is L -Lipschitz and strictly positive on all of \mathbb{R} . Prove that $|p_x(r) - 2rf(x)| \leq Lr^2$ and hence, writing p_x^{-1} for the inverse of p_x , verify that

$$|2f(x)p_x^{-1}(s) - s| \leq \frac{Ls^2}{f(x)^2}$$

for any $s \in (0, 1)$ and x such that $f(x) \geq L^{1/2}$.

Write $\hat{f}_{(k)}(x) = \frac{k}{2(n+1)\rho_{(k)}(x)}$ for the k -nearest neighbour density estimator at x . Prove that

$$\left| \mathbb{E} \left(\frac{f(x)}{\hat{f}_{(k)}(x)} \right) - 1 \right| \leq \frac{k+1}{n+2}$$

for any x such that $f(x) \geq L^{1/2}$.

2

Consider the fixed-design nonparametric regression model in which we observe $Y_i = m(x_i) + \sigma\epsilon_i$ for $i = 1, \dots, n$, where $x_i = i/n$, $\sigma \in (0, \infty)$ and $\epsilon_1, \dots, \epsilon_n \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$. For a kernel $K : \mathbb{R} \rightarrow [0, \infty)$ and bandwidth $h \in (0, \infty)$ give the definition of the *local polynomial estimator* $\hat{m}_h(x; p)$ of $m(x)$, and derive the Nadaraya–Watson (local constant) estimator $\hat{m}(x)$.

We henceforth restrict attention to the uniform kernel $K(x) = 2^{-1}\mathbb{1}_{\{|x| \leq 1\}}$. For $L > 0$ let

$$\Theta_L = \{m : |m(y) - m(x)| \leq L|x - y| \text{ for all } x, y \in [0, 1]\}$$

denote the set of L -Lipschitz functions on $[0, 1]$. Writing \mathbb{E}_m for the expectation when the true mean function is m , show that

$$\sup_{m \in \Theta_L} \mathbb{E}_m \left[\int_h^{1-h} \{\hat{m}_h(x) - m(x)\}^2 dx \right] \leq L^2 h^2 + \frac{3\sigma^2}{2nh}$$

whenever $nh \geq 1$.

Deduce that there exists n_0 depending only on L and σ^2 such that, for $n \geq n_0$,

$$\inf_{h \in (0, 1/3)} \sup_{m \in \Theta_L} \mathbb{E}_m \left[\int_h^{1-h} \{\hat{m}_h(x) - m(x)\}^2 dx \right] \leq C \left(\frac{L\sigma^2}{n} \right)^{2/3}$$

where $C > 0$ is a constant that you should specify.

Fix $x_0 \in (0, 1)$. State Le Cam's two point lemma and use it to show that there exists a constant $c > 0$ such that

$$\inf_{\tilde{m}} \sup_{m \in \Theta_L} \mathbb{E}_m [\{\tilde{m}(x_0) - m(x_0)\}^2] \geq c \left(\frac{L\sigma^2}{n} \right)^{2/3},$$

where the infimum is taken over all estimators \tilde{m} based on $\{Y_1, \dots, Y_n\}$.

[Hint: You may wish to consider the function $m_1(x) = \lambda h K\left(\frac{x-x_0}{h}\right)$, where $K(t) = \exp\left(-\frac{1}{1-t^2}\right)\mathbb{1}_{\{|t| \leq 1\}}$ and λ is a scalar to be chosen.]

3

For a non-degenerate distribution function G define the *domain of attraction* of G , denoted $D(G)$, in the context of extreme value theory for sample maxima.

Defining the notion of a *regularly varying* function, state necessary and sufficient conditions for a distribution function F to satisfy $F \in D(G)$ for the three cases of G being the Fréchet(α), the Negative Weibull(α) and the Gumbel distribution functions. State sufficient conditions in terms of the hazard function.

For a positive integer m let $f_m(x) = \frac{x^{m-1}}{(m-1)!}e^{-x}$ denote the $\Gamma(m, 1)$ density. Writing F_m for the corresponding distribution function show that

$$F_m(x + \beta_n)^n \rightarrow e^{-e^{-x}}$$

as $n \rightarrow \infty$, for all $x \in \mathbb{R}$, where $\beta_n = \log n + (m-1) \log \log n - \log((m-1)!)$. [Hint: You may first wish to show that

$$1 - F_m(x) = e^{-x} \sum_{j=0}^{m-1} \frac{x^j}{j!}$$

for all $x \in [0, \infty)$.]

4

Let Y_1, \dots, Y_n be independent, mean-zero random variables with Y_i taking values in $[a_i, b_i]$. Prove Hoeffding's inequality,

$$\mathbb{P}\left(\left|\sum_{i=1}^n Y_i\right| > \epsilon\right) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right),$$

for each $\epsilon > 0$.

Now write $M := \max_{i=1, \dots, n} \max\{-a_i, b_i\}$ and $\sigma^2 := \max_{i=1, \dots, n} \mathbb{E}Y_i^2$. Assuming that you may interchange the order of expectation and summation where necessary, prove Bennett's inequality,

$$\mathbb{P}\left(\left|\sum_{i=1}^n Y_i\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\sigma^2}{M^2} \phi\left(\frac{M\epsilon}{n\sigma^2}\right)\right),$$

for each $\epsilon > 0$, where $\phi(x) = (1+x) \log(1+x) - x$. [Hint: You may wish to bound a moment generating function using the fact that $\mathbb{E}Y_i^k \leq \sigma^2 M^{k-2}$ for all $k \geq 2$.]

Let $X \sim \text{Bin}(n, p_n)$ with $p_n \rightarrow 0$ and $np_n \rightarrow \infty$ as $n \rightarrow \infty$. For a fixed $C > 0$ and large n , which of the two inequalities above provides a better bound on $\mathbb{P}\left(\frac{|X - np_n|}{\sqrt{np_n(1-p_n)}} \geq C\right)$?

END OF PAPER