Imperial College London

IMPERIAL COLLEGE LONDON BSc and MSci EXAMINATIONS (MATHEMATICS) May-June 2012

This paper is also taken for the relevant examination for the Associateship.

M2S2

Statistical Modelling

Date: Friday, 11 May 2012 Time: 2 pm - 4 pm

Credit will be given for all questions attempted but extra credit will be given for complete or nearly complete answers.

Calculators may not be used.

Statistical tables will not be available.

- 1. (i) Suppose $\theta \in \mathbb{R}$ is the unknown parameter in a statistical model and suppose that T is an estimator for θ .
 - (a) Define the bias and the mean squared error of T.
 - (b) How are bias, variance and mean squared error related? Prove this relationship.
 - (ii) State the definition of the non-central F-distribution.
 - (iii) Can there be a two-variate random vector \boldsymbol{Y} with $cov(\boldsymbol{Y}) = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$? Justify your answer.
 - (iv) State in one or two sentences the main difference between the Bayesian approach to statistics and the classical approach to statistics.
 - (v) In a Bayesian framework, suppose you observe Y_1,\ldots,Y_n . Conditionally on τ they are independent and follow a $N(\mu,\frac{1}{\tau})$ distribution. Suppose that $\mu\in\mathbb{R}$ is a known constant and that the a-priori distribution of τ is Gamma(2,2). Identify the posterior distribution of τ given Y_1,\ldots,Y_n .

[Recall that for $\alpha > 0$ and $\beta > 0$ the probability density function (pdf) of a $Gamma(\alpha, \beta)$ -distributed random variable is $f(x) = \frac{\beta^{\alpha}}{\Gamma(\alpha)}x^{\alpha-1}e^{-\beta x}, \quad x > 0.$]

2. Let X_1, \ldots, X_n be independent and identically distributed random variables following a Pareto distribution with parameter $\alpha > 0$, i.e. for $i = 1, \ldots, n$,

$$P(X_i \le x) = 1 - x^{-\alpha}, \quad x > 1,$$

and $P(X_i \le x) = 0$ for $x \le 1$.

- (i) Write down the likelihood function.
- (ii) Find the maximum likelihood estimator (MLE) $\widehat{\alpha}_n$ of α .
- (iii) Show that the Fisher information of X_1, \ldots, X_n for α is n/α^2 .
- (iv) State a positive lower bound on the variance of any unbiased estimator for α .
- (v) State large sample properties of the MLE $\widehat{\alpha}_n$ as $n \to \infty$. You do not need to verify regularity conditions.
- (vi) Construct a two-sided asymptotic 95% confidence interval for α .

- 3. (i) Consider a linear model of the form $Y = X\beta + \epsilon$, $E(\epsilon) = 0$ for some $\beta \in \mathbb{R}^p$, where $X \in \mathbb{R}^{n \times p}$ is a deterministic matrix. Furthermore, n > p.
 - (a) State if ϵ is observable.
 - (b) State the second order assumptions (SOA).
 - (c) State the normal theory assumptions (NTA).
 - (d) Give the definition of a least squares estimator of β .
 - (e) Show that any $\widehat{\beta}$ satisfying

$$X^T X \widehat{\boldsymbol{\beta}} = X^T \boldsymbol{Y}$$

is a least squares estimator of β .

- (f) Under what conditions is the least squares estimator unique?
- (g) Define the vector of residuals and the leverage of the *i*th observation. Under second order assumptions, derive the variance of the *i*th residual in terms of the leverage.
- (ii) Suppose that the vector of observation $oldsymbol{Y}$ satisfies

$$E(\boldsymbol{Y}) = X\boldsymbol{\beta} + \boldsymbol{Z},$$

for some $\beta \in \mathbb{R}^p$, where $X \in \mathbb{R}^{n \times p}$ is a deterministic matrix of full rank and $\mathbf{Z} \in \mathbb{R}^n$ is a deterministic vector. Furthermore, n > p.

Suppose that we are interested in estimating $c^T \beta$ for some known $c \in \mathbb{R}^p$.

Assume that you estimate $oldsymbol{eta}$ from the (incorrect) linear model

$$E(\mathbf{Y}) \in \operatorname{span}(X)$$

with second order assumptions using the least squares estimator $\widehat{\beta}$.

Compute the bias and the mean squared error of $c^T \hat{\beta}$ as an estimator for $c^T \beta$.

- 4. (i) State the Fisher-Cochran Theorem.
 - (ii) Describe the usual hypotheses considered in the F-test.
 - (iii) State the test statistic of the F-test in the form that uses residual sums of squares. State its distribution under the null hypothesis.
 - (iv) Describe the decision rule for an F-test to the level 5%.
 - (v) Rewrite the test statistic of the *F*-test using projection matrices.

Starting from the Fisher-Cochran Theorem, prove the distribution of the test statistic of the F-test under the null hypothesis.

You do not need to show that certain matrices are projection matrices and work out their rank; simply state which matrices are projection matrices and give their respective rank without proof.

(vi) What is the distribution of the test statistic of the F-test if the null hypothesis is not necessarily true? Prove this.

Imperial College London

IMPERIAL COLLEGE LONDON BSc and MSci EXAMINATIONS (MATHEMATICS) May-June 2012

This paper is also taken for the relevant examination for the Associateship.

M2S2

Statistical Modelling (Solutions)

Setter's signature	Checker's signature	Editor's signature

 $1. \quad \text{(i)} \quad \text{(a)} \ \text{bias}_{\theta}(T) = \mathrm{E}_{\theta}(T) - \theta.$

2

 $\mathsf{MSE}_{\theta}(T) = \mathrm{E}_{\theta}((T-\theta)^2).$ (b) $\mathsf{MSE}_{\theta}(T) = \mathrm{Var}_{\theta}(T) + (\mathsf{bias}_{\theta}(T))^2$

2

seen ↓

Let $X = T - \theta$. Then $\operatorname{Var}_{\theta}(X) = \operatorname{E}_{\theta}(X^2) - (\operatorname{E}_{\theta}X)^2 = \operatorname{MSE}_{\theta}(T) - (\operatorname{bias}_{\theta}(T))^2$. Rearraning shows the relationship.

3

(ii) If $W_1 \sim \chi^2_{n_1}(\delta)$, $W_2 \sim \chi^2_{n_2}$ independently then

$$F = \frac{W_1/n_1}{W_2/n_2}$$

is said to have a *non-central F* distribution with (n_1,n_2) d.f. and non-centrality parameter δ .

3

(iii)

No, as the matrix is not positive semidefinite. To see the latter,

$$(1,-1)\begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ -1 \end{pmatrix} = -2 < 0.$$

sim. seen \downarrow

(iv) In the classical approach, the parameter is an unknown constant, whereas in the Bayesian approach, it is the realisation of a random variable.

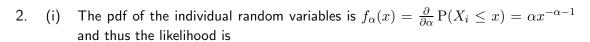
2

3

(v) $\pi(\tau) \propto \tau e^{-2\tau} \text{ and } f(Y_i|\tau) \propto \sqrt{\tau} e^{-\tau \frac{(Y_i-\mu)^2}{2}}. \text{ Hence,}$

$$\pi(\tau) \prod_{i=1}^{n} f(Y_i | \tau) \propto \tau e^{-2\tau} \prod_{i=1}^{n} \sqrt{\tau} e^{-\tau \frac{(Y_i - \mu)^2}{2}} = \tau^{(2+n/2)-1} e^{-(2+\sum_{i=1}^{n} (Y_i - \mu)^2/2)\tau}$$

Hence,
$$\tau|Y_1,\ldots,Y_n \sim \mathsf{Gamma}(2+\frac{n}{2},2+\frac{1}{2}\sum_{i=1}^n(Y_i-\mu)^2)$$



sim. seen ↓

$$L(\alpha) = \prod_{i=1}^{n} f_{\alpha}(x_i) = \alpha^n \prod_{i=1}^{n} x_i^{-\alpha - 1}.$$

3

(ii) The log-likelihood $l(\alpha)$ is

$$l(\alpha) = n \log(\alpha) - (\alpha + 1) \sum_{i=1}^{n} \log(x_i).$$

Taking derivatives gives

$$\frac{\partial}{\partial \alpha} l(\alpha) = \frac{n}{\alpha} - \sum_{i=1}^{n} \log(x_i).$$

Equating this to 0 gives $\widehat{\alpha} = \frac{n}{\sum_{i=1}^{n} \log(x_i)}$.

This indeed a maximum likelihood estimator as

 $\left(\frac{\partial}{\partial \alpha}\right)^2 l(\alpha) = -\frac{n}{\alpha^2} < 0$

2

3

(iii) First, we compute the Fisher information $I_1(\alpha)$ for a sample of size 1.

$$\log f_{\alpha}(x) = \log(\alpha) - (\alpha + 1)\log(x),$$

$$\left(\frac{\partial}{\partial \alpha}\right)^2 \log f_{\alpha}(x) = \frac{\partial}{\partial \alpha} \left(\frac{1}{\alpha} - \log(x)\right) = -\frac{1}{\alpha^2},$$

and thus

$$I_1(\alpha) = -E\left[\left(\frac{\partial}{\partial \alpha}\right)^2 \log f_{\alpha}(X_1)\right] = \frac{1}{\alpha^2}$$

Hence, due to independence, the Fisher information $I_n(\alpha)$ for α based on X_1, \ldots, X_n is $\frac{n}{\alpha^2}$.

3

[Alternatively, one can compute $I_1(\alpha)$ via

$$I_1(\alpha) = E\left[\left(\frac{\partial}{\partial \alpha}\log f_{\alpha}(X_1)\right)^2\right] = E\left[\left(\frac{1}{\alpha} - \log(X_1)\right)^2\right].$$

However, working out the expectation is somewhat lengthy.

(iv) By the Rao-Cramer inequality, any unbiased estimator T for α satisfies

$$\operatorname{Var}_{\alpha}(T) \ge \frac{1}{I_n(\alpha)} = \frac{\alpha^2}{n}.$$

3

(v)
$$\sqrt{n}(\widehat{\alpha} - \alpha) \stackrel{d}{\to} N(0, \underbrace{\frac{1}{I_1(\alpha)}}_{=\alpha^2}) \text{ as } n \to \infty.$$

(vi) From Part (v), we get

$$\sqrt{n}(\widehat{\alpha}/\alpha - 1) \stackrel{d}{\to} N(0, 1) \quad (n \to \infty)$$

Let c_1 and c_2 be such that $\Phi(c_1)=0.025$ and $\Phi(c_2)=0.975$, where Φ is the cdf of a standard normal distribution. Then

$$P(c_1 \le \sqrt{n}(\frac{\widehat{\alpha}}{\alpha} - 1) \le c_2) \to 0.95.$$

Solving the inequalities for α gives

$$\frac{\widehat{\alpha}}{n^{-\frac{1}{2}}c_2 + 1} \le \alpha \le \frac{\widehat{\alpha}}{n^{-\frac{1}{2}}c_1 + 1}$$

implying the approximate confidence interval

$$\left[\frac{\widehat{\alpha}}{n^{-\frac{1}{2}}c_2+1}, \frac{\widehat{\alpha}}{n^{-\frac{1}{2}}c_1+1}\right]$$

Alternative solution: First, show that $\widehat{\alpha} \to \alpha$: The survival function of $\log(X_i)$ is



$$P(\log(X_i) > t) = P(X_i > \exp(t)) = \exp(t)^{-\alpha} = \exp(-\alpha t)$$

Thus, by the law of large numbers $\frac{1}{n} \sum_{i=1}^{n} \log(X_i) \to E \log(X_1) = \int_0^{\infty} e^{-\alpha x} dx = 1/\alpha$.

Then, using, by Slutsky's Lemma,

$$\sqrt{n}\left(1 - \frac{\alpha}{\widehat{\alpha}}\right) = \frac{\sqrt{n}(\widehat{\alpha} - \alpha)}{\widehat{\alpha}} = \sqrt{n}(\widehat{\alpha} - \alpha)\frac{1}{n}\sum_{i=1}^{n}\log(X_i) \xrightarrow{d} N(0, 1) \quad (n \to \infty)$$

Let c_1 and c_2 be such that $\Phi(c_1) = 0.025$ and $\Phi(c_2) = 0.975$, where Φ is the cdf of a standard normal distribution. Then

$$P(c_1 \le \sqrt{n}(1 - \frac{\alpha}{\widehat{\alpha}}) \le c_2) \to 0.95.$$

Solving the inequalities for α gives

$$(1 - n^{-\frac{1}{2}}c_1)\widehat{\alpha} \ge \alpha \ge (1 - n^{-\frac{1}{2}}c_2)\widehat{\alpha}$$

giving the approximate confidence interval $[(1-n^{-\frac{1}{2}}c_2)\widehat{\alpha}, (1-n^{-\frac{1}{2}}c_1)\widehat{\alpha}].$



seen ↓

No, ϵ is not observable.



(b) Second Order Assumption (SOA): $cov(\epsilon) = \sigma^2 I_n$ for some $\sigma^2 > 0$.



(c) Normal theory assumptions (NTA): $\epsilon \sim N(\mathbf{0}, \sigma^2 I_n)$ for some $\sigma^2 > 0$.



(d) Let $S(\beta) = (Y - X\beta)^T (Y - X\beta)$. Any p-variate random vector $\widehat{\beta}$ such that

$$S(\widehat{\boldsymbol{\beta}}) = \min_{\boldsymbol{\beta} \in \mathbb{R}^p} S(\boldsymbol{\beta})$$

is called a least squares estimator.

2

(e) Let $\widehat{\boldsymbol{\beta}}$ satisfy $X^T X \widehat{\boldsymbol{\beta}} = X^T \boldsymbol{Y}$. Then for any $\boldsymbol{\beta} \in \mathbb{R}^p$,

$$S(\boldsymbol{\beta}) = (\boldsymbol{Y} - X\widehat{\boldsymbol{\beta}} + X\widehat{\boldsymbol{\beta}} - X\boldsymbol{\beta})^T (\boldsymbol{Y} - X\widehat{\boldsymbol{\beta}} + X\widehat{\boldsymbol{\beta}} - X\boldsymbol{\beta})$$

$$= S(\widehat{\boldsymbol{\beta}}) + \underbrace{(X\widehat{\boldsymbol{\beta}} - X\boldsymbol{\beta})^T (X\widehat{\boldsymbol{\beta}} - X\boldsymbol{\beta})}_{\geq 0} + 2\underbrace{(X\widehat{\boldsymbol{\beta}} - X\boldsymbol{\beta})^T (\boldsymbol{Y} - X\widehat{\boldsymbol{\beta}})}_{=(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T (X^T \boldsymbol{Y} - X^T X\widehat{\boldsymbol{\beta}}) = 0}$$

$$\geq S(\widehat{\boldsymbol{\beta}})$$

Hence, $\widehat{\boldsymbol{\beta}}$ minimises S and is therefore a least squares estimator.

3

(f) The least squares estimator is unique iff X is of full rank.

2

(g) The vector e of residuals is defined by $e = Y - X\widehat{\beta}$. Let P be the projection matrix onto the space spanned by the columns of X. The leverage of the ith observation is P_{ii} .

As
$$\boldsymbol{e} = \boldsymbol{Y} - P\boldsymbol{Y} = (I - P)\boldsymbol{Y}$$
, we have

$$cov(\mathbf{e}) = cov((I - P)\mathbf{Y}) = \sigma^2(I - P)(I - P)^T$$

Using that (I-P) is a projection matrix this simplifies to

$$cov(e) = \sigma^2(I - P)$$

Hence, $\operatorname{Var} e_i = \sigma^2 (1 - P_{ii})$.

(ii)

$$E c^{T} \hat{\boldsymbol{\beta}} = E[c^{T} (X^{T} X)^{-1} X^{T} Y] = c^{T} (X^{T} X)^{-1} X^{T} E[Y]$$
$$= c^{T} (X^{T} X)^{-1} X^{T} (X \boldsymbol{\beta} + \boldsymbol{Z}) = c^{T} \boldsymbol{\beta} + c^{T} (X^{T} X)^{-1} X^{T} \boldsymbol{Z}$$

Thus $\operatorname{bias}(\boldsymbol{c}^T\widehat{\boldsymbol{\beta}}) = \boldsymbol{c}^T(X^TX)^{-1}X^T\boldsymbol{Z}.$

Furthermore,

$$MSE(\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}) = \text{Var}(\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}) + (\text{bias}(\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}))^{2} = \boldsymbol{c}^{T} \operatorname{cov}(\widehat{\boldsymbol{\beta}})\boldsymbol{c} + (\text{bias}(\boldsymbol{c}^{T}\widehat{\boldsymbol{\beta}}))^{2}$$
$$= \boldsymbol{c}^{T}(X^{T}X)^{-1}\boldsymbol{c}\sigma^{2} + (\boldsymbol{c}^{T}(X^{T}X)^{-1}X^{T}\boldsymbol{Z})^{2}$$

4. (i) If A_1, \ldots, A_k are $n \times n$ projection matrices such that $\sum_{i=1}^k A_i = I_n$, and if $\mathbf{Z} \sim N(\boldsymbol{\mu}, I_n)$ then $\mathbf{Z}^T A_1 \mathbf{Z}, \ldots, \mathbf{Z}^T A_k \mathbf{Z}$ are independent and

 $oldsymbol{Z}^T A_i oldsymbol{Z} \sim \chi^2_{r_i}(\delta_i), \quad ext{where } r_i = ext{rank}\, A_i ext{ and } \delta_i^2 = oldsymbol{\mu}^T A_i oldsymbol{\mu}.$

3

(ii) We want to test whether a sub-model of a linear model $\mathbf{E}\, m{Y} = X m{eta}$ is true, i.e. we want to test

 $H_0 : \mathbf{E} \mathbf{Y} \in \mathrm{span}(X_0) \text{ against } H_1 : \mathbf{E} \mathbf{Y} \notin \mathrm{span}(X_0)$

for some matrix X_0 with $\mathrm{span}(X_0) \subset \mathrm{span}(X)$.

(iii) Under $H_0 : \mathbf{E} \mathbf{Y} \in \text{span}(X_0)$,

 $F = \frac{\text{RSS}_0 - \text{RSS}}{\text{RSS}} \cdot \frac{n-r}{r-s} \sim F_{r-s,n-r}$

where $r = \operatorname{rank} X$, $s = \operatorname{rank} X_0$, RSS= residual sum of squares in the full model $\operatorname{E} \boldsymbol{Y} = X\boldsymbol{\beta}$ and RSS= residual sum of squares in the sub-model $\operatorname{E} \boldsymbol{Y} = X_0 \widetilde{\boldsymbol{\beta}}$.

3

2

(iv) The null is rejected for large values of ${\cal F}.$ As we want a test to the level 0.05, we reject if

F > c,

where c is such that $P(X \ge c) = 0.05$ for $X \sim F_{r-s,n-r}$.

(v) Let P be the projection matrix onto $\operatorname{span} X$, and let Q = I - P the projection matrix onto $(\operatorname{span} X)^{\perp}$.

Likewise, let P_0 be the projection matrix onto $\operatorname{span} X_0$, and let $Q_0 = I - P_0$ the projection matrix onto $(\operatorname{span} X_0)^{\perp}$. Then,

$$RSS = \mathbf{Y}^T Q \mathbf{Y}, \quad RSS_0 = \mathbf{Y}^T Q_0 \mathbf{Y}.$$

Using this gives

$$F = \frac{\boldsymbol{Y}^T Q_0 \boldsymbol{Y} - \boldsymbol{Y}^T Q \boldsymbol{Y}}{\boldsymbol{Y}^T Q \boldsymbol{Y}} \cdot \frac{n-r}{r-s} = \frac{\boldsymbol{Y}^T (P-P_0) \boldsymbol{Y} / \sigma^2}{\boldsymbol{Y}^T (I-P) \boldsymbol{Y} / \sigma^2} \cdot \frac{n-r}{r-s}$$

Let $Z = Y/\sigma$.

Let $A_1 = I - P$, $A_2 = P - P_0$, $A_3 = P_0$.

Clearly, $A_1 + A_2 + A_3 = I$. Furthermore, A_1 , A_2 , and A_3 are projection matrices.

The Fisher-Cochran theorem now implies

*
$$\boldsymbol{Z}^T(P-P_0)\boldsymbol{Z}$$
 and $\boldsymbol{Z}^T(I-P)\boldsymbol{Z}$ are independent,

*
$$\mathbf{Z}^T(P-P_0)\mathbf{Z} \sim \chi^2_{\text{rank}(P-P_0)}(\mathbf{E}\,\mathbf{Z}^T(P-P_0)\,\mathbf{E}\,\mathbf{Z}),$$

*
$$\mathbf{Z}^T(I-P)\mathbf{Z} \sim \chi^2_{\text{rank}(I-P)}(\mathbf{E}\,\mathbf{Z}^T(I-P)\,\mathbf{E}\,\mathbf{Z}).$$

We show that the non-centrality parameters are 0.

Under H_0 , we know $\mathbf{E} \mathbf{Z} = \frac{1}{\delta} \mathbf{E} \mathbf{Y} \in \mathrm{span}(X_0) \subset \mathrm{span}(X)$ Thus,

$$(P - P_0) \to \mathbf{Z} = \underbrace{P \to \mathbf{Z}}_{= \to \mathbf{Z}} - \underbrace{P_0 \to \mathbf{Z}}_{= \to \mathbf{Z}} = \mathbf{0}.$$

Hence, $\mathrm{E}\, {m Z}^T(P-P_0)\, \mathrm{E}\, {m Z}=0.$ Furthermore,

$$E\mathbf{Z}^{T}(I-P) \mathbf{E}\mathbf{Z} = E\mathbf{Z}^{T}(\mathbf{E}\mathbf{Z} - \underbrace{P \mathbf{E}\mathbf{Z}}_{=\mathbf{E}\mathbf{Z}}) = 0.$$

Without proof: rank(I - P) = n - r, $rank(P - P_0) = r - s$.

To summarise, we have shown

$$m{Z}^T(P-P_0)m{Z}\sim\chi^2_{r-s},\,m{Z}^T(I-P)m{Z}\sim\chi^2_{n-r}$$
 independently.

Thus, by definition, $F \sim F_{r-s,n-r}$.

(vi) If H_0 is not true then the proof is still valid, except for the non-centrality parameter of ${\bf Z}^T(P-P_0){\bf Z}$. Now,

$$\mathbf{E} \mathbf{Z}^T (P - P_0) \mathbf{E} \mathbf{Z} = \frac{1}{\sigma^2} \boldsymbol{\beta}^T X^T (P - P_0) X \boldsymbol{\beta}.$$

Thus, without assuming H_0 , we get

$$F \sim F_{r-s,n-r}(\delta)$$
, where $\delta^2 = \frac{1}{\sigma^2} (X\boldsymbol{\beta})^T (P - P_0) X\boldsymbol{\beta}$.

5